

JOURNAL OF MATHEMATICAL PHYSICS

VOLUME 9, NUMBER 3

MARCH 1968

Minimum Uncertainty Product, Number-Phase Uncertainty Product, and Coherent States

ROMAN JACKIW*†

Jefferson Physical Laboratory, Harvard University, Cambridge, Massachusetts

(Received 10 April 1967)

The number-phase uncertainty products proposed by Carruthers and Nieto are studied to determine whether they are minimized by coherent states. It is found that coherent states do not minimize these products. States that do minimize some of the uncertainty products are constructed. Variational techniques for the study of arbitrary uncertainty products are developed.

I. INTRODUCTION

Recent discussions of a quantum-mechanical phase operator for harmonic oscillators have shown that a Hermitian phase operator φ_{op} does not exist.¹ Susskind and Glogower¹ (SG) have demonstrated however that Hermitian sine (S) and cosine (C) operators can be defined which have many properties that are suggested by the nomenclature. Carruthers and Nieto² (CN) have examined the matrix elements of S and C between Glauber's³ coherent states. They found that in the high-excitation (classical) limit the expectation values of S and C, in these states, behave as the sine and cosine of the phase of the harmonic oscillation.

Carruthers and Nieto have also proposed several uncertainty relations to replace the traditional expression for the number-phase uncertainty product

$$(\Delta N)^2(\Delta\varphi)^2 \geq \frac{1}{4}, \quad (1)$$

which is ill-defined. The proposed uncertainty products are given in terms of the S and C operators and have the virtue that, when evaluated with coherent

states, they approach their theoretic minimum for highly excited coherent states, and remain small for moderate excitation.

In this paper, we examine further the uncertainty products given by CN, in order to determine whether the coherent states do in fact give the smallest uncertainty product. Towards this end, we develop in Sec. II new variational techniques for determining those normalizable states which give a minimum for the uncertainty product of noncommuting Hermitian operators. In Sec. III we show that the coherent states do not minimize the number-phase uncertainty products. We also construct states which do have the desired property. In Sec. IV we examine the S-C uncertainty product.

II. MINIMUM UNCERTAINTY PRODUCT

A. When two Hermitian operators X and Y do not commute, they cannot be simultaneously diagonalized and their uncertainty product satisfies the inequality

$$(\Delta X)^2(\Delta Y)^2 \geq \frac{1}{4}\langle A \rangle^2. \quad (2)$$

Here $(\Delta X)^2 \equiv \langle X^2 \rangle - \langle X \rangle^2$, and iA is the commutator of X and Y , assumed to be nonzero. A procedure for determining the state for which the uncertainty product, appearing in (2), is minimized was given by

* Junior Fellow, Society of Fellows.

† Address during 1967-68: CERN, Geneva, Switzerland.

¹ L. Susskind and J. Glogower, *Physics* **1**, 49 (1964).

² P. Carruthers and M. M. Nieto, *Phys. Rev. Letters* **14**, 387 (1965).

³ R. J. Glauber, *Phys. Rev.* **131**, 2766 (1963).

Heisenberg when X and Y are position and momentum operators. His method holds generally when A is a c number, and is described in any textbook on quantum mechanics. However, if A is not a c number, Heisenberg's method cannot, in general, be used to obtain the minimizing state. We briefly summarize here Heisenberg's method, show that it is inapplicable when A is a q number, and then develop a variational method appropriate to the case when A is a q number. This latter method is also applicable when the uncertainty product no longer has the simple form (2).

B. Heisenberg's method consists of establishing that for any normalizable state $|\Psi\rangle$,

$$U(\Psi) \equiv (\Delta X)^2(\Delta Y)^2 = \langle \hat{X}^2 \rangle \langle \hat{Y}^2 \rangle = |\langle \hat{X} \hat{Y} \rangle|^2 + R(\Psi) \\ = \frac{1}{4}P(\Psi) + \frac{1}{4}Q(\Psi) + R(\Psi), \quad (3a)$$

where

$$\hat{X} \equiv X - \langle X \rangle, \\ P(\Psi) \equiv |\langle [\hat{X}, \hat{Y}] \rangle|^2 = |\langle [X, Y] \rangle|^2 = \langle A \rangle^2, \\ Q(\Psi) \equiv |\langle \{\hat{X}, \hat{Y}\} \rangle|^2. \quad (3b)$$

The term $R(\Psi)$ in (3a) is a positive semidefinite remainder term, arising from an application of the Schwartz inequality to $\langle \hat{X}^2 \rangle \langle \hat{Y}^2 \rangle$. $R(\Psi)$ vanishes if and only if $\hat{X}|\Psi\rangle$ is proportional to $\hat{Y}|\Psi\rangle$. $P(\Psi)$ and $Q(\Psi)$ are also positive semidefinite, depending in general on Ψ . However, if A is a c number, $P(\Psi)$ does not depend on Ψ since $\langle \Psi | \Psi \rangle = 1$. In this case

$$U(\Psi) = \frac{1}{4}A^2 + \frac{1}{4}Q(\Psi) + R(\Psi). \quad (4)$$

Clearly the minimum value for U is $\frac{1}{4}A^2$, which is reached if and only if Q and R vanish. Combining the requirement that R vanishes with the requirement that Q vanishes gives an equation for $|\Psi\rangle$:

$$\hat{X}|\Psi\rangle + i\gamma\hat{Y}|\Psi\rangle = 0, \quad \gamma \text{ real}, \quad (5a)$$

or

$$[X + i\gamma Y]|\Psi\rangle = \lambda|\Psi\rangle, \\ \lambda \equiv \langle X \rangle + i\gamma\langle Y \rangle. \quad (5b)$$

For future reference, we give another equation for $|\Psi\rangle$ which follows from (5a) upon multiplication by $\hat{X} - i\gamma\hat{Y}$:

$$[\hat{X}^2 + \gamma^2\hat{Y}^2 - \gamma A]|\Psi\rangle = 0. \quad (6a)$$

Equation (6a) shows that γ may be evaluated in terms of matrix elements of X and Y . By premultiplying (6a) by $\langle \Psi |$, we get

$$\gamma = \frac{A}{2(\Delta Y)^2} = \pm \left(\frac{(\Delta X)^2}{(\Delta Y)^2} \right)^{\frac{1}{2}}. \quad (6b)$$

In obtaining (6b), we have used the fact that A is a c number and $(\Delta X)^2(\Delta Y)^2 = \frac{1}{4}A^2$. (We assume always

that $(\Delta X)^2(\Delta Y)^2 \neq 0$; viz., $|\Psi\rangle$ is not an eigenstate of X or Y .⁴)

Equation (5b) is to be solved as an eigenvalue equation for $|\Psi\rangle$, with three free parameters: γ , $\text{Re } \lambda$, $\text{Im } \lambda$, and a subsidiary normalization condition $\langle \Psi | \Psi \rangle = 1$. The state $|\Psi\rangle$ then minimizes the uncertainty product (3c). We refer to this method as the *direct method* for obtaining the minimizing state.

C. The direct method is not in general applicable when A is a q number. In that eventuality, $P(\Psi)$ does depend on Ψ , and we cannot conclude that a minimum for U is achieved when Q and R vanish.

We now solve the minimization problem without assuming that the commutator of X and Y is a c number. For an uncertainty product of the form $(\Delta X)^2(\Delta Y)^2$, the problem is nontrivial only when the matrix elements of X between eigenstates of Y diverge, and vice versa. For if they are finite, the uncertainty product is manifestly minimized to zero when it is evaluated with eigenstates of X or Y . The subsequent analysis applies only to the nontrivial problem. However, later we generalize it to the case when the uncertainty products are not of the simple form $(\Delta X)^2(\Delta Y)^2$ and it no longer is obvious how to find the minimizing state, even though all matrix elements are finite.

Since the expression for $U(\Psi)$ given in (3c) involves the function $R(\Psi)$ about which we have no useful information, we return to the definition of $U(\Psi)$ in terms of matrix elements of X and Y . If $U(\Psi)$ is to be minimized, we may apply the variation principle and require that $U(\Psi)$ be stationary under arbitrary variations of Ψ . With the help of a Lagrange multiplier m , we also impose the subsidiary condition that $\langle \Psi | \Psi \rangle = 1$. Considering the variation of $|\Psi\rangle$ to be independent of $\langle \Psi |$, we obtain as a necessary condition for $U(\Psi)$ to be a minimum:

$$\delta U / \delta \langle \Psi | = m |\Psi\rangle. \quad (7a)$$

Since $U \equiv (\Delta X)^2(\Delta Y)^2$, we need to evaluate $\delta(\Delta X)^2 / \delta \langle \Psi |$. From the definition of $(\Delta X)^2$, we have

$$(\Delta X)^2 = \langle \Psi | X^2 | \Psi \rangle - \langle \Psi | X | \Psi \rangle^2, \\ \delta(\Delta X)^2 / \delta \langle \Psi | = X^2 | \Psi \rangle - 2X | \Psi \rangle \langle X \rangle \\ = [X - \langle X \rangle]^2 | \Psi \rangle - \langle X \rangle^2 | \Psi \rangle. \quad (7b)$$

Therefore (7a) becomes

$$(\Delta Y)^2 \hat{X}^2 | \Psi \rangle + (\Delta X)^2 \hat{Y}^2 | \Psi \rangle \\ = ((\Delta Y)^2 \langle X \rangle^2 + (\Delta X)^2 \langle Y \rangle^2 + m) | \Psi \rangle. \quad (7c)$$

⁴ If $|\Psi\rangle$ is an eigenstate of one of the two operators, say X , then $(\Delta X)^2 = 0$, and $(\Delta Y)^2$ necessarily diverges when A is a c number. Therefore the uncertainty product is of the indeterminate form $0 \cdot \infty$. We find in our subsequent analysis that such an indeterminate quantity can be sometimes evaluated; see Appendix A.

Finally, by taking matrix elements of the above with $\langle \Psi |$, we evaluate m , and discover that the coefficient on the right-hand side is $2(\Delta X)^2(\Delta Y)^2$, which we assume to be nonzero. Thus we obtain an Euler-Lagrange (EL)-type equation for $|\Psi\rangle$ which must be satisfied if U is to be a minimum:

$$\left[\frac{\hat{X}^2}{(\Delta X)^2} + \frac{\hat{Y}^2}{(\Delta Y)^2} - 2 \right] |\Psi\rangle = 0. \quad (8)$$

We call the states $|\Psi\rangle$ which solve (8) *critical states*.⁵ Evidently every critical state $|\Psi\rangle$ makes $U(\Psi)$ stationary.

This simple equation provides the appropriate generalization of the direct method to the case that $[X, Y]$ is a q number. For future reference we call this the *analytic method*. Equation (8) is to be solved as an eigenvalue equation with four free, real parameters, viz.,

$$\left[\frac{(X - \alpha)^2}{a^2} + \frac{(Y - \beta)^2}{b^2} - 2 \right] |\Psi\rangle = 0. \quad (9a)$$

Once $|\Psi\rangle$ has been obtained from (9a), the four parameters are determined self-consistently by setting⁶

$$\begin{aligned} \alpha &= \langle X \rangle, & \beta &= \langle Y \rangle, \\ a^2 &= \langle X^2 \rangle - \langle X \rangle^2, & b^2 &= \langle Y^2 \rangle - \langle Y \rangle^2. \end{aligned} \quad (9b)$$

In general, since Eq. (9a) is an eigenvalue equation with self-consistency conditions (9b), we expect to obtain solutions only when a special (eigenvalue) relation exists between the parameters. Nevertheless, we may expect to obtain more than one solution, since Eq. (8) serves equally well to determine other stationary points of U : further minima, maxima, or "points of inflection" of U . One must therefore examine $(\Delta X)^2(\Delta Y)^2 = a^2b^2$ for the various critical states, to determine which gives the smallest value. (If it is not evident that a minimum has indeed been attained, one might compute the second variation of U to determine the nature of the stationary point.)

Having established a necessary condition on $|\Psi\rangle$ for $U(\Psi)$ to be a minimum, we may examine the direct method critically to establish its precise relation to the

⁵ Eigenstates of X and Y pose a special problem. For suppose we take $|\Psi\rangle$ to be an eigenstate of X , and assume that $(\Delta Y)^2$ diverges so that the problem is nontrivial. Then Eq. (8) has the indeterminate form $0/0$ $|\Psi\rangle + \hat{Y}^2/\infty |\Psi\rangle - 2|\Psi\rangle = 0$. Evidently an effective point of view is to ignore those solutions of (8) which are eigenstates of X and Y , and evaluate U separately with the eigenstates to determine whether these minimize U .

⁶ In the direct method, the parameters λ and γ need not be evaluated separately since their value is set by the form of Eq. (5a). Indeed the four conditions in (9b) are redundant since the form of Eq. (9a) assures that one relation between the parameters exists, viz.,

$$\frac{\langle X^2 \rangle - 2\alpha\langle X \rangle + \alpha^2}{a^2} + \frac{\langle Y^2 \rangle - 2\beta\langle Y \rangle + \beta^2}{b^2} = 2.$$

analytic method. Suppose we set out to determine a state $|\Psi\rangle$ by the direct method (regardless of the nature of A). Then (6a) is valid, which may now be written as

$$\left[\frac{\hat{X}^2}{(\Delta X)^2} + \frac{\gamma^2 \hat{Y}^2}{(\Delta X)^2} - \frac{\gamma A}{(\Delta X)^2} \right] |\Psi\rangle = 0. \quad (10a)$$

The parameter γ may be again evaluated by taking matrix elements and remembering that (within the direct method) $(\Delta X)^2(\Delta Y)^2 = \frac{1}{4}A^2$. Then [compare (6b)],

$$\gamma = \frac{\langle A \rangle}{2(\Delta Y)^2} = \pm \left(\frac{(\Delta X)^2}{(\Delta Y)^2} \right)^{\frac{1}{2}} \quad (10b)$$

and (10a) becomes

$$\left[\frac{\hat{X}^2}{(\Delta X)^2} + \frac{\hat{Y}^2}{(\Delta Y)^2} - \frac{2A}{\langle A \rangle} \right] |\Psi\rangle = 0. \quad (10c)$$

Comparing this to (8), we see that the direct method determines a critical state $|\Psi\rangle$ which corresponds to a stationary value of $U(\Psi)$, if and only if $|\Psi\rangle$ is an eigenstate of A .

In conclusion we note that even when A is a c number and the direct method is applicable, Eq. (5) may not have a solution. Then U never achieves its minimum of $\frac{1}{4}A^2$. Nevertheless, it may achieve some lowest value which is greater than $\frac{1}{4}A^2$; and the analytic method may be used to determine the states for which this occurs.

D. In order to exhibit the workings of our analytic method, we solve the classic problem of minimizing the position-momentum uncertainty product

$$(\Delta x)^2(\Delta p)^2.$$

In obtaining this old result, we find all the critical states for which $(\Delta x)^2(\Delta p)^2$ is stationary.

According to (8) we must solve ($\hbar = 1$)

$$\left[\frac{(x - \alpha)^2}{a^2} + \frac{[(1/i)\partial/\partial x - \beta]^2}{b^2} \right] \Psi(x) = 2\Psi(x). \quad (11a)$$

The solution to (11a) can be found by comparison to the Schrödinger equation for a harmonic oscillator. Therefore (11a) possesses normalized solutions only when

$$|ab| = \frac{1}{2}(2n + 1). \quad (11b)$$

The normalized solutions are

$$\Psi_n(x) = e^{i\beta x} U_n(x - \alpha), \quad (11c)$$

where $U_n(x)$ is a normalized harmonic-oscillator eigenfunction, with mass $\frac{1}{2}b^2$, stiffness constant $2/a^2$ and energy 2. The self-consistency requirements (9b) set no further conditions beyond (11b), and all the Ψ_n 's are critical states for $(\Delta x)^2(\Delta p)^2$. Evidently the

minimum uncertainty product is $\frac{1}{4}$, which is attained with the state Ψ_0 :

$$\Psi_0(x) = [2\pi(\Delta x)^2]^{-\frac{1}{4}} \exp\{-[(x - \langle x \rangle)^2/(\Delta x)^2]\}. \quad (11d)$$

The fact that an oscillator ground-state wavefunction minimizes the position-momentum uncertainty product is well known, and has been considered to be a fortuitous coincidence. It is seen from the present analysis that this result is a natural consequence of our analytic method. Moreover, we have obtained the further knowledge that all the harmonic-oscillator wavefunctions are critical states which make $(\Delta x)^2(\Delta p)^2$ stationary.

E. We continue with our discussion of the minimization problem for uncertainty products by discussing objects which are of a form more complicated than $(\Delta X)^2(\Delta Y)^2$. (Such uncertainty products have been proposed by CN.)

If the commutator of X and Y is not a c number, it may be of consequence to consider an uncertainty product of the form

$$U_1(\Psi) \equiv \frac{(\Delta X)^2(\Delta Y)^2}{\langle A \rangle^2} = \frac{U(\Psi)}{\langle A \rangle^2}. \quad (12a)$$

By applying the variation principle, we immediately obtain the necessary condition on $|\Psi\rangle$, for which $U_1(\Psi)$ is stationary:

$$\left[\frac{\hat{X}^2}{(\Delta X)^2} + \frac{\hat{Y}^2}{(\Delta Y)^2} - \frac{2A}{\langle A \rangle} \right] |\Psi\rangle = 0. \quad (12b)$$

This equation is the same as Eq. (10b) which follows from the direct method. Indeed, that the direct method is applicable, may be seen by reference to Eq. 3c).

According to that expression,

$$U_1(\Psi) = \frac{1}{4} + \frac{1}{4} \frac{Q(\Psi)}{P(\Psi)} + \frac{R(\Psi)}{P(\Psi)}. \quad (12c)$$

Thus when we arrange for Q and R to vanish, as is done in the direct method, U_1 attains its minimum. [We must of course examine separately the situation if the direct method yields a solution for which $P(\Psi) = 0$.]

When the expression for the uncertainty product is even more complicated, for example if it involves the matrix elements of more than two operators, the direct method, even if applicable, will not in general yield solutions. The variation principle may nevertheless be used to give a (complicated) necessary condition.

III. NUMBER-PHASE UNCERTAINTY PRODUCTS

A. We now turn to the number-phase uncertainty products proposed by CN. Following SG,¹ we

consider a harmonic oscillator described by creation and annihilation operators a and a^\dagger , respectively, which obey $[a, a^\dagger] = 1$. The number operator $N_{op} \equiv a^\dagger a$ has number states $|n\rangle$ as eigenvectors; and $a|n\rangle = n^{\frac{1}{2}}|n-1\rangle$, $a^\dagger|n\rangle = (n+1)^{\frac{1}{2}}|n+1\rangle$, $a^\dagger a|n\rangle = n|n\rangle$. The eigenvectors of a are the coherent states $|\alpha\rangle$: $a|\alpha\rangle = \alpha|\alpha\rangle$. They have the property that $N \equiv \langle \alpha | N_{op} | \alpha \rangle = |\alpha|^2 = \langle \alpha | N_{op}^2 | \alpha \rangle - \langle \alpha | N_{op} | \alpha \rangle^2 \equiv (\Delta N)^2$. In terms of number states, the coherent states are given by

$$|\alpha\rangle = e^{-\frac{1}{2}|\alpha|^2} \sum_{n=0}^{\infty} \frac{\alpha^n}{(n!)^{\frac{1}{2}}} |n\rangle. \quad (14)$$

Evidently each coherent state may be described by two parameters: amplitude and phase of α . Thus we frequently write $|N\varphi\rangle$ for $|\alpha\rangle$ where $\alpha = N^{\frac{1}{2}}e^{i\varphi}$. To define the sine and cosine operators, we define first the number state raising and lowering operators E_\pm :

$$E_- \equiv (N_{op} + 1)^{-\frac{1}{2}}, \\ E_+ \equiv a^\dagger (N_{op} + 1)^{-\frac{1}{2}} = (E_-)^\dagger. \quad (15a)$$

These satisfy

$$E_\pm |n\rangle = |n \pm 1\rangle n \neq 0, \\ E_+ |0\rangle = |1\rangle, \quad E_- |0\rangle = 0, \\ E_- E_+ = I, \quad E_+ E_- = I - P, \\ [E_-, E_+] = P, \quad P |n\rangle = \delta_{n0} |0\rangle. \quad (15b)$$

The S and C operators then are defined by

$$C = \frac{1}{2}[E_- + E_+], \quad S = \frac{1}{2}i^{-1}[E_- - E_+], \quad (16a)$$

$$[N_{op}, S] = iC, \quad [N_{op}, C] = -iS, \quad [S, C] = \frac{1}{2i}P. \quad (16b)$$

For coherent states, the matrix elements of C , S , C^2 , S^2 are given by

$$I_c \equiv \langle N\varphi | C | N\varphi \rangle = I(N) \cos \varphi, \\ I_s \equiv \langle N\varphi | S | N\varphi \rangle = I(N) \sin \varphi,$$

$$I(N) \equiv N^{\frac{1}{2}} e^{-N} \sum_n \frac{N^n}{n!(n+1)^{\frac{1}{2}}}, \quad (17a)$$

$$J_c \equiv \langle N\varphi | C^2 | N\varphi \rangle \\ = \frac{1}{2} - \frac{1}{4}e^{-N} + \frac{1}{2}(\cos^2 \varphi - \sin^2 \varphi)J(N),$$

$$J_s \equiv \langle N\varphi | S^2 | N\varphi \rangle \\ = \frac{1}{2} - \frac{1}{4}e^{-N} - \frac{1}{2}(\cos^2 \varphi - \sin^2 \varphi)J(N),$$

$$J(N) \equiv N e^{-N} \sum_{n=0}^{\infty} \frac{N^n}{n!((n+1)(n+2))^{\frac{1}{2}}}. \quad (17b)$$

The functions $I(N)$ and $J(N)$ have the asymptotic (large N) forms

$$I(N) \approx 1 - \frac{1}{8N}, \quad J(N) \approx 1 - \frac{1}{2N}. \quad (17c)$$

Hence for large N ,

$$\begin{aligned} I_c &\approx \cos \varphi, & I_s &\approx \sin \varphi, \\ J_c &\approx \cos^2 \varphi, & J_s &\approx \sin^2 \varphi. \end{aligned} \quad (17d)$$

It is seen that the expressions involving S may be obtained from those involving C by replacing φ by $\frac{1}{2}\pi - \varphi$.

The uncertainty relations proposed by CN are

$$\begin{aligned} U_1(\Psi) &\equiv (\Delta N)^2(\Delta C)^2/\langle S \rangle^2 \geq \frac{1}{4}, \\ U_2(\Psi) &\equiv (\Delta N)^2(\Delta S)^2/\langle C \rangle^2 \geq \frac{1}{4}, \\ U_3(\Psi) &\equiv (\Delta N)^2 \frac{(\Delta S)^2 + (\Delta C)^2}{[\langle S \rangle^2 + \langle C \rangle^2]} \geq \frac{1}{4}. \end{aligned} \quad (18)$$

These relations have the virtues that (i) they represent plausible generalizations of the imprecise statement $(\Delta N)^2(\Delta \varphi)^2 \geq \frac{1}{4}$; (ii) for highly excited coherent states they closely approach their theoretical lower limit $\frac{1}{4}$ and remain small for moderate excitations. The last uncertainty product is independent of φ when evaluated with coherent states.

B. It is demonstrably true that the coherent states do not permit the U_i 's to attain their theoretical lower limit $\frac{1}{4}$. It may nevertheless be the case that no normalizable states exist for which $U_i = \frac{1}{4}$; and the coherent states give the lowest attainable minimum. We establish that (i) the coherent states are not critical states, viz., they do not make the uncertainty products stationary; therefore *a fortiori* they do not minimize the uncertainty products; and (ii) normalizable states exist which allow some of the U_i 's to reach their theoretical lower limit of $\frac{1}{4}$.

C. We first study U_1 . The critical states, which make U_1 stationary, satisfy according to (12b)

$$0 = \left[\frac{[N_{\text{op}} - \langle N \rangle]^2}{(\Delta N)^2} + \frac{[C - \langle C \rangle]^2}{(\Delta C)^2} - \frac{2S}{\langle S \rangle} \right] |\Psi\rangle. \quad (19)$$

Expanding $|\Psi\rangle$ in number states $|\Psi\rangle = \sum_n a_n |n\rangle$, we find that the coefficients a_n must satisfy the recursion relation

$$\begin{aligned} \frac{1}{4}a_2 + \left(\frac{ib^2}{\gamma} - \beta \right) a_1 + \left(\frac{b^2}{a^2} \alpha^2 + \frac{1}{4} + \beta^2 \right) a_0 &= 0, \\ \frac{1}{4}a_{n+2} + \left(\frac{ib^2}{\gamma} - \beta \right) a_{n+1} + \left(\frac{b^2}{a^2} (n - \alpha)^2 + \frac{1}{2} + \beta^2 \right) a_n \\ - \left(\frac{ib^2}{\gamma} + \beta \right) a_{n-1} + \frac{1}{4}a_{n-2} &= 0 \quad n \geq 1, \quad a_{-1} = 0, \end{aligned} \quad (20a)$$

subject to the subsidiary conditions

$$\begin{aligned} 1 &= \langle \Psi | \Psi \rangle, & \alpha &\equiv \langle \Psi | N_{\text{op}} | \Psi \rangle, \\ \beta &\equiv \langle \Psi | C | \Psi \rangle, & \gamma &\equiv \langle \Psi | S | \Psi \rangle, \\ a^2 &\equiv (\Delta N)^2 = \langle \Psi | N_{\text{op}}^2 | \Psi \rangle - \alpha^2, \\ b^2 &\equiv (\Delta C)^2 = \langle \Psi | C^2 | \Psi \rangle - \beta^2. \end{aligned} \quad (20b)$$

Whether the coherent states satisfy these equations can be easily checked by setting $|\Psi\rangle = |N\varphi\rangle$; viz.,

$$a_n = e^{-\frac{1}{2}N} [N^{\frac{1}{2}n} e^{in\varphi} / (n!)^{\frac{1}{2}}], \quad \alpha = a^2 = N.$$

For simplicity we also assume $\varphi = \frac{1}{2}\pi$; viz., $\beta = 0$. Then (20a) becomes

$$\frac{N}{\sqrt{2}} + \frac{4b^2}{\gamma} N^{\frac{1}{2}} - (4b^2N + 1) = 0, \quad (21a)$$

$$\begin{aligned} \frac{N^2}{((n+2)(n+1)n(n-1))^{\frac{1}{2}}} + 4 \frac{b^2 N^{\frac{3}{2}}}{\gamma((n+1)n(n-1))^{\frac{1}{2}}} \\ - \left(\frac{4b^2}{N} [n - N]^2 + 2 \right) \frac{N}{(n(n+1))^{\frac{1}{2}}} \\ + \frac{4b^2}{\gamma} \frac{N^{\frac{1}{2}}}{(n)^{\frac{1}{2}}} + 1 = 0. \end{aligned} \quad (21b)$$

Equations (21) are manifestly not satisfied; hence the coherent states are not critical states, and do not minimize the uncertainty product.

We may however demonstrate that for large N the coherent states do satisfy (21b) approximately. For large N , $(\Delta N)^2(\Delta C)^2 \sim \frac{1}{4}\langle S \rangle^2 \sim \frac{1}{4}\sin^2 \varphi = \frac{1}{4}$; thus $4b^2 = 4(\Delta C)^2 \sim \langle S \rangle^2 / \langle \Delta N \rangle^2 \sim 1/N$; and $\gamma \sim \sin \varphi = 1$ [see (17)]. Also for large N the most important number states $|n\rangle$, contributing to the coherent state $|N\varphi\rangle$, are those with $n \sim N$; since for these values of n , $N^{\frac{1}{2}}(n!)^{-\frac{1}{2}}$ assumes its maximum. Therefore for large N and $n \sim N$ the left-hand side of (21b) becomes $O(1/N)$, and (21b) is approximately satisfied. This argument cannot be given when $\langle C \rangle \equiv \beta \neq 0$.

It is evident that the analysis of U_2 proceeds in the same fashion towards the same conclusion, except that the condition $\langle C \rangle = 0$ is now replaced by $\langle S \rangle = 0$.

D. States that do minimize the uncertainty product U_1 and allow it to achieve its theoretical lower limit of $\frac{1}{4}$ may be easily constructed (under certain restrictions). The discussion in Sec. IIE shows that we may use the direct method to determine these states. Accordingly we wish to solve

$$\begin{aligned} (N_{\text{op}} + i\gamma C) |\Psi\rangle &= \lambda |\Psi\rangle, \\ \langle \Psi | \Psi \rangle &= 1. \end{aligned} \quad (22)$$

For simplicity, we again confine ourselves to the case $\langle \Psi | C | \Psi \rangle = 0$. This makes λ real and equal to $\langle \Psi | N_{\text{op}} | \Psi \rangle$. Expanding in number states leads to the recursion relation

$$\begin{aligned} (\lambda - n)a_n &= \frac{1}{2}i\gamma(a_{n+1} + a_{n-1}), \\ a_{-1} &= 0. \end{aligned} \quad (23a)$$

To put this in a more transparent form, we define

$a_n = (-i)^n b_n$; then the b_n 's satisfy

$$(2/\gamma)(n - \lambda)b_n = (b_{n-1} - b_{n+1}),$$

$$b_{-1} = 0. \quad (23b)$$

This recursion relation is well known.⁷ We do not examine it in detail as it is sufficient for our purposes to extract one solution. A solution to (23b) is

$$b_n = \nu I_{n-\lambda}(\gamma), \quad (24a)$$

where $I_\mu(Z)$ is a modified Bessel function of the first kind of order μ .⁸ We also require $I_{-1-\lambda}(\gamma) = 0$. (This forces λ to satisfy $2k + 1 > \lambda > 2k$, where $k = 0, 1, \dots$.) The multiplicative constant ν is obtained from the normalization condition

$$|\nu|^2 \sum_{n=0}^{\infty} I_{n-\lambda}^2(\gamma) = 1. \quad (24b)$$

In Appendix B we prove that the series in (24b) converges, that $\langle \Psi | C | \Psi \rangle = 0$, and $\langle \Psi | S | \Psi \rangle \neq 0$. Thus the desired solution of (22) for which $U_1 = \frac{1}{4}$, is

$$|\Psi\rangle = \nu \sum_{n=0}^{\infty} (-i)^n I_{n-\lambda}(\gamma) |n\rangle,$$

$$\lambda = \langle N_{\text{op}} \rangle, \quad (25)$$

$$(\Delta N)^2 (\Delta C)^2 / \langle S \rangle^2 = \frac{1}{4}.$$

Unfortunately these states do not seem to have any physical significance.

The recursion relation (23b) is also solved by the number states. These however do not minimize U_1 , as we demonstrate explicitly in Appendix A.

It is clear that states which allow U_2 to reach $\frac{1}{4}$ can also be constructed.

E. We now examine the symmetric uncertainty product U_3 . We first show that no states exist for which U_3 attains its minimum value of $\frac{1}{4}$. According to (3a) we have

$$(\Delta N)^2 (\Delta C)^2 = \frac{1}{4} \langle S \rangle^2 + \frac{1}{4} Q_1(\Psi) + R_1(\Psi),$$

$$(\Delta N)^2 (\Delta S)^2 = \frac{1}{4} \langle C \rangle^2 + \frac{1}{4} Q_2(\Psi) + R_2(\Psi). \quad (26a)$$

Therefore for U_3 to be $\frac{1}{4}$ we must have

$$0 = U_3 - \frac{1}{4} = [\langle S \rangle^2 + \langle C \rangle^2]^{-1} [\frac{1}{4} Q_1(\Psi) + \frac{1}{4} Q_2(\Psi) + R_1(\Psi) + R_2(\Psi)]. \quad (26b)$$

Since each term on the right-hand side is positive semidefinite, $Q_{1,2}$ and $R_{1,2}$ must vanish separately, which according to (5b) requires

$$[N_{\text{op}} + i\gamma_1 C] |\Psi\rangle = \lambda_1 |\Psi\rangle,$$

$$[N_{\text{op}} + i\gamma_2 S] |\Psi\rangle = \lambda_2 |\Psi\rangle, \quad (26c)$$

where γ_1 and γ_2 are real and nonzero. Evidently the commutator of the operators appearing on the left-hand side of (26c) must annihilate the state $|\Psi\rangle$. This sets the condition

$$\left[\frac{1}{\gamma_2} S + \frac{1}{\gamma_1} C + \frac{i}{2} P \right] |\Psi\rangle = 0. \quad (26d)$$

Equation (26c) may be used again to evaluate $S|\Psi\rangle$ and $C|\Psi\rangle$. Therefore (26d) becomes

$$[\gamma_1^{-2}(N_{\text{op}} - \lambda_1) + \gamma_2^{-2}(N_{\text{op}} - \lambda_2) + \frac{1}{2}P] |\Psi\rangle = 0. \quad (26e)$$

Expanding $|\Psi\rangle$ in number states yields the conditions

$$\left[-\frac{\lambda_1}{\gamma_1^2} - \frac{\lambda_2}{\gamma_2^2} + \frac{1}{2} \right] a_0 = 0, \quad (26f)$$

$$[\gamma_1^{-2}(n - \lambda_1) + \gamma_2^{-2}(n - \lambda_2)] a_n = 0.$$

These recursion relations are solved only by the number state $|\lambda\rangle$, where $\lambda_1 = \lambda_2 = \lambda = \text{integer}$; $\langle C \rangle = \langle S \rangle = 0$. We demonstrate in Appendix A that the number states do not minimize U_3 .

Thus the direct method does not yield any solutions, and we are led to consider U_3 by the analytic method. The procedure to follow is the same as for U_1 . The variation principle gives an EL equation which represents a necessary condition which must be satisfied if U_3 is to be minimized. With this condition, it can easily be verified that the coherent states are not critical states and do not minimize U_3 . Again it is found that, for large N , the coherent states approximately satisfy the necessary condition, but now no restriction is set on φ . The EL equation is too complicated to serve to determine the states that do minimize U_3 ; hence we do not present the details of this calculation. (The recursion relation which follows from the EL equation actually is elementary, but the imposition of the subsidiary conditions is complicated. In any case the solution, if it exists, surely has no physical significance.)

Since the first variation of U_3 does not vanish for coherent states, there exist states, arbitrarily close to the coherent states, for which U_3 is smaller than it is when evaluated with coherent states. For example, with the state $|\Psi_1\rangle = \omega[|N\varphi\rangle + \epsilon e^{-\frac{1}{2}N} |0\rangle]$, where ω is a normalization factor, ϵ a positive small parameter, and N large, $U_3(\Psi_1)$ is smaller by an amount $2\epsilon N^{\frac{3}{2}} e^{-N}$ than with the coherent state of the same excitation.

IV. SINE-COSINE UNCERTAINTY PRODUCT

Since S and C do not commute, limitations exist on the simultaneous measurement of these two quantities. However in the classical limit these limitations

⁷ G. N. Watson, *A Treatise on the Theory of Bessel Functions* (Cambridge University Press, London, 1952), p. 294.

⁸ Reference 7, p. 172.

must disappear. Thus we are led to consider the uncertainty product

$$U_4 = (\Delta S)^2(\Delta C)^2. \quad (27)$$

In their original discussion of the S and C operators, SG demonstrated explicitly that there exist *unnormalizable* states

$$|\theta\rangle = \sum_{n=0}^{\infty} e^{in\theta} |n\rangle \quad (28)$$

for which $(\Delta S)^2 = 0 = (\Delta C)^2$, hence $U_4 = 0$. Carruthers and Nieto have shown that for the normalizable coherent states U_4 goes rapidly to zero for large N , and to $\frac{1}{16}$ for small N .

We now wish to use our analytic method to investigate whether there exist normalizable states which minimize U_4 , and whether the coherent states are critical states for U_4 . We find that no normalizable critical states exist.

To establish this result, we use (8d)

$$\left[\frac{S^2}{(\Delta S)^2} + \frac{C^2}{(\Delta C)^2} - 2 \right] |\Psi\rangle = 0. \quad (29a)$$

For simplicity we confine ourselves to the symmetric case

$$\langle S \rangle = \langle C \rangle, \quad \langle S^2 \rangle = \langle C^2 \rangle.$$

Thus we need to solve

$$[(S - \alpha)^2/a^2 + (C - \alpha)^2/a^2 - 2] |\Psi\rangle = 0, \quad (29b)$$

with the subsidiary conditions

$$\begin{aligned} \langle \Psi | \Psi \rangle &= 1, \quad \alpha = \langle S \rangle = \langle C \rangle, \\ a^2 = \langle S^2 \rangle - \alpha^2 &= \langle C^2 \rangle - \alpha^2. \end{aligned} \quad (29c)$$

Equation (29b) may be simplified into

$$[v - \frac{1}{2}P - ue^{\frac{1}{2}i\pi}E_+ - ue^{-\frac{1}{2}i\pi}E_-] |\Psi\rangle = 0, \quad (29d)$$

where

$$\begin{aligned} u &= \sqrt{2}\alpha, \\ v &= 1 + u^2 - 2a^2. \end{aligned} \quad (29e)$$

Expanding in number states gives the recursions

$$(v - \frac{1}{2})a_0 - ue^{-\frac{1}{2}i\pi}a_1 = 0, \quad (30a)$$

$$va_n - u(e^{\frac{1}{2}i\pi}a_{n-1} + e^{-\frac{1}{2}i\pi}a_{n+1}) = 0, \quad n \geq 1. \quad (30b)$$

This is obviously not satisfied by the coherent states, except approximately for large N and $n \sim N$. The general solution of (30b) is given by

$$a_n = e^{i(\pi/4)n} [Ap^n + Bp^{-n}], \quad (31a)$$

$$\frac{v}{u} = p + \frac{1}{p}. \quad (31b)$$

We assume $v \neq 0$. If $|\Psi\rangle$ is to be normalizable, we

must have $\sum_n |a_n|^2 = 1$; therefore A is zero if $|p| > 1$, and B is zero if $|p| < 1$. Taking the latter case and imposing (30a) and (31b) determines $p = 2u$ and sets $v = \frac{1}{2} + 2u^2$. The normalization can now be determined, and we obtain as a solution

$$|\Psi\rangle = (1 - 4u^2)^{\frac{1}{2}} \sum_{n=0}^{\infty} e^{i(\pi/4)n} (2n)^n |n\rangle, \quad 4u^2 < 1. \quad (32a)$$

Imposing now the subsidiary condition

$$u/\sqrt{2} = \alpha = \langle \Psi | C | \Psi \rangle = \sqrt{2}u \quad (32b)$$

gives $u = 0$ and no nontrivial normalizable solution is obtained.

Similarly when $4u^2 > 1$, no normalizable solution is obtained. Therefore we conclude that U_4 cannot be minimized by normalizable states; and only the unnormalizable states (28) minimize U_4 . For these states $\langle N_{op} \rangle$ diverges and they obviously represent the high excitation limit. [Such states cannot be determined by the EL equation (29d) since that equation was derived under the assumption that the solutions are normalizable.]

V. SUMMARY

In conclusion, we summarize our results. We have developed new variational techniques for the determination of states that minimize the uncertainty product of operators. By the use of these techniques we have demonstrated that the coherent states do not minimize the various uncertainty relations which can be given for number and phase operators. Although normalizable states do exist that minimize some of the number-phase uncertainty products, we do not believe that these states have any physical significance as they are strongly dependent on the specific form of the uncertainty product. Moreover, the coherent states do not even make any of the uncertainty products stationary. Thus the coherent states have no unique relevance to the classical limit of the phase operators. Indeed any state, which, when expanded in number states $|n\rangle$, has expansion coefficients a_n , which for large $\langle N_{op} \rangle$ are strongly peaked and constant at $n \sim \langle N_{op} \rangle$, serves to minimize approximately the uncertainty products. An example of such a state was given at the end of Sec. III E.

APPENDIX A

Throughout our analysis we have ignored the fact that the number states are solutions to some of the various equations we studied. We examine here whether these eigenstates of N_{op} minimize the various uncertainty products U_1, U_2, U_3 [Eq. (18)].

For number states, of course, $(\Delta N)^2 = 0$, $\langle S \rangle^2 = 0 = \langle C \rangle^2$, while $(\Delta S)^2 \neq 0 \neq (\Delta C)^2$. Therefore the

uncertainty products have the indeterminate form 0/0. To obtain the value for this we proceed as follows. Consider the *excited coherent state*

$$|N\varphi n\rangle = (E_+)^n |N\varphi\rangle. \quad (\text{A1})$$

These states are normalized and the number states $|n\rangle$ are reached as $N \rightarrow 0$ in these states. We therefore evaluate all matrix elements with the states $|N\varphi n\rangle$ and then let $N \rightarrow 0$.

The relevant matrix elements are easily evaluated. One finds

$$\begin{aligned} \langle N\varphi n | N_{\text{op}} | N\varphi n \rangle &= N + n, \\ \langle N\varphi n | N_{\text{op}}^2 | N\varphi n \rangle &= (N + n)^2 + N, \\ \langle N\varphi n | T | N\varphi n \rangle &= \langle N\varphi | T | N\varphi \rangle, \end{aligned} \quad (\text{A2})$$

where T is any of the operators S , S^2 , C , C^2 . From these it follows that

$$U_i(N\varphi n) = U_i(N\varphi) \quad (\text{A3})$$

and

$$U_i(n) = \lim_{N \rightarrow 0} U_i(N\varphi). \quad (\text{A4})$$

According to the definitions of the U_i 's [Eq. (18)] and using the formulas (17) for the matrix elements of the S and C operators between coherent states, we have

$$\begin{aligned} U_1(N\varphi) &= N[I^2(N) \sin^2 \varphi]^{-1} \\ &\times [\tfrac{1}{2} - \tfrac{1}{4}e^{-N} + \tfrac{1}{2}(\cos^2 \varphi - \sin^2 \varphi) \\ &\times J(N) - I^2(N) \cos^2 \varphi], \\ U_2(N\varphi) &= U_1(\tfrac{1}{2}N\pi - \varphi), \\ U_3(N\varphi) &= U_1(\tfrac{1}{4}N\pi). \end{aligned} \quad (\text{A5})$$

For small N , $I^2 \sim N$ and $J \sim N/\sqrt{2}$. Therefore

$$\begin{aligned} U_1(n) &= 1/4 \sin^2 \varphi, \\ U_2(n) &= 1/4 \cos^2 \varphi, \\ U_3(n) &= \tfrac{1}{2}. \end{aligned} \quad (\text{A6})$$

It is seen that $U_1(N\varphi n)$ and $U_2(N\varphi n)$ do not approach

a unique limit, viz., a limit independent of φ ; therefore $U_1(n)$ and $U_2(n)$ do not exist. For U_3 we conclude that either $U_3(n) = \tfrac{1}{2}$, or if there exist other ways of approaching the number states, leading to a different value of $U_3(n)$, the limit does not exist. In any case, the number states do not minimize the uncertainty products.

APPENDIX B

We wish to prove that the state $|\Psi\rangle$,

$$|\Psi\rangle = \nu \sum_{n=0}^{\infty} (-i)^n I_{n-\lambda}(\gamma) |n\rangle, \quad (\text{B1})$$

is normalizable, viz.,

$$\sum_{n=0}^{\infty} I_{n-\lambda}^2(\gamma) < \infty. \quad (\text{B2})$$

For large enough n , $n - \lambda$ is positive and the following integral representation for I_μ is valid⁸:

$$I_\mu(Z) = \frac{(\tfrac{1}{2}Z)^\mu}{\pi^{1/2} \Gamma(\mu + \tfrac{1}{2})} \int_{-1}^1 (1 - t^2)^{\mu - \frac{1}{2}} e^{\pm 2it} dt; \quad \text{Re } \mu > -\tfrac{1}{2}. \quad (\text{B3})$$

Evidently

$$I_{n-\lambda}^2(\gamma) \leq \frac{(\tfrac{1}{4}\gamma^2)^{n-\lambda}}{[\Gamma(n - \lambda + \tfrac{1}{2})]^2} M(\gamma), \quad n > \lambda - \tfrac{1}{2}, \quad (\text{B4})$$

where M is positive and independent of n and λ . Therefore the series (B2) converges.

We also need to show that $\langle \Psi | C | \Psi \rangle = 0$. This is readily established from the formulas (16a) and (B1).

Finally we establish that $\langle \Psi | S | \Psi \rangle \neq 0$. Recall that S is proportional to the commutator of N_{op} and C . According to the general discussion of Sec. IIC, we know that the expectation value of the commutator is proportional to $\gamma(\Delta C)^2$. Since γ is nonzero, we may prove that $\langle S \rangle$ is nonzero by showing that $(\Delta C)^2$ does not vanish. However $(\Delta C)^2 = \langle C^2 \rangle$ since $\langle C \rangle = 0$. But $\langle C^2 \rangle$ is nonzero since the operator C manifestly does annihilate $|\Psi\rangle$.

Onset of ODLRO and the Phase Transition in the Ideal Boson Gas

MARSHALL LUBAN*

Department of Physics, Bar-Ilan University, Ramat Gan, Israel

AND

MICHAEL REVZEN

Department of Physics, Technion-Israel Institute of Technology, Haifa, Israel

(Received 6 September 1967)

The divergence of the constant-pressure specific heat C_P and the isothermal compressibility K_T as one lowers the temperature of the ideal boson gas to the transition temperature T_c is discussed in terms of the onset of off-diagonal long-range order (ODLRO) of the one-particle density matrix ρ_1 .

The phase transition which occurs in the ideal boson gas is well known to every student of statistical mechanics.¹ If N identical bosons, each of mass m , are in a container of volume V , then below the critical temperature $T_c = 3.313\hbar^2\rho^{2/3}/(k_B m)$, where k_B is Boltzmann's constant, a finite fraction N_0/N of the particles have zero momentum in the thermodynamic limit ($N, V \rightarrow \infty$ with $\rho = N/V$ held fixed). This so-called Einstein condensation reflects itself in the behavior of the thermodynamic functions, for example, that the specific heat at constant volume C_V is a continuous function of T but has a discontinuous derivative at T_c . Although the mathematical derivation of these and other well-known properties of the system is straightforward, it is our contention that the physical basis for some of these properties is as yet not fully understood. As an example we cite the easily derived results that the specific heat at constant pressure C_P and the isothermal compressibility K_T diverge as $(T - T_c)^{-1}$ when T is lowered to T_c for fixed ρ .² A satisfactory physical explanation must clarify the origin of this behavior of C_P and K_T which occurs even though N_0/N is strictly zero in the temperature range in question.

Several years ago Yang³ suggested that in a quantum many-body system it is possible to have off-diagonal long-range order (ODLRO) of certain reduced-density matrices in the coordinate representation, and that this order characterizes a new thermodynamic phase of quantum-mechanical origin. A second and vital plank of the program is the requirement that the thermodynamic functions must explicitly reflect

the existence of ODLRO. In this paper we generalize Yang's program to write several thermodynamic functions describing the ideal boson gas for temperatures just above T_c (where there is no ODLRO) in a form so that they explicitly reflect the onset of ODLRO as $T \rightarrow T_c + 0$. In particular, we recall a simple relationship for this system between the pair correlation function and the one-particle density matrix ρ_1 , and we use it together with a fluctuation-dissipation theorem to relate the divergence of K_T and C_P as $T \rightarrow T_c + 0$ to the onset of ODLRO in ρ_1 .

For a system of bosons of number density ρ in thermodynamic equilibrium and described by the grand canonical ensemble, the one-particle density matrix in the coordinate representation and the pair distribution function are defined as

$$\rho_1(\mathbf{r}, \mathbf{r}') = \langle \psi^\dagger(\mathbf{r})\psi(\mathbf{r}') \rangle, \quad (1)$$

$$g(\mathbf{r}, \mathbf{r}') = \rho^{-2} \langle \psi^\dagger(\mathbf{r})\psi^\dagger(\mathbf{r}')\psi(\mathbf{r}')\psi(\mathbf{r}) \rangle, \quad (2)$$

respectively. In these equations $\psi(\mathbf{r})$ denotes a field operator satisfying the Bose-Einstein commutation relations

$$[\psi(\mathbf{r}), \psi^\dagger(\mathbf{r}')] = \delta^3(\mathbf{r} - \mathbf{r}'), \quad (3a)$$

$$[\psi(\mathbf{r}), \psi(\mathbf{r}')] = 0; \quad (3b)$$

the average of any operator \mathcal{O} is given as

$$\langle \mathcal{O} \rangle = (Z_G)^{-1} \text{Tr} [e^{-\beta(H - \mu N)} \mathcal{O}], \quad \beta = (k_B T)^{-1}, \quad (4)$$

where the quantities Z_G , μ , H , and N denote the grand partition function,

$$Z_G = \text{Tr} e^{-\beta(H - \mu N)}, \quad (5)$$

the chemical potential, Hamiltonian, and total particle number operator of the system, respectively. For a translationally invariant system, ρ_1 and g are functions of $\mathbf{r} - \mathbf{r}'$, and if one employs orthonormalized single-particle plane-wave states $\varphi_{\mathbf{k}} = V^{-1/2} \exp(i\mathbf{k} \cdot \mathbf{r})$ satisfying periodic boundary conditions with respect to the normalization volume V , then (1) and (2) can be

* Work performed while the author was a Fellow of the John Simon Guggenheim Memorial Foundation and on leave of absence from the University of Pennsylvania, Philadelphia, Pa.

¹ For example, see K. Huang, *Statistical Mechanics* (John Wiley & Sons, Inc., New York, 1963), Chap. 12.

² F. London, *Superfluids* (John Wiley & Sons, Inc., New York, 1954), Vol. II, p. 53, obtains C_P by taking the appropriate derivative of the enthalpy and shows that the former diverges as $T \rightarrow T_c + 0$. The explicit form of the divergence is not given by London but it is trivial to calculate using several of his formulas.

³ C. N. Yang, *Rev. Mod. Phys.* **34**, 694 (1962).

written as

$$\rho_1(\mathbf{r} - \mathbf{r}') = V^{-1} \sum_{\mathbf{k}} \langle N_{\mathbf{k}} \rangle e^{i\mathbf{k} \cdot (\mathbf{r} - \mathbf{r}')} \quad (6)$$

$$g(\mathbf{r} - \mathbf{r}') = \langle N \rangle^{-2} \sum_{\mathbf{p}, \mathbf{q}, \mathbf{k}} \langle a_{\mathbf{p}}^\dagger a_{\mathbf{q}}^\dagger a_{\mathbf{q}-\mathbf{k}} a_{\mathbf{p}+\mathbf{k}} \rangle e^{i\mathbf{k} \cdot (\mathbf{r} - \mathbf{r}')} \quad (7)$$

In these equations $a_{\mathbf{k}}^\dagger$ and $a_{\mathbf{k}}$ are the usual boson creation and destruction operators for the state $\varphi_{\mathbf{k}}$, and $N_{\mathbf{k}} = a_{\mathbf{k}}^\dagger a_{\mathbf{k}}$. Integration of (7) over the volume V gives the relation

$$1 + \rho \int d^3r [g(\mathbf{r}) - 1] = \frac{1}{\langle N \rangle} (\langle N^2 \rangle - \langle N \rangle^2) \quad (8)$$

Now the right-hand side is related to the isothermal compressibility K_T by the well-known relation⁴

$$\langle N^2 \rangle - \langle N \rangle^2 = \langle N \rangle k_B T \rho K_T \quad (9)$$

so that we obtain the formula

$$k_B T \rho K_T = 1 + \rho \int d^3r [g(\mathbf{r}) - 1] \quad (10)$$

In the case of the ideal boson gas, nonzero contributions to the sums in (7) arise only when $\mathbf{k} = 0$ and $\mathbf{k} = \mathbf{q} - \mathbf{p}$. It is then straightforward to show that

$$\begin{aligned} g(\mathbf{r}) &= 1 + (\rho_1(\mathbf{r})/\rho)^2 \\ &+ \langle N \rangle^{-2} \sum_{\mathbf{k}} (\langle N_{\mathbf{k}}^2 \rangle - 2\langle N_{\mathbf{k}} \rangle^2 - \langle N_{\mathbf{k}} \rangle) \\ &= 1 + (\rho_1(\mathbf{r})/\rho)^2 \end{aligned} \quad (11)$$

The second equality is a consequence of the relations

$$\langle N_{\mathbf{k}}^2 \rangle - 2\langle N_{\mathbf{k}} \rangle^2 - \langle N_{\mathbf{k}} \rangle = 0, \quad (12)$$

$$\langle N_{\mathbf{k}} \rangle = (e^{\beta \epsilon_{\mathbf{k}}} - 1)^{-1}, \quad \epsilon_{\mathbf{k}} = (2m)^{-1} \hbar^2 k^2 - \mu, \quad (13)$$

which holds for this system. In the thermodynamic limit and for temperatures $T > T_c$, the function $\rho_1(\mathbf{r})$ is given as⁵

$$\rho_1(\mathbf{r}) = \frac{1}{\pi \lambda^3} \left(\frac{\lambda}{r} \right) \int_0^\infty dy \frac{\sin(2\pi^{\frac{1}{2}} y^{\frac{1}{2}} r/\lambda)}{e^{y-\beta\mu} - 1}, \quad (14)$$

where $\lambda = (2\pi \hbar^2 \beta/m)^{\frac{1}{2}}$ is the thermal de Broglie wavelength. In the regime $r \ll \lambda$, $\rho_1(\mathbf{r})$ as defined by (6) reduces to ρ and indeed (14) correctly reduces to the standard relation

$$\rho = \frac{2}{\pi^{\frac{1}{2}}} \frac{1}{\lambda^3} \int_0^\infty dy \frac{y^{\frac{1}{2}}}{e^{y-\beta\mu} - 1}, \quad (15)$$

which applies for $T > T_c$. In the opposite limit, $r \gg \lambda$, the dominant contribution to the integrand in (14) arises from small y and so⁶

$$\begin{aligned} \rho_1(\mathbf{r}) &\sim \frac{1}{\pi \lambda^3} \frac{\lambda}{r} \int_0^\infty dy \frac{1}{y - \beta\mu} \sin(2\pi^{\frac{1}{2}} y^{\frac{1}{2}} r/\lambda) \\ &= \lambda^{-3} (\lambda/r) \exp[-2\pi^{\frac{1}{2}} (-\beta\mu)^{\frac{1}{2}} r/\lambda]. \end{aligned} \quad (16)$$

⁴ See p. 167 of Ref. 1.

⁵ For $T < T_c$ the $\mathbf{k} = 0$ term of the sum in (5) is nonvanishing in the thermodynamic limit and thus ρ_1 is given as the sum of $\langle N_0 \rangle/V$ plus the integral of (14). Also see (20) below.

⁶ The next term in the asymptotic expansion of $\rho_1(r)$ is $2\lambda^{-3} (\lambda/r) e^{-2\pi r/\lambda} \cos(2\pi r/\lambda)$ and thus may be ignored for the purpose of obtaining the leading term of the asymptotic expansion of $g(r) - 1 = \rho_1(r)^2$.

Now, as shown in the Appendix, in the temperature range $0 < T - T_c \ll T_c$

$$(-\beta\mu)^{\frac{1}{2}} = 1.105(T - T_c)/T_c, \quad (17)$$

and so we see that the range of $\rho_1(\mathbf{r})$ increases as $(T - T_c)^{-1}$ as $T \rightarrow T_c + 0$. These formulas quantitatively describe the onset of ODLRO in the density matrix ρ_1 . Using (11), (16), and (17), we have⁷

$$g(\mathbf{r}) \sim 1 + 0.146(\lambda/r)^2 \exp[-7.84(r/\lambda)(T - T_c)/T_c] \quad (18)$$

for large r . This relation combined with (10) shows explicitly that K_T diverges as $(T - T_c)^{-1}$ for $T \rightarrow T_c + 0$. Finally, by recalling the thermodynamic relation

$$C_P = C_V + vT(\partial P/\partial T)_v^2 K_T, \quad (v = 1/\rho), \quad (19)$$

we see that C_P also diverges as $(T - T_c)^{-1}$ since C_V and $(\partial P/\partial T)_v$ remain finite.

Below T_c the term $\mathbf{k} = 0$ on the right-hand side of (5) does not vanish in the thermodynamic limit and we have

$$\rho_1(\mathbf{r}) \sim (\langle N_0 \rangle/V) + \lambda^{-3} (\lambda/r) \quad (20)$$

for large r . That is, $\lim_{r \rightarrow \infty} \rho_1(\mathbf{r}) = \langle N_0 \rangle/V \neq 0$ and thus ρ_1 displays ODLRO.

The above treatment provides a new way of viewing the origin of the phase transition in the ideal boson gas. We have seen that the onset of the phase transition as T is lowered to T_c , and in particular the singular behavior of K_T and C_P , can be described in terms of the onset of ODLRO in ρ_1 . In turn, ODLRO and its onset in ρ_1 is a particularly useful concept, for it enables one to give quantitative expression to the intuitive idea that there exists "order" in the system in the quantum regime, that is, for T below and just above T_c . The order, of course, is the manifestation of the requirement of using symmetric wavefunctions to describe a boson system, and it persists until masked by the increasing disorder of the heat reservoir as T is raised to the vicinity of T_c .

A useful intuitive concept in discussing the behavior of a ferromagnetic spin system above the Curie point T_c is that of the onset of long-range spin order as $T \rightarrow T_c + 0$. In the case of the two-dimensional Ising lattice, the onset of long-range spin order is accompanied by a logarithmic divergence of the specific heat. The onset of ODLRO in ρ_1 , although

⁷ It is worth recalling that the Ornstein-Zernicke result for the large r behavior of the pair distribution function of a classical fluid near the critical point is $g(r) - 1 \sim (A/r)e^{-r/R(T)}$, where the range parameter $R(T) \propto (T - T_c)^{-\frac{1}{2}}$. On the other hand, the Landau theory applied to such a fluid predicts $K_T \sim B(T - T_c)^{-1}$ for $T \rightarrow T_c + 0$ along the critical isochore. See M. Fisher, *J. Math. Phys.* **5**, 944 (1964), Eqs. (1.4), (3.15), and (3.19); L. P. Kadanoff *et al.*, *Rev. Mod. Phys.* **39**, 395 (1967), Sec. II.

more abstract, is the analogous concept for the ideal boson gas.

The lambda transition in liquid ^4He also features a divergence of C_P and K_T as one lowers T to the lambda T_λ , although in that case the divergences are logarithmic instead of the form $(T - T_\lambda)^{-1}$. Presumably these divergences also reflect the onset of ODLRO in ρ_1 but with a less rapid increase as $T \rightarrow T_\lambda + 0$ than for the ideal boson gas.⁸

ACKNOWLEDGMENT

A portion of the work of one of the authors (M. L.) was performed at the Department of Physics of the

⁸ An alternate suggestion for explaining the divergence of C_P and K_T for $T \rightarrow T_\lambda + 0$ in the case of liquid helium has been advanced by W. H. Keesom and A. P. Keesom, *Physica* **2**, 557 (1935). According to this suggestion, in the vicinity of T_λ the liquid is to be thought of as divided into small, weakly interacting domains which, because of thermal fluctuations, can have different temperatures. The domains with temperatures less than T_λ give the anomalously large contributions to the specific heat. The net effect is to cause the specific heat in the vicinity of T_λ to be a symmetric function of $T - T_\lambda$.

Hebrew University of Jerusalem, and he wishes to acknowledge their kind hospitality.

APPENDIX

We derive here (17) of the text, giving μ as a function of T just above T_c . Evaluating (15) for $T = T_c$, where $\mu = 0$, one obtains

$$\rho = \zeta\left(\frac{3}{2}\right) \left(\frac{mk_B T_c}{2\pi\hbar^2}\right)^{\frac{3}{2}}, \quad (\text{A1})$$

where $\zeta\left(\frac{3}{2}\right) = 2.612$. Just above T_c the value of μ is very small and thus we can approximate the right-hand side of (15) as

$$\begin{aligned} \rho &\simeq \frac{2}{\pi^{\frac{1}{2}}} \frac{1}{\lambda^3} \int_0^\infty dx \frac{x^{\frac{1}{2}}}{e^x - 1} - \frac{2}{\pi^{\frac{1}{2}}} \frac{\beta|\mu|}{\lambda^3} \int_0^\infty dx \frac{1}{x^{\frac{1}{2}}(x + \beta|\mu|)} \\ &\simeq \rho[1 + \frac{3}{2}(T - T_c)/T_c] - \frac{2\pi^{\frac{1}{2}}}{\zeta\left(\frac{3}{2}\right)} \rho(\beta|\mu|)^{\frac{1}{2}}, \end{aligned}$$

and thus

$$(\beta|\mu|)^{\frac{1}{2}} = 1.105(T - T_c)/T_c. \quad (\text{A2})$$

Electron Correlations, Magnetic Ordering, and Mott Insulation in Solids

PETER RICHMOND

Department of Physics, Queen Mary College, London, England

AND

GEOFFREY L. SEWELL*

University of Rochester, Rochester, New York

(Received 15 June 1967)

A model is constructed for the purpose of investigating electron correlations pertinent to magnetic ordering and Mott insulation in solids. The model consists of an assembly of interacting itinerant electrons in a periodic atomic lattice, such that the intra-atomic coupling between electrons is extremely strong. The correlations due to this latter coupling serve to prevent electrons of opposite spin from occupying the same atomic state, except in virtual transitions. Thus their net effect is to renormalize certain interactions and, also, to confine the state vectors of the entire system to a subspace, \mathfrak{H} , of the Hilbert space, \mathfrak{H}_0 , that is kinematically available to them. The observables are thus represented by operators on \mathfrak{H} , whose algebraic properties are different from those of the corresponding operators on \mathfrak{H}_0 . Thus, the correlations due to intra-atomic forces are imbedded in the theory in the form of the new algebra. In cases of one electron per atom, these correlations lead simply to both magnetic ordering and Mott insulation. In cases of nonintegral number of electrons per atom, they can lead to magnetic ordering, subject to specified conditions.

I. INTRODUCTION

In recent years a number of authors¹⁻³ have formulated theories of magnetic ordering in solids on the basis of itinerant electron models. Each of these

models represents an assembly of electrons in a periodic lattice, the electrons being confined to states formed from atomic orbitals for incomplete shells. They differ from the earlier model of Stoner⁴ in that they take account of the *intra*-atomic forces between electrons. These forces can be very strong⁵ in cases of some transition or rare-earth metals and their

* Permanent address: Queen Mary College, London, England.

¹ J. Hubbard, *Proc. Roy. Soc. (London)* **A276**, 238 (1963).

² J. Kanamori, *Progr. Theoret. Phys. (Kyoto)* **30**, 275 (1963); T. Izouyama and R. Kubo, *J. Appl. Phys.* **35**, 1074 (1964); M. C. Gutzwiller, *Phys. Rev.* **137**, 1726 (1956).

³ D. C. Mattis, *Phys. Rev.* **132**, 2521 (1963).

⁴ E. C. Stoner, *Proc. Roy. Soc. (London)* **A165**, 372 (1938).

⁵ N. F. Mott, *Advan. Phys.* **13**, 325 (1964).

more abstract, is the analogous concept for the ideal boson gas.

The lambda transition in liquid ^4He also features a divergence of C_P and K_T as one lowers T to the lambda T_λ , although in that case the divergences are logarithmic instead of the form $(T - T_\lambda)^{-1}$. Presumably these divergences also reflect the onset of ODLRO in ρ_1 but with a less rapid increase as $T \rightarrow T_\lambda + 0$ than for the ideal boson gas.⁸

ACKNOWLEDGMENT

A portion of the work of one of the authors (M. L.) was performed at the Department of Physics of the

⁸ An alternate suggestion for explaining the divergence of C_P and K_T for $T \rightarrow T_\lambda + 0$ in the case of liquid helium has been advanced by W. H. Keesom and A. P. Keesom, *Physica* **2**, 557 (1935). According to this suggestion, in the vicinity of T_λ the liquid is to be thought of as divided into small, weakly interacting domains which, because of thermal fluctuations, can have different temperatures. The domains with temperatures less than T_λ give the anomalously large contributions to the specific heat. The net effect is to cause the specific heat in the vicinity of T_λ to be a symmetric function of $T - T_\lambda$.

Hebrew University of Jerusalem, and he wishes to acknowledge their kind hospitality.

APPENDIX

We derive here (17) of the text, giving μ as a function of T just above T_c . Evaluating (15) for $T = T_c$, where $\mu = 0$, one obtains

$$\rho = \zeta\left(\frac{3}{2}\right) \left(\frac{mk_B T_c}{2\pi\hbar^2}\right)^{\frac{3}{2}}, \quad (\text{A1})$$

where $\zeta\left(\frac{3}{2}\right) = 2.612$. Just above T_c the value of μ is very small and thus we can approximate the right-hand side of (15) as

$$\begin{aligned} \rho &\simeq \frac{2}{\pi^{\frac{1}{2}}} \frac{1}{\lambda^3} \int_0^\infty dx \frac{x^{\frac{1}{2}}}{e^x - 1} - \frac{2}{\pi^{\frac{1}{2}}} \frac{\beta|\mu|}{\lambda^3} \int_0^\infty dx \frac{1}{x^{\frac{1}{2}}(x + \beta|\mu|)} \\ &\simeq \rho[1 + \frac{3}{2}(T - T_c)/T_c] - \frac{2\pi^{\frac{1}{2}}}{\zeta\left(\frac{3}{2}\right)} \rho(\beta|\mu|)^{\frac{1}{2}}, \end{aligned}$$

and thus

$$(\beta|\mu|)^{\frac{1}{2}} = 1.105(T - T_c)/T_c. \quad (\text{A2})$$

Electron Correlations, Magnetic Ordering, and Mott Insulation in Solids

PETER RICHMOND

Department of Physics, Queen Mary College, London, England

AND

GEOFFREY L. SEWELL*

University of Rochester, Rochester, New York

(Received 15 June 1967)

A model is constructed for the purpose of investigating electron correlations pertinent to magnetic ordering and Mott insulation in solids. The model consists of an assembly of interacting itinerant electrons in a periodic atomic lattice, such that the intra-atomic coupling between electrons is extremely strong. The correlations due to this latter coupling serve to prevent electrons of opposite spin from occupying the same atomic state, except in virtual transitions. Thus their net effect is to renormalize certain interactions and, also, to confine the state vectors of the entire system to a subspace, \mathfrak{H} , of the Hilbert space, \mathfrak{H}_0 , that is kinematically available to them. The observables are thus represented by operators on \mathfrak{H} , whose algebraic properties are different from those of the corresponding operators on \mathfrak{H}_0 . Thus, the correlations due to intra-atomic forces are imbedded in the theory in the form of the new algebra. In cases of one electron per atom, these correlations lead simply to both magnetic ordering and Mott insulation. In cases of nonintegral number of electrons per atom, they can lead to magnetic ordering, subject to specified conditions.

I. INTRODUCTION

In recent years a number of authors¹⁻³ have formulated theories of magnetic ordering in solids on the basis of itinerant electron models. Each of these

models represents an assembly of electrons in a periodic lattice, the electrons being confined to states formed from atomic orbitals for incomplete shells. They differ from the earlier model of Stoner⁴ in that they take account of the *intra*-atomic forces between electrons. These forces can be very strong⁵ in cases of some transition or rare-earth metals and their

* Permanent address: Queen Mary College, London, England.

¹ J. Hubbard, *Proc. Roy. Soc. (London)* **A276**, 238 (1963).

² J. Kanamori, *Progr. Theoret. Phys. (Kyoto)* **30**, 275 (1963); T. Izouyama and R. Kubo, *J. Appl. Phys.* **35**, 1074 (1964); M. C. Gutzwiller, *Phys. Rev.* **137**, 1726 (1956).

³ D. C. Mattis, *Phys. Rev.* **132**, 2521 (1963).

⁴ E. C. Stoner, *Proc. Roy. Soc. (London)* **A165**, 372 (1938).

⁵ N. F. Mott, *Advan. Phys.* **13**, 325 (1964).

compounds. If sufficiently strong, they lead naturally to nonzero atomic magnetic moments (Hund's rule). Once these moments are established, they can be aligned by cooperative effects with the result that the electronic system becomes magnetically ordered. Thus, the magnetic ordering arises as a result of both intra- and inter-electronic correlations.

A second kind of electronic ordering which can occur in cases of assemblies based on incomplete atomic shells is that corresponding to Mott⁶ insulation. This can occur only if the number of electrons per atom in the assembly is integral. This type of ordering also results from correlation effects that are not taken into account in the conventional band model. It is worth noting that the phenomenon of Mott insulation is often accompanied by magnetic ordering, e.g., in compounds of transition and rare-earth metals. Thus, it is of interest to understand whether, and how, these phenomena can arise from the same correlative effects.

The object of the present paper is to introduce a sharpened treatment of the electronic correlation problem, with reference to magnetic ordering and Mott insulation. This will be based on a model of interacting electrons in a periodic atomic lattice, for which the intra-atomic forces are extremely strong. We shall employ a simplification introduced by previous authors^{1,2} in that the electronic states will be considered to be formed from combinations of a set of equivalent nondegenerate atomic wavefunctions—even though the corresponding atomic functions in real magnetic systems are degenerate (*3d* or *5f*). Thus, in our model, c , the number of electrons per atom, is < 2 . Further, we may restrict our analysis to cases where $c < 1$, without loss of generality, since the cases $c > 1$ could equivalently be treated in terms of holes in the band, whose number per atom would then be $(2 - c) < 1$.

The essential new feature of our treatment is that it takes account, *ab initio*, of the fact that intra-atomic forces, if sufficiently strong, will prevent two electrons of opposite spin from occupying the same atomic state, except in virtual transitions which merely lead to certain renormalizations. Thus the intra-atomic forces serve to reduce the "phase space" available to the electronic system. In order to formulate this effect we first represent the states that are kinematically available to the system as vectors in a Hilbert space \mathfrak{H}_0 . Correspondingly, the observables are represented by a set of operators $\{Q_0\}$ on \mathfrak{H}_0 . We then take account of the fact that the intra-atomic forces exclude the

state vectors from a well-defined set in \mathfrak{H}_0 , which means that they confine the states to a subspace, \mathfrak{H} , of \mathfrak{H}_0 . Consequently, the observables are now represented by operators $\{Q\}$ on \mathfrak{H} . These are related to the corresponding primitive operators $\{Q_0\}$ by the formula

$$Q = PQ_0P,$$

where P is the projection operator for \mathfrak{H} . It is evident that the algebra of the set $\{Q\}$ is quite different from that of $\{Q_0\}$. This difference represents the changes in the properties of the system due to correlations engendered by intra-atomic forces. In other words, the effects of these forces are built into our formalism through the new algebra of the operators on the reduced Hilbert space \mathfrak{H} . Once the new algebra is formulated, it becomes a relatively simple matter to investigate the properties of the model. In particular, it leads to a simple theory of both Mott insulation and magnetic ordering for the case of a single electron per atom, and it also leads to a criterion for ferromagnetism in the case of nonintegral number of electrons per atom.

We shall set out the theory as follows. In Sec. II we shall construct our formalism for the model. Thus we start with a Hamiltonian, H_0 (on \mathfrak{H}_0), for an assembly of electrons in a band which interact via both intra- and inter-atomic forces with one another, and which are also subject to static forces that govern tunneling motion. We then take account of the strong intra-atomic forces by confining the state vectors to \mathfrak{H} and reformulating the model in terms of a simpler Hamiltonian operator H (on \mathfrak{H}). We also formulate the algebraic properties of new creation and annihilation from which all operators on \mathfrak{H} are generated. The theory that follows will be based on these operators on \mathfrak{H} .

In Sec. III we shall investigate the properties of the model for the case where c , the number of electrons per atom, is equal to unity. It will be shown, by an exact treatment of the model, that in this case the system is both a Mott insulator and a ferromagnet, or antiferromagnet, depending on the exchange forces.

In Sec. IV we shall obtain a criterion for ferromagnetism for the case where $c < 1$. This will be based on a treatment of appropriate quantum-mechanical Green functions in which a decoupling approximation is introduced. The criterion we obtain will be seen to be similar to that of Hubbard.¹ In Sec. V we shall summarize our conclusions.

II. THE MODEL

As stated above, the model will represent an assembly of interacting electrons. These particles are

⁶ N. F. Mott, *Phil. Mag.* 6, 287 (1961).

represented by a quantized wave operator

$$\psi_\lambda(x) = \sum_r \alpha_{r\lambda} \varphi_r(x),$$

where the φ_r 's are an orthonormal set of atomic wavefunctions (Wannier functions) localized at lattice sites r , λ ($= \pm 1$) labels the electronic spin, and $\alpha_{r\lambda}^*$, $\alpha_{r\lambda}$ are creation and annihilation operators on \mathfrak{H}_0 , satisfying the Fermion anticommutation rules

$$\left. \begin{aligned} \{\alpha_{r\lambda}^*, \alpha_{r'\lambda'}\}_+ &= 0 \\ \{\alpha_{r\lambda}^*, \alpha_{r'\lambda'}\}_+ &= \delta_{rr'} \delta_{\lambda\lambda'} \end{aligned} \right\} \quad (1)$$

The vacuum state vector $|o\rangle$ is defined by

$$\alpha_{r\lambda} |o\rangle = 0 \quad \text{for all } r, \lambda.$$

The set of vectors formed by application of the operators $\{\alpha_{r\lambda}^*\}$ and all their products to the vacuum-state vector, span the Hilbert space \mathfrak{H}_0 . The observables and state vectors of the system are represented by operators and vectors, respectively, in \mathfrak{H}_0 . In particular, the operator representing the number of electrons of spin λ at the site r is

$$v_{r\lambda} = \alpha_{r\lambda}^* \alpha_{r\lambda}, \quad (2)$$

while the operators representing the total electronic number and twice the electronic spin (in units where $\hbar = 1$) at that site are

$$\left. \begin{aligned} v_r &= \sum_\lambda v_{r\lambda} \\ \text{and} \\ \sigma_r &= (\sigma_r^{(x)}, \sigma_r^{(y)}, \sigma_r^{(z)}) \\ &= \sum_\lambda (\alpha_{r\lambda}^* \alpha_{r,-\lambda}, i\lambda \alpha_{r\lambda}^* \alpha_{r,-\lambda}, \lambda \alpha_{r\lambda}^* \alpha_{r\lambda}), \end{aligned} \right\} \quad (3)$$

respectively. The Hamiltonian for the system will be assumed to be (cf. Refs. 1, 2)

$$\begin{aligned} H_0 &= \sum'_{\Delta, r, \lambda} t_\Delta \alpha_{r+\Delta, \lambda}^* \alpha_{r\lambda} + \sum'_{\Delta, r} j_\Delta \sigma_{r+\Delta} \cdot \sigma_r \\ &\quad + \sum'_{\Delta, r} k_\Delta v_{r+\Delta} v_r + I_0 \sum_r v_{r1} v_{r,-1}, \end{aligned} \quad (4)$$

where the primes over the first three \sum 's indicate exclusion of terms with $\Delta = 0$. In this Hamiltonian, $I_0 (> 0)$ represents the interaction energy between two electrons on the same atom; j_Δ , k_Δ represent the strengths of exchange and spin-independent couplings between electrons on different sites, and t_Δ is the parameter governing the transfer of an electron from r to $r + \Delta$ by tunnel effect. It should be noted that the above form of H_0 depends on the neglect of interactions involving overlap of more than two atomic wavefunctions, (e.g., terms such as

$$\text{const} \times \alpha_{r_1 \lambda_1}^* \alpha_{r_2 \lambda_2}^* \alpha_{r_3 \lambda_3} \alpha_{r_4 \lambda_4}$$

with r_1, r_2, r_3, r_4 all unequal). This neglect is justified in cases where the atomic wavefunctions φ_r are highly localized, as in very narrow bands. It should also be noted that the parameters I_0, j, k, t are phenomenological quantities. In relating their values to properties of real solids, one should realize that these parameters contain contributions due to indirect interactions involving other bands; for example, the interaction between "magnetic" electrons in a metal is screened by conduction electrons in a higher conduction band. Further, the tunneling parameter t_Δ may be strongly dependent on the coupling between electrons and lattice vibrations⁷ as this can reduce t_Δ to much less than its value for a rigid lattice.

Our treatment of the model will be based on the assumption that is valid for very narrow bands⁵—namely, that the intra-atomic interaction parameter, I_0 , is much greater in absolute magnitude than the interatomic couplings $t_\Delta, j_\Delta, k_\Delta$. Now, by formula (4), I_0 represents the energy required to bring two electrons into the same atomic (spatial) state. Consequently, for sufficiently large I_0 ($\gg |t_\Delta|, |j_\Delta|, |k_\Delta|$) such doubly occupied atomic states cannot occur in the low-lying states of the entire system, except in virtual transitions. It is well known that, in general, such transitions lead to renormalizations of the interactions in a system. In the present context, this means that they lead to renormalizations⁸ of the parameters $t_\Delta, j_\Delta, k_\Delta$. It is readily seen, for example, that the renormalization of j_Δ , by processes involving virtual transitions to doubly occupied atomic states, is exactly equivalent to the introduction of Anderson's⁹ superexchange coupling.

Thus, as transitions to doubly occupied atomic states occur only in virtual processes, we may take account of them by renormalizing the parameters $t_\Delta, j_\Delta, k_\Delta$, and thereafter excluding such processes from the theory. This exclusion is evidently equivalent to confining the state vectors of the system to \mathfrak{H} , the subspace of \mathfrak{H}_0 spanned by those vectors $|\rangle$ for which

$$v_{r1} v_{r,-1} |\rangle = 0, \quad \text{for all } r.$$

Thus, the projection operator for \mathfrak{H} is

$$P = \prod_r (I - v_{r1} v_{r,-1}). \quad (5)$$

Our program, then, is to reformulate the model so as to represent the observables and states of the system

⁷ G. L. Sewell, *Phil. Mag.* **3**, 1361 (1958).

⁸ One may derive formulas for T, J, K in terms of t, j, k by taking account of virtual transitions from \mathfrak{H} to $\mathfrak{H}_0 - \mathfrak{H}$ by standard-field theoretical methods. [cf. W. Heitler, *Quantum Theory of Radiation* (Clarendon Press, Oxford, England, 1954)].

⁹ P. W. Anderson, *Solid State Phys.* **14**, 99 (1963).

by operators and vectors, respectively, in \mathfrak{H} . For this purpose we note that, if Q_0 is an operator on \mathfrak{H}_0 , then its component in \mathfrak{H} is PQ_0P . Thus we define creation and annihilation operators in \mathfrak{H} as the components of $\alpha_{r\lambda}^*$, $\alpha_{r\lambda}$ given by

$$a_{r\lambda}^* = P\alpha_{r\lambda}^*P; \quad a_{r\lambda} = P\alpha_{r\lambda}P; \quad (6)$$

and, likewise, we define number and spin operators on \mathfrak{H} by

$$n_{r\lambda} = P\nu_{r\lambda}P; \quad n_r = P\nu_rP; \quad s_r = P\sigma_rP.$$

It follows from these equations, together with (1)–(3), (5), and (6), that

$$\left. \begin{aligned} n_{r\lambda} &= a_{r\lambda}^* a_{r\lambda}, \quad n_r = \sum_{\lambda} a_{r\lambda}^* a_{r\lambda} \\ s_r &= (s_r^{(x)}, s_r^{(y)}, s_r^{(z)}) \\ &= \sum_{\lambda} (a_{r\lambda}^* a_{r,-\lambda}, i\lambda a_{r\lambda}^* a_{r,-\lambda}, \lambda a_{r\lambda}^* a_{r\lambda}). \end{aligned} \right\} \quad (7)$$

The operators a^* , a do not satisfy the normal Fermion anticommutation rules. In fact, it follows from Eqs. (1), (5)–(7), that the anticommutative algebra of these operators is given by

$$\left. \begin{aligned} \{a_{r\lambda}^*, a_{r'\lambda'}^*\}_+ &= 0, \\ a_{r\lambda}^* a_{r'\lambda'}^* &= 0; \end{aligned} \right\} \quad (8)$$

and

$$\{a_{r\lambda}^*, a_{r'\lambda'}\}_+ = (P - n_{r,-\lambda})\delta_{r'r'}\delta_{\lambda\lambda'} + \frac{1}{2}(s_r^{(x)} - i\lambda s_r^{(y)})\delta_{r'r'}\delta_{\lambda,-\lambda'}.$$

As we shall henceforth be concerned only with operators on \mathfrak{H} , and as P is the unit operator in that space, we may replace P by unity on the rhs of this last equation. Thus

$$\{a_{r\lambda}^*, a_{r'\lambda'}\}_+ = (1 - n_{r,-\lambda})\delta_{r'r'}\delta_{\lambda\lambda'} + \frac{1}{2}(s_r^{(x)} - i\lambda s_r^{(y)})\delta_{r'r'}\delta_{\lambda,-\lambda'}. \quad (9)$$

Further important algebraic relations that follow from (7)–(9) are

$$\left. \begin{aligned} [a_{r\lambda}^*, n_{r'\lambda'}]_- &= -a_{r\lambda}^* \delta_{r'r'} \delta_{\lambda\lambda'} \\ [a_{r\lambda}^*, s_r]_- &= -(a_{r,-\lambda}^*, i\lambda a_{r,-\lambda}^*, \lambda a_{r\lambda}^*); \end{aligned} \right\} \quad (10)$$

and, also, the usual spin commutation rules

$$\left. \begin{aligned} [s_r^{(x)}, s_r^{(y)}]_- &= i s_r^{(z)} \delta_{r'r'}; \quad [s_r^{(y)}, s_r^{(z)}]_- = i s_r^{(x)} \delta_{r'r'} \\ [s_r^{(z)}, s_r^{(x)}]_- &= i s_r^{(y)} \delta_{r'r'}. \end{aligned} \right\} \quad (11)$$

In order to formulate the effective Hamiltonian for the model as an operator on \mathfrak{H} , we first note that the

component of H_0 in that subspace is

$$PH_0P = \tilde{H}, \quad \text{say,}$$

i.e., by Eqs. (1)–(7),

$$\tilde{H} = \sum'_{\Delta,r,\lambda} t_{\Delta} a_{r+\Delta,\lambda}^* a_{r\lambda} + \sum'_{\Delta,r} j_{\Delta} s_{r+\Delta} \cdot s_r + \sum'_{\Delta,r} k_{\Delta} n_{r+\Delta} n_r.$$

This is not the effective Hamiltonian for the model, as it takes no account of virtual transitions from \mathfrak{H} to the complimentary space $\mathfrak{H}_0 - \mathfrak{H}$ that lead to renormalizations of the interaction parameters t_{Δ} , j_{Δ} , and k_{Δ} . The effective Hamiltonian, H , is simply obtained by replacing those parameters by the corresponding renormalized quantities T_{Δ} , J_{Δ} , and K_{Δ} in the above formula for \tilde{H} . Thus

$$H = H_T + H_J + H_K \quad (12)$$

with

$$\begin{aligned} H_T &= \sum'_{\Delta,r,\lambda} T_{\Delta} a_{r+\Delta,\lambda}^* a_{r\lambda}; \quad H_J = \sum'_{\Delta,r} J_{\Delta} s_{r+\Delta} \cdot s_r; \\ H_K &= \sum'_{\Delta,r} K_{\Delta} n_{r+\Delta} n_r. \end{aligned} \quad (13)$$

The Hamiltonian H , like H_0 , is appropriate when the atomic wavefunctions φ_r are so highly localized that effective interactions arising from overlap of more than two of them may be neglected. In fact, Eqs. (12) and (13) represent the most general form for a Hamiltonian governing a system whose states are confined to \mathfrak{H} , and for which overlap between three or more different atomic states is negligible. We shall henceforth treat T_{Δ} , J_{Δ} , and K_{Δ} as the basic parameters of the model, rather than express⁸ these quantities in terms of t_{Δ} , j_{Δ} , k_{Δ} —recall that these latter parameters are themselves highly complicated quantities, whose values involve certain renormalizations. Thus we would regard T_{Δ} , J_{Δ} , and K_{Δ} as phenomenological parameters whose values might be obtained by matching the properties of our model, as represented by Eqs. (7)–(13), to experimental data.

The model is now entirely defined by Eqs. (7)–(13). It is evident that correlations due to the strong intra-atomic forces that confine the state vectors to \mathfrak{H} , are built into the formulation of the model through the new anticommutation rules (8) and (9).

III. THE CASE $c = 1$: MAGNETIC ORDERING AND MOTT INSULATION

As no atomic state can be doubly occupied, it follows that in the case where the total numbers of electrons and atomic sites of the model are equal, i.e., where $c = 1$, the system is confined to states where each site is occupied by precisely one electron. In other words, the state vectors of the system are confined to a subspace of \mathfrak{H} , namely $\mathfrak{H}^{(1)}$, spanned by the set of vectors $|\rangle$ for which $n_r |\rangle = |\rangle$, for

all r . Hence, by (7), (8), (10), and (13),

$$\begin{aligned} H_T | \rangle &= \sum'_{\Delta, r, \lambda} T_{\Delta} a_{r+\Delta, \lambda}^* a_{r, \lambda} n_{r+\Delta} | \rangle \\ &= \sum'_{\Delta, r, \lambda} T_{\Delta} a_{r+\Delta, \lambda}^* n_{r+\Delta} a_{r, \lambda} | \rangle \\ &= \sum'_{\Delta, r, \lambda, \lambda'} T_{\Delta} a_{r+\Delta, \lambda}^* a_{r+\Delta, \lambda'}^* a_{r+\Delta, \lambda} a_{r, \lambda} | \rangle = 0; \end{aligned}$$

and

$$H_K | \rangle = \sum'_{\Delta, r} K_{\Delta} | \rangle.$$

Thus, by (12), H reduces to H_J , apart from an additive constant, when acting on $\mathfrak{H}^{(1)}$. Further, it is evident from Eqs. (11) and (13) that H_J is simply the Heisenberg spin Hamiltonian which represents an insulating ferromagnet or antiferromagnet according to the value of J_{Δ} .

Consequently, in the case $c = 1$, the model is both a Mott-type insulator and a Heisenberg ferromagnet (or antiferromagnet). The reason for the insulation is simply that, in the present model, electrons cannot move onto sites that are already occupied, which they all are in the case $c = 1$.

IV. THE CASE $c < 1$: CRITERION FOR FERROMAGNETISM

We now consider the properties of the model for the case $c < 1$. For definiteness we restrict ourselves to consideration of ferromagnetic and paramagnetic phases where the states of the system possess the translational symmetry of the atomic lattice—this automatically excludes spin-density waves¹⁰ and other antiferromagnetic-type configurations from the theory. Our aim, then, is to obtain a criterion for which the ferromagnetic, rather than the paramagnetic, state, is stable.

In view of the translational symmetry of the model, it is convenient to express its properties in terms of extended rather than localized wave operators. Thus we define

$$\left. \begin{aligned} a_{k\lambda} &= N^{-\frac{1}{2}} \sum_r a_{r\lambda} e^{ik \cdot r}, & a_{k\lambda}^* &= N^{-\frac{1}{2}} \sum_r a_{r\lambda}^* e^{-ik \cdot r}, \\ \text{i.e.,} & & & \\ a_{r\lambda} &= N^{-\frac{1}{2}} \sum_k a_{k\lambda} e^{-ik \cdot r}, & a_{r\lambda}^* &= N^{-\frac{1}{2}} \sum_k a_{k\lambda}^* e^{ik \cdot r}, \end{aligned} \right\} \quad (14)$$

where N is the total number of lattice sites and k represents a wave vector in the first Brillouin zone. We also define number and spin-density wave operators by the equations

$$\left. \begin{aligned} n_{k\lambda} &= N^{-1} \sum_r n_{r\lambda} e^{ik \cdot r}, & n_k &= N^{-1} \sum_r n_r e^{ik \cdot r}, \\ \text{and} & & & \\ s_k &= N^{-1} \sum_r s_r e^{ik \cdot r}. \end{aligned} \right\} \quad (15)$$

It follows from (8)–(10), (14), and (15) that the extended wave operators satisfy the algebraic relations

$$\left. \begin{aligned} \{a_{k\lambda}^*, a_{k'\lambda'}\}_+ &= (\delta_{kk'} - n_{k-k', -\lambda}) \delta_{\lambda\lambda'} \\ &\quad + \frac{1}{2}(s_{k-k'}^{(x)} - i\lambda s_{k-k'}^{(y)}) \delta_{\lambda, -\lambda'}, \\ \{a_{k\lambda}^*, a_{k'\lambda'}\}_+ &= 0, \\ [a_{k\lambda}^*, n_{k'\lambda'}]_- &= -a_{k-k', \lambda}^* \delta_{\lambda\lambda'}, \\ [a_{k\lambda}^*, s_{k'}]_- &= -(a_{k-k', -\lambda}^*, i\lambda a_{k-k', -\lambda}^*, \lambda a_{k-k', \lambda}^*). \end{aligned} \right\} \quad (16)$$

Further, by (13)–(15), the Hamiltonian components H_T , H_J , and H_K may be expressed in terms of the extended wave operators by the equations

$$\begin{aligned} H_T &= \sum_{k\lambda} T(k) a_{k\lambda}^* a_{k\lambda}; & H_J &= N \sum_k J(k) s_k \cdot s_{-k}; \\ H_K &= N \sum_k K(k) n_k n_{-k}, \end{aligned} \quad (17)$$

with

$$\begin{aligned} T(k) &= \sum'_{\Delta} T_{\Delta} e^{ik \cdot \Delta}; & J(k) &= \sum'_{\Delta} J_{\Delta} e^{ik \cdot \Delta}; \\ K(k) &= \sum'_{\Delta} K_{\Delta} e^{ik \cdot \Delta} \end{aligned} \quad (18)$$

The operator representing the total number of electrons of spin λ is $n_{\lambda} = \sum_r n_{r\lambda}$, i.e., by (7), (14), and (15),

$$n_{\lambda} = \sum_k a_{k\lambda}^* a_{k\lambda} = N(n_{k\lambda})_{k=0}. \quad (19)$$

The equilibrium statistical operator for the system will be taken to be

$$\rho = \Omega \exp \left\{ -\beta \left(H - \sum_{\lambda} \mu_{\lambda} n_{\lambda} \right) \right\}, \quad (20)$$

where Ω is a normalizing constant, $\beta = (\kappa_B T)^{-1}$ and μ_{λ} is the chemical potential governing the number of particles of spin λ . The thermal average of an observable Q will then be

$$\bar{Q} \equiv \langle Q \rangle \equiv \text{Tr}(\rho Q). \quad (21)$$

We shall assume, as we may do without loss of generality, that the total electronic spin is $N\xi$ directed along Oz . As the mean number of electrons per atom is c , it follows from (15) and (19) and the assumed translational symmetry that

$$\left. \begin{aligned} \langle s_k \rangle &= (0, 0, \xi) \delta_{k0} \\ \text{and} & \\ \langle n_{k\lambda} \rangle &= \frac{1}{2}(c + \lambda \xi) \delta_{k0} \equiv c_{\lambda} \delta_{k0}. \end{aligned} \right\} \quad (22)$$

The investigation of the properties of the model will be centered on the Green functions

$$\left. \begin{aligned} G_{k\lambda}(t) &= \langle a_{k\lambda}^*(t) a_{k\lambda} \rangle \\ \text{and} & \\ G_{k\lambda}(t) &= \langle \langle a_{k\lambda}^* | a_{k\lambda} \rangle \rangle_t, \end{aligned} \right\} \quad (23)$$

where $A(t)$ is the Heisenberg operator $e^{iHt} A e^{-iHt}$,

¹⁰ A. W. Overhauser, Phys. Rev. **128**, 1437 (1962).

in units where $\hbar = 1$, and

$$\langle\langle A | B \rangle\rangle_t \equiv \begin{cases} i\langle\langle A(t), B \rangle\rangle_+, & \text{for } t \geq 0 \\ 0, & \text{for } t < 0. \end{cases} \quad (24)$$

Following the general theory of Green functions (cf. Ref. 10, Chap. 1, Sec. 4), we express $G_{k\lambda}(t)$ in terms of the Fourier transform of G^+ . Thus

$$G_{k\lambda}(t) = 2 \int_{-\infty}^{\infty} d\omega e^{i\omega t} \text{Im } G_{k\lambda}^+(\omega) [1 + \exp \beta(\omega - \mu_\lambda)]^{-1} \quad (25)$$

with

$$G_{k\lambda}^+(\omega) = (2\pi)^{-1} \int_{-\infty}^{\infty} dt e^{-i\omega t} G_{k\lambda}(t). \quad (26)$$

We shall henceforth restrict the theory to the absolute zero of temperature ($\beta \rightarrow \infty$), as this will suffice for our purpose of obtaining a criterion for ferromagnetism. Thus we now rewrite (25) as

$$G_{k\lambda}(t) = 2 \int_{-\infty}^{\mu_\lambda} d\omega e^{i\omega t} G_{k\lambda}^+(\omega). \quad (25')$$

The energy of the system may be expressed in terms of the function $G^+(\omega)$, for it follows from (7)–(10) and (12)–(14) that

$$\begin{aligned} & \left(-i \frac{\partial}{\partial t} \sum_{k\lambda} a_{k\lambda}^*(t) a_{k\lambda} \right)_{t=0} \\ & \equiv \left(-i \frac{\partial}{\partial t} \sum_{r\lambda} a_{r\lambda}^*(t) a_{r\lambda} \right)_{t=0} \equiv \sum_{r\lambda} [H, a_{r\lambda}^*]_{-} a_{r\lambda} \\ & \equiv H_T + 2H_J + 2H_K. \end{aligned}$$

Hence, by (12), (17), and (23), the total energy per atomic site is

$$\begin{aligned} E &= N^{-1} \langle H_T + H_J + H_K \rangle \\ & \equiv (2N)^{-1} \sum_{k\lambda} [T(k)G_{k\lambda}(t) - i\dot{G}_{k\lambda}(t)]_{t=0}, \end{aligned}$$

i.e., by (25')

$$E = N^{-1} \sum_{k\lambda} \int_{-\infty}^{\mu_\lambda} d\omega (\omega + T(k)) \text{Im } G_{k\lambda}^+(\omega). \quad (27)$$

Similarly, it follows from (19), (23), and (25') that the total number of electrons of spin λ per atom is

$$\frac{2}{N} \sum_k \int_{-\infty}^{\mu_\lambda} d\omega \text{Im } G_{k\lambda}^+(\omega).$$

Hence it follows from (19) and (22) that

$$c_\lambda = \frac{1}{2}(c + \lambda\xi) = \frac{2}{N} \sum_k \int_{-\infty}^{\mu_\lambda} d\omega \text{Im } G_{k\lambda}^+(\omega). \quad (28)$$

This equation serves to relate the chemical potential to c and ξ .

It now remains for us to evaluate G^+ and thence to express E as a function of ξ by means of Eqs. (27)

and (28). For this purpose, we formulate the equation of motion for $G^+(t)$, which may be written, in view of (23) and (24), as

$$\begin{aligned} -i(\partial/\partial t)G_{k\lambda}^+(t) &= \langle\langle [H, a_{k\lambda}^*]_{-} | a_{k\lambda} \rangle\rangle_t \\ & \quad + \langle\langle a_{k\lambda}^*, a_{\lambda k} \rangle\rangle_+ \delta(t). \end{aligned}$$

Hence, by (12), (16), and (17),

$$\begin{aligned} -i \frac{\partial}{\partial t} G_{k\lambda}^+(t) &= \sum_{k'} T(k') \langle\langle a_{k\lambda}^* (\delta_{kk'} - n_{k-k',-\lambda}) | a_{k\lambda} \rangle\rangle_t \\ & \quad + \frac{1}{2} \sum_{k'} T(k') \langle\langle a_{k,-\lambda}^* (s_{k-k'}^{(x)} - i\lambda s_{k-k'}^{(y)}) | a_{k\lambda} \rangle\rangle_t \\ & \quad + 2\lambda \sum_{k'} J(k') \langle\langle s_{-k'}^{(z)} a_{k-k',\lambda}^* | a_{k\lambda} \rangle\rangle_t \\ & \quad + 2 \sum_{k'} J(k') \langle\langle (s_{-k'}^{(x)} + i\lambda s_{-k'}^{(y)}) a_{k-k',\lambda}^* | a_{k\lambda} \rangle\rangle_t \\ & \quad + 2 \sum_{k'} K(k') \langle\langle n_{-k'} a_{k-k',\lambda}^* | a_{k\lambda} \rangle\rangle_t + \langle 1 - n_{k,-\lambda} \rangle \delta(t). \end{aligned}$$

We now employ a familiar decoupling procedure, replacing the field variables that multiply a^* on the rhs of this equation by their average values. Hence, by (22) and (23), the equation simplifies to the form

$$\begin{aligned} -i(\partial/\partial t)G_{k\lambda}^+(t) &= [T(k)(1 - c_{-\lambda}) + 2\lambda\xi J(0) + 2cK(0)]G_{k\lambda}^+(t). \end{aligned}$$

The Fourier transform of this equation is simply

$$\begin{aligned} G_{k\lambda}^+(\omega) &= \frac{1}{2\pi} [(1 - c_{-\lambda})/\omega - T(k)(1 - c_{-\lambda}) \\ & \quad - 2\lambda\xi J(0) - 2cK(0) + i\delta], \quad (29) \end{aligned}$$

where δ is an infinitesimal real, positive, constant arising from the requirement that $G^+(t)$ is zero for $t < 0$. It follows from this equation that, as the single-particle energies are given by the poles of $G^+(\omega)$, these energies are simply

$$(1 - c_{-\lambda})T(k) + 2\lambda\xi J(0) + 2cK(0).$$

In this expression, the factor $(1 - c_{-\lambda})$ represents a "band narrowing" due to restrictions imposed on the electronic motion by the exclusion of doubly occupied atomic states, while the terms $2\xi J(0)$ and $2cK(0)$ represent additional contributions to the energy of a particle due to its exchange and spin-independent couplings, respectively, to the rest of the electrons.

It follows from (27)–(29) that E and c_λ are related to the chemical potentials by the formulas

$$\begin{aligned} E &= N^{-1} \sum_{k\lambda} \int_{-\infty}^{\mu_\lambda} d\omega [(1 - \frac{1}{2}c_{-\lambda})\omega - \lambda\xi J(0) - cK(0)] \\ & \quad \times \delta\{\omega - 2\lambda\xi J(0) - 2cK(0) - (1 - c_{-\lambda})T(k)\} \end{aligned}$$

and

$$c_\lambda = N^{-1} \sum_k (1 - c_{-\lambda}) \int_{-\infty}^{\mu_\lambda} d\omega \times \delta\{\omega - 2\lambda\xi J(0) - 2cK(0) - (1 - c_{-\lambda})T(k)\}.$$

By changing the variable of integration, these equations may be rewritten as

$$\left. \begin{aligned} E &= \sum_\lambda (1 - c_{-\lambda}) \int_{-\infty}^{\mu'_\lambda} d\omega P(\omega) \\ &\quad \times \left[(1 - \frac{1}{2}c_{-\lambda})\omega + \lambda\xi J(0) + cK(0) \right] \\ \text{and} \\ c_\lambda &= (1 - c_{-\lambda}) \int_{-\infty}^{\mu'_\lambda} d\omega P(\omega), \end{aligned} \right\} (30)$$

where

$$P(\omega) = N^{-1} \sum_k \delta(\omega - T(k)) \quad (31)$$

and the explicit form of the relationship between μ'_λ and μ_λ need not concern us further.

Now $P(\omega)$, as defined by (31), is the normalized density of states for a system of noninteracting particles, with Hamiltonian $\sum T_{\Delta} \alpha_{\tau+\Delta, \lambda}^* \alpha_{\tau, \lambda}$, which we shall regard as a reference system. Equations (30) and (31) enable us to express the properties of the model of interest in terms of those of the reference system. Thus we note that the energy and particle number, per atom, corresponding to a spin direction governed by chemical potential μ in the reference system are

$$\left. \begin{aligned} \epsilon &= \int_{-\infty}^{\mu} \omega P(\omega) d\omega, \\ \text{and} \\ y &= \int_{-\infty}^{\mu} P(\omega) d\omega. \end{aligned} \right\} (32)$$

These equations define a functional dependence of μ and ϵ on y , i.e.,

$$\left. \begin{aligned} \epsilon &= \epsilon(y) \\ \text{and} \\ \mu &= \mu(y). \end{aligned} \right\} (33)$$

It follows from (32) and (33) that

$$d\epsilon/dy = \mu(y); \quad d^2\epsilon/dy^2 = [P\{\mu(y)\}]^{-1}. \quad (34)$$

On comparing (30) with (32) it follows that, corresponding to (33), the solution of (30) is

$$\mu'_\lambda = \mu\{c_\lambda/(1 - c_{-\lambda})\}$$

and

$$E = \sum_\lambda \left\{ (1 - c_{-\lambda}) \left(1 - \frac{1}{2}c_{-\lambda} \right) \epsilon\{c_\lambda/(1 - c_{-\lambda})\} + (1 - c_{-\lambda})(\lambda\xi J(0) + cK(0)) \right\},$$

i.e., by (22),

$$\begin{aligned} E &= \xi^2 J(0) + (1 - \frac{1}{2}c - \frac{1}{2}\xi) \\ &\quad \times (1 - \frac{1}{4}c - \frac{1}{4}\xi) \\ &\quad \times \epsilon\{(c - \xi)/(2 - c - \xi)\} \\ &\quad + (1 - \frac{1}{2}c + \frac{1}{2}\xi)(1 - \frac{1}{4}c + \frac{1}{4}\xi) \\ &\quad \times \epsilon\{(c + \xi)/(2 - c + \xi)\}, \end{aligned} \quad (35)$$

together with an additive term independent of ξ .

The system will be ferromagnetic if E takes its minimum value for nonzero ξ . A sufficient condition for ferromagnetism is therefore that

$$(d^2E/d\xi^2)_{\xi=0} < 0.$$

By Eqs. (34) and (35), this condition is equivalent to

$$2J(0) + \frac{(1 - c)^2(1 - \frac{1}{4}c)}{2P(\mu_0)(1 - \frac{1}{2}c)^3} + \frac{1}{2}\epsilon_0 + \frac{(1 - c)}{(2 - c)}\mu_0 < 0 \quad (36)$$

where

$$\mu_0 = \mu\{c/(2 - c)\}; \quad \epsilon_0 = \epsilon\{c/(2 - c)\}. \quad (37)$$

Let us now apply this criterion to two cases.

Case (a)

Consider the case where $J(0) \neq 0$ and $P(\omega)$ is uniform, i.e.,

$$\left. \begin{aligned} P(\omega) &= (2T)^{-1}, \quad \text{for } -T \leq \omega \leq T \\ &= 0 \quad \text{elsewhere.} \end{aligned} \right\} (38)$$

It follows from (32) and (38) that, in this case,

$$y = 1/2T(\mu + T), \quad \epsilon = 1/4T(\mu^2 - T^2)$$

and, consequently,

$$\mu(y) = T(2y - 1), \quad \epsilon(y) = Ty(y - 1).$$

Using these formulas for $\mu(y)$, $\epsilon(y)$, Eqs. (37) take the forms

$$\left. \begin{aligned} \mu_0 &= T[(3c - 2)/(2 - c)], \\ \epsilon_0 &= [2Tc(c - 1)]/(2 - c)^2. \end{aligned} \right\} (39)$$

On inserting the values for P , μ_0 , and ϵ_0 given by Eqs. (38) and (39) into the condition (36), we see that this condition reduces to

$$J(0) + [T(1 - c)^2]/[4(1 - \frac{1}{2}c)^2] < 0. \quad (40)$$

This signifies that, for ferromagnetism, the exchange parameter $J(0)$ must be negative and large enough in absolute magnitude to exceed the tunneling parameter, reduced by a factor proportional to $(1 - c)^2$ as a result of the exclusion of doubly occupied atomic states.

Case (b)

We now consider the case where $J(0) = 0$. The purpose of this is to investigate whether, even in the absence of interatomic exchange forces, the combined effects due to intra-atomic correlations and electronic

itineracy can lead to magnetic ordering. It will be shown that, even with $J(0) = 0$, the condition (36) can be fulfilled for certain forms of $P(\omega)$. In particular, we shall show that the condition can be fulfilled if the form of $P(\omega)$ has two suitably separated sharp peaks. This case is of physical interest since electronic densities of states in transition metals¹¹ exhibit similar energy dependences. Our conclusions about this case will be in general agreement with Hubbard's.¹

We assume, then, that P take the form

$$P(\omega) = \frac{1}{2\Delta} \left[p\left(\frac{\omega - T_0}{\Delta}\right) + p\left(\frac{\omega + T_0}{\Delta}\right) \right] \quad (41)$$

where $p(\theta)$ takes its maximum value at $\theta = 0$ and possesses the following properties:

$$p(0) = 1; \quad (42)$$

$$\int_{-\infty}^{\infty} d\theta p(\theta) = 1; \quad \int_{-\infty}^{\infty} d\theta \theta p(\theta) = 0; \quad (43)$$

$$\int_{-\infty}^{\infty} d\theta \theta^2 p(\theta) = 1;$$

and

$$p\left(\frac{\omega \pm T_0}{\Delta}\right) = 0 \quad \text{for } |\omega| > T, \quad (44)$$

where $(-T, T)$ is the energy range of the band for the reference system. We also assume that

$$T_0 \gg \Delta. \quad (45)$$

It should be pointed out that although we have assumed that the maximum and the dispersion of $p(\theta)$ are both equal to unity, it would actually suffice for our purposes if they were both of the order of unity, as distinct from T_0/Δ , say. In fact, throughout the following analysis, it will suffice to consider the orders of magnitude of the terms arising in the condition (36), bearing in mind the inequality (45).

It follows from (41)–(45) that we have defined $P(\omega)$ so that it consists of two sharp peaks, each of height Δ^{-1} and width Δ , centered at $\omega = \pm T_0$. Further, as it follows from this form of P and Eqs. (32) and (33), that $\mu(y)$, the Fermi energy corresponding to concentration y of the reference system, will lie within Δ of one of the centers, $\pm T_0$, for all but very small ranges of y . Specifically, it will lie in the range $(-T_0 - \Delta, -T_0 + \Delta)$ for

$$y < \frac{1}{2},$$

$$y, \quad \frac{1}{2} - y \gtrsim \Delta/T_0.$$

Hence, by (37) and (45), μ_0 will lie in $(-T_0 - \Delta, -T_0 + \Delta)$ for

$$c < \frac{2}{3}$$

and

$$c, \quad \frac{2}{3} - c > \Delta/T_0.$$

We shall assume henceforth that these latter conditions on c are satisfied and, thus, that

$$\mu = -T_0 + 0(\Delta). \quad (46)$$

Likewise, it follows from (32), (33), and (37), that, under the above conditions

$$\epsilon = -T_0 + 0(\Delta). \quad (47)$$

Let us now examine the terms involved in the condition (36). The first term is zero, as we are here considering cases where $J(0) = 0$. By Eqs. (41)–(45), the second term of (36) is

$$\simeq \frac{(1-c)^2(1-\frac{1}{4}c)}{(1-\frac{1}{2}c)^3} \Delta, \quad (48)$$

while by (46) and (47), the third and fourth terms of (36) add up to

$$\simeq -(4-3c)/(2-c)T_0. \quad (49)$$

Hence, as $T_0 \gg \Delta$ —it is seen that the sum of the terms (48) and (49) is negative, which means that the condition (36) is fulfilled.

V. CONCLUSION

We have constructed a model of interacting electrons in an atomic lattice that is designed for cases of strong intra-atomic coupling. This coupling serves to confine the electron-state vectors to a subspace, \mathfrak{H} , of the Hilbert space, \mathfrak{H}_0 , kinematically available to them. Thus, the correlations due to intra-atomic forces are simply represented in terms of the algebra of operators on \mathfrak{H} .

These correlations lead directly to both magnetic ordering and Mott insulation for the case where c , the number of electrons per atom, is equal to unity. In cases where $c < 1$ (or $1 < c < 2$), they tend to favor magnetic ordering, and, for certain forms of the energy dependence of the electronic density of states, they can lead to ferromagnetism without the aid of exchange forces. One such form [case (b) of Sec. IV] is qualitatively similar to those of "magnetic" electrons in transition metals.

ACKNOWLEDGMENTS

One of us, (P. R.), wishes to express his gratitude to Science Research Council for financial assistance. The other, (G. L. S.), would like to thank Professor E. W. Montroll for the hospitality he received during a pleasant and stimulating visit to the University of Rochester.

¹¹ V. L. Bonch-Bruевич and S. V. Tyablikov, *The Green Function Method in Statistical Mechanics* (North-Holland Publishing Co., Amsterdam, 1962), Chap. 1.

Foundations for Quantum Statistics*

ROBIN GILES

Department of Mathematics, Queen's University, Kingston, Ontario

(Received 9 August 1967)

A derivation of Fermi and Bose statistics is given, based on the general structure of quantum mechanics, together with a simple axiom of direct physical significance. The axiom concerns an operation, denoted by \circ , forming the *union* of two states; $\Psi \circ \Phi$ denotes the state of a compound system whose parts are in the states Ψ and Φ . Let Ψ and Φ denote 1-particle states and assume: (a) $\Psi \circ \Phi$ exists whenever $\Psi \perp \Phi$; (b) $\Psi \circ \Phi = \Phi \circ \Psi$; (c) the transition probability between $\Psi \circ \Phi$ and $\Psi' \circ \Phi'$ is zero if $\Psi \perp \Psi'$ and $\Phi \perp \Phi'$, and (d) the product of the transition probabilities from Ψ to Ψ' and from Φ to Φ' if $\Psi \perp \Phi'$ and $\Phi \perp \Psi'$. It is then shown that, at least in so far as 2-particle states are concerned, the particles obey either Fermi or Bose statistics.

1. INTRODUCTION

In a recent review¹ Dresden has drawn attention to the fact that the arguments generally put forward as establishing the existence of just two kinds of quantum statistics (Fermi-Dirac and Bose-Einstein) are usually invalid, or at best depend on assumptions which are not explicitly stated. He examines various arguments in detail and finds no compelling reason for the exclusion of "parastatistics," but remarks (p. 383): "It is conceivable that a more detailed and rigorous investigation of the mathematical character of the theory would show the impossibility of other statistics. It is also conceivable that an additional physical principle (preferably one which has a direct and transparent physical interpretation) leads to the desired exclusion."

The present paper is intended as a contribution in this direction. We give a derivation of Fermi and Bose statistics from a minimum of assumptions. First, we assume the structure of general quantum mechanics with superselection rules. This is entirely conventional, and is given together with some mathematical preliminaries in Sec. 2.

In Sec. 3 the characteristic feature of the present approach is introduced. As Dresden remarks (Ref. 1, p. 395), it is a regrettable fact that most treatments of quantum statistics start with a description in which the particles are regarded as identifiable, while the essential indistinguishability is introduced only later: for instance, if $|\psi_1(x)\rangle$ and $|\psi_2(x)\rangle$ are two 1-particle wavefunctions, then $|\psi(x_1, x_2)\rangle = |\psi_1(x_1)\psi_2(x_2)\rangle$ represents a state of a 2-particle system in which the first particle is in the state ψ_1 and the second in the state ψ_2 . It is assumed that every 2-particle wave-

function is a linear combination of "product wavefunctions" of this type. Only then are arguments presented to show that, for any one kind of particle, either only symmetric or only antisymmetric linear combinations of product wavefunctions occur. To express the same thing more abstractly. The Hilbert space H'' representing states of the 2-particle system is initially *assumed* to be a subspace of the tensor product $H' \otimes H'$, where H' is the Hilbert space of 1-particle states; reasons are then given to show that H'' must coincide with either the symmetric (Bose case) or the antisymmetric (Fermi case) subspace of $H' \otimes H'$. Not only is it clearly undesirable to pretend, even initially, that indistinguishable particles can be labeled, but it is unsatisfactory also to *assume* any direct relation between H'' and $H' \otimes H'$. Both these defects are avoided in the present approach; indeed, even in treating two *distinguishable* systems, with Hilbert spaces J and K , we do not postulate any connection between the Hilbert space of states of the joint system (formed by uniting the two systems) and the tensor product $J \otimes K$. Instead the existence, under certain conditions, of a natural isomorphism between these two Hilbert spaces is deduced as a theorem.

Naturally, we must replace the conventional assumption about the relation between H'' and H' by some alternative. It is necessary to embody in the mathematics the physical fact that H'' corresponds to states of pairs of particles whose states as individuals are represented by elements of H' . This is done in the following way.

First, observe that strictly speaking, a state is represented not by a vector ψ , but by the *ray* Ψ on which that vector lies. Now suppose that Ψ and Φ are two 1-particle states, i.e., elements of H' (by which we denote the set of all rays of H'). We then denote by $\Psi \circ \Phi$ the 2-particle state (if such exists) in which there is a particle in state Ψ and a particle in state Φ ; $\Psi \circ \Phi$ is thus a ray of H'' . Our procedure is to put

* This work was supported by a grant from the National Research Council of Canada. A large part of it was carried out while the author was a member of the Summer Research Institute of the Canadian Mathematical Congress.

¹M. Dresden, 1963 Brandeis University Summer Institute Lectures in Theoretical Physics, Vol. 2.

forward as axioms certain physically plausible properties of \circ ; for instance, we clearly have $\psi \circ \varphi = \varphi \circ \psi$.

Before proceeding to other properties of the operation \circ , a remark should be made. It is clearly unwise to assume that $\psi \circ \varphi$ always exists (consider, for instance, the Fermi case with $\psi = \varphi$). We make instead the weaker assumption that for certain systems which we call *simple* systems, $\psi \circ \varphi$ exists whenever $\psi \perp \varphi$. This is reasonable, since there is, in a certain sense, no interference between orthogonal states. We can conceive of the possibility of both states being prepared simultaneously; such a simultaneous preparation constitutes a preparation of $\psi \circ \varphi$. We assume, or better, *define* H'' to be the smallest Hilbert space which contains all rays of the form $\psi \circ \varphi$ with ψ and φ rays of H' .

The concept of the transition probability amplitude between states ψ and φ is defined in the usual way as $\langle \psi, \varphi \rangle = |\langle \psi, \varphi \rangle|$, where ψ and φ are unit vectors on the rays ψ and φ ; thus $\langle \psi, \varphi \rangle^2$ is the probability that a system in the state ψ will, when subjected to a φ measurement, be found in the state φ . Our assumptions concerning the operation \circ are expressed in the following axiom.

Axiom 1.1: If $\psi, \psi', \varphi, \varphi'$ are states of a simple system with $\psi \perp \varphi$ and $\psi' \perp \varphi'$, then

- (a) $\psi \circ \varphi = \varphi \circ \psi$,
- (b) $\langle \psi' \circ \varphi', \psi \circ \varphi \rangle = 0$ whenever $\psi' \perp \psi$ and $\varphi' \perp \varphi$,
- (c) $\langle \psi' \circ \varphi', \psi \circ \varphi \rangle = \langle \psi', \psi \rangle \langle \varphi', \varphi \rangle$ whenever $\psi' \perp \varphi$ and $\varphi' \perp \psi$.

The physical motivation for these postulates is clear. Thus $\langle \psi' \circ \varphi', \psi \circ \varphi \rangle^2$ is the probability that the two particles will make transitions into the states ψ and φ , respectively. If both ψ' and φ' are orthogonal to ψ , then neither particle can make the transition to ψ , so this probability must be zero. Similarly, if $\psi' \perp \varphi$ and $\varphi' \perp \psi$, only the particle in state ψ' can make a transition to ψ and only that in φ' can make one to φ , the respective probabilities of these events being $\langle \psi', \psi \rangle^2$ and $\langle \varphi', \varphi \rangle^2$. Our postulate then follows from the hypothesis that these events are independent, which is here particularly credible, since both ψ and ψ' are orthogonal to both φ and φ' .

In Sec. 3 a somewhat more detailed discussion of the operation \circ is given. The rest of the paper is devoted to showing that, at least in so far as the relation between H' and H'' is concerned, every simple system obeys either Fermi or Bose statistics. More precisely, we show (Theorems 6.5 and 6.6) that there

is a natural isomorphism U_{\pm} between the symmetric (antisymmetric) subspace of the tensor product $H' \otimes H'$ and H'' such that, for any orthogonal vectors ψ and φ of H' , $\psi \circ \varphi$ is the ray determined by the vector $U_{\pm}[(\psi \otimes \varphi \pm \varphi \otimes \psi)/\sqrt{2}]$. (The upper and lower signs refer to the Bose and Fermi cases, respectively.)

The route to this result is somewhat devious, essentially because of the weakness of our axiom, which makes no assertion about the value of $\langle \psi' \circ \varphi', \psi \circ \varphi \rangle$, except in cases where the rays involved exhibit certain orthogonality relations.

First, in Sec. 4, we introduce the concept of *disjoint systems*. Two systems, with Hilbert spaces J and K , are disjoint if $\psi \circ \varphi$ exists whenever ψ and φ are states of J and K , respectively. For such systems we show (Theorem 4.2) that the usual representation of the union of the systems by the tensor product $J \otimes K$ is justified. In proving this we require that the conditions (b) and (c) of Axiom 1.1 hold whenever ψ and ψ' are states of J and φ and φ' are states of K . An important case, and the only case in which we use this result, is where J and K are orthogonal subspaces of a simple system H' . Such a pair of systems is by definition disjoint. Moreover, in this case the requirement just mentioned is automatically satisfied as a consequence of the definition and Axiom 1.1.

The technique involved in proving Theorem 4.2 is similar to that used in the original proof of Wigner's theorem (Theorem 2.4). This theorem² allows one to "lift" a given isometry V between ray spaces H_1 and H_2 to a (linear or antilinear) isometry V between the corresponding Hilbert spaces, i.e., V is constructed so that $\psi_2 = V\psi_1$ whenever $\varphi_2 = V\varphi_1$. Here we have a mapping $U: J \times K \rightarrow H$ given by $U(\psi, \varphi) = \psi \circ \varphi$ and want to "lift" it to a corresponding mapping $U_{JK}: J \otimes K \rightarrow H$. By fixing either the first or the second variable, we get from U an isometry to which Wigner's theorem can be applied; then we use the vector mappings obtained in this way to construct the required mapping U_{JK} .

In this construction U_{JK} is defined by giving its action on the elements of a (orthonormal) basis.³ This treatment, though mathematically inelegant, is convenient for two reasons: first, owing to the form of Axiom 1.1, we are frequently concerned with orthogonal vectors; second, it suits the simplest description of the tensor product $J \otimes K$. If $\{\psi_j\}$ and $\{\varphi_k\}$ are bases for J and K , respectively, then $\{\psi_j \otimes \varphi_k\}$ is a

² E. P. Wigner, *Group Theory* (Academic Press Inc., New York, 1959), pp. 233-236. See also V. Bargmann, *J. Math. Phys.* 5, 862 (1964) and further references given there.

³ Throughout the paper the term "basis" will be used as an abbreviation for "orthonormal basis."

basis for $J \otimes K$. Similar remarks apply to the later sections of the paper.

In Sec. 5 the same methods are applied to the case of two indistinguishable simple systems. Here the states of each system are represented by the rays of the *same* Hilbert space H' . (To deny this would be to admit the existence of an observable which takes different values for the two systems, i.e., to admit their distinguishability.) One cannot treat this case by simply putting $\mathbf{J} = \mathbf{K} = \mathbf{H}'$, since the mapping U which was, in Sec. 4, defined on $\mathbf{J} \times \mathbf{K}$ is not now defined on $\mathbf{H}' \times \mathbf{H}'$. Its domain consists only of those pairs $\{\psi, \varphi\}$ of rays of \mathbf{H}' such that $\psi \perp \varphi$. For this reason the mapping U does not give us direct access to the tensor product $H' \otimes H'$. If $B = \{\psi_i\}$ is a basis for H' , U gives the direction of the rays corresponding to the vectors $\psi_i \otimes \psi_j$ only for $i \neq j$. Our procedure is therefore to work first only with the subspace (denoted $H' \otimes_B H'$) of $H' \otimes H'$ for which these vectors form a basis. By methods similar to those of Sec. 4, though more complicated, we construct (Theorem 5.6) a linear mapping $U_B: H' \otimes_B H' \rightarrow H''$ which "lifts" the mapping U , in that $U_B(\psi \otimes \varphi)$ always lies on the ray $\psi \circ \varphi$. Of course, U_B is not an isometry. For example, $U_B(\varphi \otimes \varphi)$ and $U_B(\varphi \otimes \psi)$ both lie on $\psi \circ \varphi$, whereas $\psi \otimes \varphi$ and $\varphi \otimes \psi$ are orthogonal in $H' \otimes_B H'$. We show, however (Theorem 5.7), that *either*, for every B , ψ , and φ , $U_B(\psi \otimes \varphi) = -U_B(\varphi \otimes \psi)$ or, for every B , ψ , and φ , $U_B(\psi \otimes \varphi) = U_B(\varphi \otimes \psi)$. In this way the Fermi and Bose cases appear.

Section 6 is devoted to the construction of the mappings U_- and U_+ in the Fermi and Bose cases, respectively. The Fermi case is much the easier, owing to the fact that $H' \otimes_B H'$ contains a basis for the antisymmetric tensor product $H' \otimes_- H'$, which is the domain of the mapping U_- . The construction of U_- is thus trivial ($U_- = U_B | H' \otimes_- H'$), although some work has to be done to show that it has the desired property even when $\psi \otimes \varphi$ does not lie in $H' \otimes_B H'$. In the Bose case we have the problem, in constructing U_+ , of defining $U_+(\psi_i \otimes \psi_i)$, since $\psi_i \otimes \psi_i$ does not lie in the domain of U_B . The difficulty is overcome by expressing $\psi_i \otimes \psi_i$ as a linear combination of two vectors which lie in the domains of $U_{B'_{ij}}$ and $U_{B''_{ij}}$, respectively, where B'_{ij} and B''_{ij} are certain bases closely related to B .

It should be noted that our derivation of Fermi and Bose statistics from Axiom 1.1 applies only to the relation between 1-particle and 2-particle states. Axiom 1.1 does not imply anything significant about a system of N particles for $N > 2$. In this connection, however, two points should be noted. First, Dresden

(Ref. 1, p. 467) has shown that parastatistics must manifest itself already in the behavior of 2-particle states. It would seem, therefore, that parastatistics is incompatible with Axiom 1.1. Second, there is a natural generalization of Axiom 1.1 which may well suffice for a derivation of Fermi and Bose statistics valid for all N . We first make the physically natural postulate that \circ is associative and commutative whenever it is defined, and then define a simple system as one which, for every positive integer N , $\psi_1 \circ \cdots \circ \psi_N$ exists whenever ψ_1, \cdots, ψ_N is an orthogonal set. We then postulate.

If ψ_1, \cdots, ψ_N and $\varphi_1, \cdots, \varphi_N$ are two orthogonal sets of states of a simple system (N being any positive integer), then

- (i) $\langle \psi_1 \circ \cdots \circ \psi_N, \varphi_1 \circ \cdots \circ \varphi_N \rangle = 0$
if $\psi_i \perp \varphi_i$ for $1 \leq i \leq N$,
- (ii) $\langle \psi_1 \circ \cdots \circ \psi_N, \varphi_1 \circ \cdots \circ \varphi_N \rangle$
 $= \langle \psi_1, \varphi_1 \rangle \langle \psi_2, \varphi_2 \rangle \cdots \langle \psi_N, \varphi_N \rangle$
if $\psi_i \perp \varphi_j$ whenever $i \neq j$.

2. QUANTUM MECHANICS WITH SUPERSELECTION RULES

First some mathematical preliminaries.

Definitions 2.1: Let H be a complex Hilbert space with inner product $\langle \cdot, \cdot \rangle$ (assumed linear in the second variable). Any vector ψ in H determines a 1-dimensional subspace or *ray* denoted ray ψ or ψ . Let \mathbf{H} denote the set of all rays in H . In \mathbf{H} we define a *scalar product* by

$$\langle \psi, \varphi \rangle = |\langle \psi, \varphi \rangle| / (\langle \psi, \psi \rangle \langle \varphi, \varphi \rangle)^{\frac{1}{2}}$$

and the *distance* or *angle* $d(\psi, \varphi)$ between ψ and φ by $d(\psi, \varphi) = \cos^{-1} \langle \psi, \varphi \rangle$. We say ψ and φ are *orthogonal* if $\langle \psi, \varphi \rangle = 0$ and write $\psi \perp \varphi$. \mathbf{H} is then a complete metric space. Any metric space which is isomorphic to \mathbf{H} for some H will be called a *ray space*, and H will be called a *representative* of \mathbf{H} . By the *dimension* $\dim \mathbf{H}$ of the ray space \mathbf{H} we mean $\dim H$ and by a *subspace* of a ray space we mean a subset which is also a ray space.

Definitions 2.2: Let H and K be Hilbert spaces and U a mapping from H into K . Then,

- (a) U is *linear* if $U(a\psi + b\varphi) = aU\psi + bU\varphi$,
 - (b) U is *antilinear* if $U(a\psi + b\varphi) = \bar{a}U\psi + \bar{b}U\varphi$,
 - (c) U is an *isometry* if $\|U\psi - U\varphi\| = \|\psi - \varphi\|$,
- where in each case the relation holds for all ψ and φ in H and all numbers a and b .

Definition 2.3: Let H be a Hilbert space and let $\psi \rightarrow \psi^*$ be a one-to-one correspondence of H onto a

set H^* . Now make H^* a Hilbert space by setting

$$\begin{aligned}\psi^* + \varphi^* &= (\psi + \varphi)^*, \\ a\psi^* &= (\bar{a}\psi)^*, \\ \langle \psi^*, \varphi^* \rangle &= \overline{\langle \psi, \varphi \rangle},\end{aligned}$$

for all numbers a and elements ψ, φ of H . Then H^* is the *complex conjugate* of H . Clearly $\psi \rightarrow \psi^*$ is an antilinear isometry. One may regard H and H^* as the same set with the same law of addition, but with the complex conjugate law of multiplication by numbers and inner product. With this understanding the Hilbert spaces H and H^* determine the same ray space, $\mathbf{H} = \mathbf{H}^*$.

If H and K are Hilbert spaces, a linear or antilinear isometry $U: H \rightarrow K$ defines in an obvious way an isometry $\mathbf{U}: \mathbf{H} \rightarrow \mathbf{K}$. Conversely,

Theorem 2.4 (Wigner's theorem): Let H and K be Hilbert spaces with $\dim H \geq 2$. Let \mathbf{T} be an isometry of \mathbf{H} into \mathbf{K} . Then there exists either a linear isometry U or an antilinear isometry U (but not both) from H onto a subspace K' of K such that $\mathbf{U} = \mathbf{T}$. Moreover, U is unique up to a phase factor. We say that $\mathbf{T}: \mathbf{H} \rightarrow \mathbf{K}$ is *lifted* to $U: H \rightarrow K$.

For brevity we omit a proof.^{2,4}

Corollary 2.5: Let K be a Hilbert space and \mathbf{H} be a ray space contained in \mathbf{K} . Then $K' = \{\psi \in K: \psi \in \mathbf{H}\}$ is a closed subspace of K .

Proof: The identity map of \mathbf{H} into \mathbf{K} is an isometry and so can be lifted to a linear or antilinear isometry of H onto the closed subspace K' of K .

Taking $\mathbf{H} = \mathbf{K}$, we obtain,

Corollary 2.6: There are, up to isomorphism, exactly two Hilbert spaces representing any given ray space of dimension ≥ 2 , and these are complex conjugates of each other.

Let \mathbf{H} and \mathbf{K} be ray spaces with $\dim \mathbf{H} \geq 2$, and \mathbf{U} be an isometry of \mathbf{H} into \mathbf{K} . If specific Hilbert spaces H and K representing \mathbf{H} and \mathbf{K} have been selected, we can assign a *type* (linear or antilinear) to \mathbf{U} according to whether the mapping $U: H \rightarrow K$, obtained by lifting \mathbf{U} , is linear or antilinear. In the absence of a specific choice of H and K , no type can be assigned to a single \mathbf{U} , but we can still say, of two isometries \mathbf{U}_1 and \mathbf{U}_2 , whether or not they are of the same type.

⁴ More details of some of this work are given in mimeographed lecture notes available from the author on request.

The situation is sharper if \mathbf{H} is a subspace of \mathbf{K} . Then we naturally require that H be chosen as the corresponding subspace of K and each isometry $\mathbf{U}: \mathbf{H} \rightarrow \mathbf{K}$ has a unique type. The identity map is clearly of linear type. We might suspect from this that every isometry of \mathbf{H} into \mathbf{K} which differs little from the identity is also of linear type. This conjecture is established in a strong form in the following theorem.

Theorem 2.7: Let \mathbf{K} be a ray space, \mathbf{H} a subspace of \mathbf{K} of dimension ≥ 2 , and $\mathbf{U}: \mathbf{H} \rightarrow \mathbf{K}$ an isometry of antilinear type. Then there is a ray ψ in \mathbf{H} with $\mathbf{U}\psi \perp \psi$.

Proof: Let ψ_1 and ψ_2 be fixed linearly independent vectors in H and put $\psi = a\psi_1 + b\psi_2$. Then $\langle U\psi, \psi \rangle$ is a homogeneous quadratic function of a and b and therefore vanishes for some pair of complex numbers a, b not both zero.

We now state an axiom which describes the structure of conventional quantum mechanics with superselection rules. Because we shall be discussing the process of "uniting" systems to form larger systems, the ray space \mathbf{H} which is now introduced should be regarded as a "global ray space" containing rays corresponding to every state of every system.

Axiom 2.8: There is a Hilbert space H which is a direct sum $H = \bigoplus_{i \in I} H_i$ of a set of orthogonal subspaces called *coherent subspaces* such that:

(a) Every state of every system is represented by a unique ray in H lying in a coherent subspace; every ray which lies in a coherent subspace represents a physical state.⁵ (We shall call these rays *physical rays*.)

(b) Every observable is represented by a Hermitian operator on H which commutes with all the projections on the coherent subspaces; conversely, every such Hermitian operator represents an observable.

(c) The expectation value of an observable A in a state ψ is $\langle \psi, A\psi \rangle$, where ψ is a unit vector on the ray ψ .

If A is the projection on a physical ray φ , then by (c) its expectation value is $|\langle \psi, \varphi \rangle|^2 = \langle \psi, \varphi \rangle^2$; thus $\langle \psi, \varphi \rangle^2$ gives the probability that a system prepared in the state ψ will, when the appropriate measurement is made, make a transition into the state φ .

It should be noted that the physically significant object is the ray space \mathbf{H} . In ordinary quantum mechanics, where the set I contains only one element,

⁵ It is possible⁴ to deduce this structure of physical rays from the following hypothesis: If two physical rays are nonorthogonal, then every ray in the subspace spanned by them is physical.

there are (by Corollary 2.6) two possible choices for the Hilbert space H representing \mathbf{H} . For quantum mechanics with superselection rules, the arbitrariness in H is much greater; for each i in I we may independently choose for H_i one of the two possible Hilbert spaces which represent \mathbf{H}_i and then set $H = \bigoplus_{i \in I} H_i$.

A final word concerning our use of the term "system." For mathematical purposes it is sufficient to identify a system with the set of all its states, i.e., with a subset of \mathbf{H} . In agreement with the convention in quantum mechanics that the states of a system are represented by the rays of a Hilbert space, we adopt the following definition.

Definition 2.9: A system is a closed subspace of \mathbf{H} consisting entirely of physical rays (i.e., it is a closed subspace of a coherent subspace).

3. UNION OF STATES

In any experiment in which the interaction of two systems is under study, for example in a scattering experiment, the first step consists of the preparation of the two systems. This involves the simultaneous preparation of two states ψ and φ , one for each system. We must assume that procedures for the preparation of these states are known; the new feature consists of the requirement that the two states be prepared simultaneously. The design of the experiment thus requires, among other things, the discovery of suitable methods of preparation of the two systems—methods which do not interfere with each other. For many pairs of states ψ, φ , such a pair of compatible methods of preparation exists; we then denote by $\psi \circ \varphi$ the state which is prepared in this way by simultaneously preparing the states ψ and φ and call it the *union* of ψ and φ .⁶ When $\psi \circ \varphi$ exists we shall call the states ψ and φ *disjoint*.⁷

It is likely that, unless ψ and φ are orthogonal, $\psi \circ \varphi$ will not exist (see Sec. 1). We shall accordingly only contemplate the formation of unions in the case of orthogonal states. However, orthogonal states are not necessarily disjoint. Suppose, for instance, that ψ and φ are states of a hydrogen atom; suppose that in both cases the proton is in a given state (for instance an eigenstate of momentum), while the electron is, in ψ , in the ground state, and, in φ , in the first excited state. Then $\psi \perp \varphi$, but it seems probable that $\psi \circ \varphi$ does not exist, since its preparation would

involve the simultaneous preparation of two protons in the same state.

This possibility—that two states can be orthogonal but not disjoint—clearly arises from the fact that the system involved is a composite one; for a single electron, for instance, no analogous example presents itself. We shall regard this as a characteristic feature of a composite system and adopt the following definition.

Definition 3.1: A system is *simple* if every pair of orthogonal states is disjoint.

Axiom 1.1, on which our derivation of quantum statistics depends, applies only to simple systems. It is easy to give examples of composite systems for which it does not hold⁴; thus the derivation does not apply to such systems. On the other hand, this is as it should be since composite systems only obey Fermi or Bose statistics, if at all, in an approximate way. An assembly of hydrogen atoms, for instance, obeys Bose statistics only in so far as the exchange of electrons can be neglected.

4. UNION OF DISJOINT SYSTEMS

Extending the terminology introduced in Sec. 3, we shall call two systems (with ray spaces) \mathbf{J} and \mathbf{K} *disjoint* if $\psi \circ \varphi$ exists whenever ψ is in \mathbf{J} and φ is in \mathbf{K} . For the purposes of *this section only* we shall assume:

Axiom 4.1: The conditions of Axiom 1.1 are satisfied whenever ψ and ψ' belong to \mathbf{J} , and φ and φ' belong to \mathbf{K} , \mathbf{J} and \mathbf{K} being disjoint systems.

In subsequent sections we shall be concerned only with the case where \mathbf{J} and \mathbf{K} are orthogonal closed subspaces of a simple system \mathbf{H}' ; in this case the assertion of Axiom 4.1 follows from Axiom 1.1.

This section is devoted to proving the following theorem, which justifies the customary representation of the Hilbert space corresponding to the union of two disjoint systems as the tensor product of the Hilbert spaces corresponding to the separate systems. Throughout, H denotes the global Hilbert space introduced in Axiom 2.8.

Theorem 4.2: Let \mathbf{J} and \mathbf{K} be disjoint systems and for each ψ in \mathbf{J} and φ in \mathbf{K} , let $U(\psi, \varphi) = \psi \circ \varphi$. Then Hilbert spaces J and K representing \mathbf{J} and \mathbf{K} can be chosen in such a way that there is a linear isometry $U_{JK}: J \otimes K \rightarrow H$ such that, for every ψ in J and φ in K with $\psi \neq 0$ and $\varphi \neq 0$,

$$U(\psi, \varphi) = \text{ray } U_{JK}(\psi \otimes \varphi),$$

and U_{JK} is unique up to a phase factor. Further, for

⁶ See, R. Giles, *Mathematical Foundations of Thermodynamics* (Pergamon Press, Inc., Oxford, England, 1964), where essentially the same operation of union is denoted by "+."

⁷ This is a generalization of the concept *strongly orthogonal*, used in work on the many-body problem in quantum mechanics. See, for instance, A. J. Coleman, *Rev. Mod. Phys.* **35**, 68 (1963).

any ψ, ψ' in J and φ, φ' in K ,

$$\langle U_{JK}(\psi \otimes \varphi), U_{JK}(\psi' \otimes \varphi') \rangle = \langle \psi, \psi' \rangle \langle \varphi, \varphi' \rangle.$$

Proof: The proof is via a series of lemmas. Throughout, ψ and φ , with or without embellishments, denote arbitrary elements of J and K , respectively. We first observe that J and K , being disjoint, are orthogonal so that Axiom 1.1 gives

Lemma 1: For any $\psi_1, \psi_2, \varphi_1, \varphi_2$,

$$\langle \psi_1 \circ \varphi_1, \psi_2 \circ \varphi_2 \rangle = \langle \psi_1, \psi_2 \rangle \langle \varphi_1, \varphi_2 \rangle.$$

Lemma 2: The range of U lies in a single coherent subspace.

Proof: Let $\psi_1, \psi_2, \varphi_1, \varphi_2$ be given. Choose ψ nonorthogonal to both ψ_1 and ψ_2 , and φ nonorthogonal to both φ_1 and φ_2 . Then, by Lemma 1, each ray in the following sequence is nonorthogonal to the next: $\psi_1 \circ \varphi_1, \psi \circ \varphi, \psi_2 \circ \varphi_2$. Thus all these rays must lie in the same coherent subspace.

Lemma 3: Let G denote the ray space spanned by the range of U . Let $U(\psi \cdot): K \rightarrow G$ be defined for each ψ in J by $U(\psi \cdot)\varphi = \psi \circ \varphi$. Then $U(\psi \cdot)$ is an isometry. Moreover, the type (linear or antilinear) of $U(\psi \cdot)$ is the same for every ψ in J .

Proof: That $U(\psi \cdot)$ is an isometry follows from Lemma 1. Next, for any ψ_1 and ψ_2 , $U(\psi_1 \cdot)[U(\psi_2 \cdot)]^{-1}$ is a mapping of the subspace $U(\psi_2 \cdot)K$ of G into G and if $\theta = \psi_2 \circ \varphi$ is any element of $U(\psi_2 \cdot)K$, then

$$\begin{aligned} \langle U(\psi_1 \cdot)[U(\psi_2 \cdot)]^{-1}\theta, \theta \rangle &= \langle \psi_1 \circ \varphi, \psi_2 \circ \varphi \rangle \\ &= \langle \psi_1, \psi_2 \rangle. \end{aligned}$$

Now Lemma 3 clearly holds if $\dim K = 1$, so we may assume $\dim K \geq 2$. But then, by Theorem 2.7, if $\langle \psi_1, \psi_2 \rangle \neq 0$, $U(\psi_1 \cdot)[U(\psi_2 \cdot)]^{-1}$ is of linear type, so that $U(\psi_1 \cdot)$ and $U(\psi_2 \cdot)$ are of the same type. The same must apply even if $\langle \psi_1, \psi_2 \rangle = 0$, since we can always choose ψ_3 nonorthogonal to both ψ_1 and ψ_2 .

In the same way we prove

Lemma 4: Let $U(\cdot\varphi): J \rightarrow G$ be defined for each φ in K by $U(\cdot\varphi)\psi = \psi \circ \varphi$. Then $U(\cdot\varphi)$ is an isometry and its type is independent of the choice of the ray φ .

Now let a Hilbert space G representing G be chosen arbitrarily. Then choose Hilbert spaces J and K representing J and K in such a way that all the mappings $U(\psi \cdot)$ and $U(\cdot\varphi)$ are of linear type.⁸

⁸ If, instead, the spaces J and K are prescribed in advance, the mapping U will belong to one of four types. If all the mappings $U(\cdot\varphi)$ are of linear (antilinear) type, we say U is of linear (antilinear) type in the first variable, and similarly for the second variable.

Next, for each unit vector ψ in J , let $U_{JK}(\psi \cdot): K \rightarrow G$ denote a linear isometry obtained by lifting the isometry $U(\psi \cdot)$; i.e., $U_{JK}(\psi \cdot)$ satisfies

$$\text{ray } U_{JK}(\psi \cdot)\varphi = U(\psi \cdot)\varphi = \psi \circ \varphi,$$

for every vector φ in K . Similarly, for each unit vector φ in K , let $U_{JK}(\cdot\varphi): J \rightarrow G$ denote a linear isometry obtained by lifting the isometry $U(\cdot\varphi)$. By Wigner's theorem the operators $U_{JK}(\psi \cdot)$ and $U_{JK}(\cdot\varphi)$ all exist, and each is arbitrary up to a phase factor.

We now use these operators to construct the operator $U_{JK}: J \otimes K \rightarrow G$ of the theorem. Let $\{\psi_j: j \in \mathfrak{J}\}$ and $\{\varphi_k: k \in \mathfrak{K}\}$ be bases in J and K .⁹ Then $\{\psi_j \circ \varphi_k: j \in \mathfrak{J}, k \in \mathfrak{K}\}$ is an orthogonal set of rays in G .⁹ Let θ_{jk} denote a unit vector on the ray $\psi_j \circ \varphi_k$. We now choose the phases of the θ_{jk} and of the linear isometries $U_{JK}(\psi_j \cdot)$ and $U_{JK}(\cdot\varphi_k)$ as follows:

Choose θ_{11} arbitrarily.

Choose $U_{JK}(\cdot\varphi_1)$ so that $U_{JK}(\cdot\varphi_1)\psi_1 = \theta_{11}$.

For $j \neq 1$, choose θ_{j1} so that $U_{JK}(\cdot\varphi_1)\psi_j = \theta_{j1}$.

For each j , choose $U_{JK}(\psi_j \cdot)$ so that $U_{JK}(\psi_j \cdot)\varphi_1 = \theta_{j1}$.

For $k \neq 1$, choose θ_{jk} so that $U_{JK}(\psi_j \cdot)\varphi_k = \theta_{jk}$.

Henceforth let $\psi = \sum_j a_j \psi_j$ and $\varphi = \sum_k b_k \varphi_k$ denote arbitrary unit vectors in J and K , respectively, and let $\theta = \sum_{j,k} c_{jk} \theta_{jk} + \tilde{\theta}$, where for every j and k $\theta_{jk} \perp \tilde{\theta}$, be a unit vector on the ray $\theta = \psi \circ \varphi$. We shall show that by suitably choosing the phase of θ we have $c_{jk} = a_j b_k$ for every ψ and φ .

Lemma 5: For each j, k , $c_{jk} = t_{jk} a_j b_k$, where $|t_{jk}| = 1$; also $\tilde{\theta} = 0$.

Proof: By Lemma 1, $\langle \psi_j \circ \varphi_k, \psi \circ \varphi \rangle = \langle \psi_j, \psi \rangle \langle \varphi_k, \varphi \rangle$, which gives $|c_{jk}| = |a_j| |b_k|$. It follows that $\sum_{jk} c_{jk} \theta_{jk}$ is a unit vector, whence $\tilde{\theta} = 0$.

The phase factor t_{jk} may be expected to depend on j and k . If the pair $\langle \psi, \varphi \rangle$ is such that it does not, then a suitable choice of the phase of θ will yield $t_{jk} = 1$ and so $c_{jk} = a_j b_k$. Let us denote by Ψ the set of all pairs $\langle \psi, \varphi \rangle$ for which this is so. Since $U_{JK}(\cdot\varphi_1)$ is linear and $U_{JK}(\cdot\varphi_1)\psi_j = \theta_{j1}$ for every j , $\langle \psi, \varphi_1 \rangle \in \Psi$ for every ψ . Since $U_{JK}(\psi_j \cdot)$ is linear and $U_{JK}(\psi_j, \varphi_k) = \theta_{jk}$ for every k , $\langle \psi_j, \varphi \rangle \in \Psi$ for every φ . We can now prove

Lemma 6: For every ψ and φ , $\langle \psi, \varphi \rangle \in \Psi$.

⁹ Throughout this section an index j will range always over the set \mathfrak{J} and k over \mathfrak{K} . We do not assume these sets to be countable; nevertheless, we denote by 1 a particular element (chosen once and for all) of either set.

Proof: Case 1: $b_1 \neq 0$. Since, for any j , $\langle \psi_j, \varphi \rangle \in \Psi$, the equation $\langle \psi_j \circ \varphi, \psi \circ \varphi \rangle = \langle \psi_j, \psi \rangle$ yields

$$\left| \left\langle \sum_k b_k \theta_{jk}, \sum_{jk} t_{jk} a_j b_k \theta_{jk} \right\rangle \right| = |a_j|,$$

i.e.,

$$\left| \sum_k \bar{b}_k b_k t_{jk} \right| |a_j| = |a_j|.$$

Now t_{jk} is only defined if $a_j \neq 0$ and $b_k \neq 0$, so we may restrict attention to such j and k . Also $\sum_k \bar{b}_k b_k = 1$. It follows that, for each (such) j , t_{jk} is independent of k . Let us write $t_{jk} = t_j$.

In the same way we can show, using the fact that $\{\psi, \varphi_1\} \in \Psi$, that t_j is independent of j . Thus t_{jk} is independent of j and k .

Case 2: $b_1 = 0$. This may be treated as a limiting case of Case 1 or we may proceed as follows. Let $\varphi' = (\varphi_1 + \varphi)/\sqrt{2}$. Then $\{\psi, \varphi'\} \in \Psi$, by Case 1. Using this fact, an argument of the same type as that used for Case 1 shows that t_{jk} is independent of j and k .

The proof of Theorem 4.2 is now easily completed. Define a linear isometry $U_{JK}: J \otimes K \rightarrow G$ by setting $U_{JK}(\psi_j \otimes \varphi_k) = \theta_{jk}$. Then for every ψ and φ

$$U_{JK}(\psi \otimes \varphi) = \sum_{jk} a_j b_k U_{JK}(\psi_j \otimes \varphi_k) = \sum_{jk} a_j b_j \theta_{jk},$$

which, by Lemma 6, lies in the ray $\psi \circ \varphi$. Thus U_{JK} has the required property.

To establish the uniqueness (up to a phase factor) of U_{JK} , we need merely note that if U'_{JK} is any linear isometry satisfying the condition of Theorem 4.2, then the restriction of U'_{JK} to the subspace $J \otimes \varphi_1$ of $J \otimes K$ determines the same ray mapping of J into G as does $U_{JK}(\cdot \otimes \varphi_1)$. It follows from Wigner's theorem that, by altering U'_{JK} (if necessary) by a phase factor, we can make these isometries coincide so that, in particular

$$U'_{JK}(\psi_j \otimes \varphi_1) = U_{JK}(\cdot \otimes \varphi_1) \psi_j = \theta_{j1} \quad \text{for all } j.$$

Similarly, for each j , the restriction of U'_{JK} to the subspace $\psi_j \otimes K$ of $J \otimes K$ determines the same ray mapping of K into G as does $U_{JK}(\psi_j \cdot)$. Moreover, these isometries now agree at $\varphi = \varphi_1$, so they coincide. In particular

$$U'_{JK}(\psi_j \otimes \varphi_k) = U_{JK}(\psi_j \cdot) \varphi_k = \theta_{jk} = U_{JK}(\psi_j \otimes \varphi_k).$$

Thus the linear isometries U'_{JK} and U_{JK} agree on a basis of $J \otimes K$ and consequently coincide. This completes the proof of uniqueness.

Lastly, the final statement in Theorem 4.2 now follows at once from the definition and linearity of U_{JK} .

5. SIMPLE SYSTEMS

In the rest of the paper H' denotes a simple system, and ψ, φ denotes arbitrary rays in H' . We shall study the mapping U defined by

$$U(\psi, \varphi) = \psi \circ \varphi,$$

whose domain is the subset $\{(\psi, \varphi) : \psi \in H', \varphi \in H', \psi \perp \varphi\}$ of $H' \times H'$. Let H' be the subspace of H spanned by the range of U . We shall call H' and H'' 1-particle and 2-particle space, respectively. Let H' and H'' be Hilbert spaces representing the ray spaces H' and H'' .¹⁰

Let J and K be orthogonal proper subspaces of H' , and let U_{JK} denote the restriction of U to $J \times K$. J and K are disjoint systems so the mapping U_{JK} can be assigned a "type."⁸ Let J^\perp denote the orthogonal complement of J in H' . We shall say J is of linear (antilinear) type if U_{JJ^\perp} is of linear (antilinear) type in the first variable. If $\dim J = 1$, J is of both linear and antilinear type; if $\dim J \geq 2$, J is of one type only.

If J_1 and J_2 are any subspaces of H' , the types of J_1 and J_2 of course depend on the choice of H' and H'' , but the validity of the statement " J_1 and J_2 are of the same type" is independent of that choice. We now prove:

Theorem 5.1: If $\dim H' \geq 4$ and J_1, J_2 are any proper subspaces of H' , then J_1 and J_2 are of the same type.

Proof: Since the type of J is that of the mapping from J into H'' given by $\psi \rightarrow \psi \circ \varphi$, where φ is an arbitrary element of J^\perp , any proper subspace of H' containing J has the same type as J . Let J'_1, J'_2 be subspaces of H' of codimension 1 containing J_1 and J_2 , respectively. Then $\dim J'_1 \cap J'_2 \geq 2$. It now follows that in the following sequence each subspace is of the same type as the next: $J_1, J'_1, J'_1 \cap J'_2, J'_2, J_2$. Since all these spaces except possibly J_1 and J_2 have dimension exceeding 1, J_1 and J_2 are of the same type.

We shall henceforth assume the representative H'' of H' to be chosen so that every proper subspace of H' is of linear type.

¹⁰ By Corollary 2.6 there is a two-fold choice in the selection of both H' and H'' ; however, as we shall see shortly, the choice of H'' will be matched to that of H' in a certain way. It follows that if H' and H'' lie in the same coherent subspace and if H' is chosen to be isomorphic to a subspace of the global Hilbert space H , then H'' may (in so far as our axioms are concerned) turn out to be isomorphic to the complex conjugate of a subspace of H . However, considerations beyond the scope of this paper (for instance, the existence of a lower bound to the energy) may be invoked to rule out this possibility.

Theorem 5.2: Let J, K be orthogonal proper subspaces of H' . Then there exists a linear isometry $U_{JK}: J \otimes K \rightarrow H''$ such that, for any ψ in J and φ in K with $\psi \neq 0$ and $\varphi \neq 0$,

$$\Psi \circ \Phi = \text{ray } U_{JK}(\psi \otimes \varphi).$$

Moreover, U_{JK} is unique up to a phase factor.

Proof: Since \mathbf{J} is of linear type, U_{JK} is of linear type in the first factor. Since $\mathbf{U}(\Psi, \Phi) = \mathbf{U}(\Phi, \Psi)$, it is also of linear type in the second factor. The result now follows from Theorem 4.2.

Definition 5.3: Let B be any orthonormal set in H' . We say a subspace J of H' is *compatible with B* if a subset of B is a basis for J . We say $\{J, K\}$ is an *orthogonal pair of subspaces compatible with B* if J and K are orthogonal proper subspaces of H' , both compatible with B . We say $\{\psi, \varphi\}$ is an *orthogonal pair of vectors compatible with B* if there exists an orthogonal pair of subspaces $\{J, K\}$ compatible with B , such that $\psi \in J$ and $\varphi \in K$.

We are now going to amalgamate into a single mapping U_B the mappings U_{JK} for all orthogonal pairs of subspaces $\{J, K\}$ compatible with B .

Definition 5.4: We shall say that two linear mappings *agree* if their restrictions to the intersection of their domains coincide. (Note that “agree” is not an equivalence relation.) The term *agree up to a phase factor* will be used in the same way.

Theorem 5.5: Let $B = \{\psi_i: i \in \mathfrak{J}\}$ be an orthonormal set in H' . Let $\{J, K\}$ and $\{J', K'\}$ be two orthogonal pairs of subspaces compatible with B . Then,

- (i) $(J \otimes K) \cap (J' \otimes K') = (J \cap J') \otimes (K \cap K')$ and
- (ii) U_{JK} and $U_{J'K'}$ agree up to a phase factor.

Proof: (i) Let $\{\psi_j: j \in \mathfrak{J}\}$ be a basis for J and define $\mathfrak{J}', \mathfrak{K}, \mathfrak{K}'$ similarly. Then a basis for $J \otimes K$ is $\{\psi_j \otimes \psi_k: j \in \mathfrak{J}, k \in \mathfrak{K}\}$ and similarly for $J' \otimes K'$. Hence a basis for $(J \otimes K) \cap (J' \otimes K')$ is $\{\psi_j \otimes \psi_k: j \in \mathfrak{J} \cap \mathfrak{J}', k \in \mathfrak{K} \cap \mathfrak{K}'\}$ and this is also a basis for $(J \cap J') \otimes (K \cap K')$.

(ii) By (i) the intersection of the domains is $(J \cap J') \otimes (K \cap K')$. But on this restricted domain the mappings U_{JK} and $U_{J'K'}$ both induce the ray mapping $\mathbf{U}_{J \cap J', K \cap K'}$. The result now follows from Theorem 4.2.

We shall make use of the following notation. Let B be any orthonormal set in a Hilbert space L . Then we

denote by $L \otimes_B L$ the closed subspace of $L \otimes L$ spanned by $\{\psi \otimes \varphi: \psi \in B, \varphi \in B, \psi \neq \varphi\}$.

Theorem 5.6: Let $\dim H' \geq 5$. Let B be an orthonormal set in H' with at least three elements. Then there exists a linear mapping $U_B: H' \otimes_B H' \rightarrow H''$ such that, whenever $\{\psi, \varphi\}$ is an orthogonal pair of unit vectors compatible with B , then $U_B(\psi \otimes \varphi)$ is a unit vector on the ray $\Psi \circ \Phi$.¹¹ Further, U_B is unique up to a phase factor.

Proof: Let $B = \{\psi_i: i \in \mathfrak{J}\}$. Let L be the subspace spanned by B , and for each i in \mathfrak{J} let L_i be the 1-dimensional subspace spanned by ψ_i .¹²

We first establish uniqueness. Let U_B and U'_B both satisfy the conditions of the theorem. Then, for $j \neq k$, $U'_B(\psi_j \otimes \psi_k) = t_{jk} U_B(\psi_j \otimes \psi_k)$, where $|t_{jk}| = 1$.

We shall show that all the phase factors t_{jk} ($j \neq k$) are equal. First note that, for any $j' \neq k$, $t_{jk} = t_{j'k}$; this follows from the fact that both U_B and U'_B agree up to a phase factor with $U_{L_j \oplus L_{j'}, L_k}$. Similarly, for any $k' \neq j$, $t_{jk} = t_{jk'}$. Then note that we can, by a sequence of two or three changes, each affecting only the first or last index, pass from any pair j, k to any other pair j', k' . (Here essential use is made of the condition $\dim L \geq 3$; otherwise we could not establish that $t_{12} = t_{21}$.) It follows that all the phase factors t_{jk} are equal, so that U_B and U'_B agree up to a phase factor.

To establish the existence of U_B , we first observe that it is sufficient to consider the case $\dim L \geq 5$, for if $B' \subset B$, then $U_B|_{H' \otimes_{B'} H'}$ satisfies the conditions required in the theorem for $U_{B'}$. Henceforth, then, we assume $\dim L \geq 5$. Let ψ_1, \dots, ψ_5 denote five elements of B .

We shall choose the phase of U_{JK} for various orthogonal pairs $\{J, K\}$ of subspaces compatible with B in such a way that the selected U_{JK} agree with each other.

First choose arbitrarily a unit vector θ_{12} on the ray $\theta_{12} = \Psi_1 \circ \Psi_2$. We shall say that the pair $\{J, K\}$ is of *type 1* if $J \supset L_1$ and $K \supset L_2$.¹³ For each such $\{J, K\}$, fix the phase of U_{JK} by setting $U_{JK}(\psi_1 \otimes \psi_2) = \theta_{12}$. By Theorem 5.5 we have at once:

Lemma 1: If $\{J, K\}$ and $\{J', K'\}$ are of type 1, then U_{JK} and $U_{J'K'}$ agree.

¹¹ By Theorem 5.2 this is equivalent to the condition that, whenever $\{J, K\}$ is an orthogonal pair of subspaces compatible with B , U_{JK} agrees up to a phase factor with U_B .

¹² Henceforth the indices i, j, k take values in \mathfrak{J} .

¹³ In the rest of the proof the expression “orthogonal pair of subspaces compatible with B ” will normally be abbreviated to “pair.”

We shall say the pair $\{J, K\}$ is of *type 2* if, for some $\{\tilde{J}, \tilde{K}\}$ of type 1, $J \subset \tilde{J}$ and $K \subset \tilde{K}$. Equivalently $\{J, K\}$ is of type 2 if $L_2 \not\subset J$ and $L_1 \not\subset K$. (Note that every pair of type 1 is also of type 2.) We choose the phase of U_{JK} for each $\{J, K\}$ of type 2 by first choosing $\{\tilde{J}, \tilde{K}\}$ of type 1 with $J \subset \tilde{J}$ and $K \subset \tilde{K}$ and then putting $U_{JK} = U_{\tilde{J}\tilde{K}}|_{J \otimes K}$. By Lemma 1 this specification of U_{JK} is independent of the choice of \tilde{J} and \tilde{K} .

Lemma 2: If $\{J, K\}$ and $\{J', K'\}$ are of type 2, then U_{JK} and $U_{J'K'}$ agree.

Proof: Let $\{\tilde{J}, \tilde{K}\}$ and $\{\tilde{J}', \tilde{K}'\}$ be corresponding pairs of type 1 selected as above. Then $U_{\tilde{J}\tilde{K}}$ and $U_{\tilde{J}'\tilde{K}'}$ agree by Lemma 1 and their domains contain those of U_{JK} and $U_{J'K'}$, respectively.

We shall say the pair $\{J, K\}$ is of *type 3* if, for some $\{\tilde{J}, \tilde{K}\}$ of type 2, $J = \tilde{J} \oplus L_2$ and $K = \tilde{K} \oplus L_1$. (\tilde{J} and \tilde{K} are uniquely fixed by J and K and neither \tilde{J} nor \tilde{K} contains L_1 or L_2 .) For each pair $\{J, K\}$ of type 3, we choose the phase of U_{JK} so that

$$U_{JK}|_{\tilde{J} \otimes \tilde{K}} = U_{\tilde{J}\tilde{K}}.$$

Lemma 3: If $\{J, K\}$ is of type 3 and $\{J', K'\}$ is of type 2, then U_{JK} and $U_{J'K'}$ agree.

Proof: Let $J = \tilde{J} \oplus L_2$ and $K = \tilde{K} \oplus L_1$. The intersection of the domains of U_{JK} and $U_{J'K'}$ is contained in $\tilde{J} \otimes \tilde{K}$. By Lemma 2 and the definition of U_{JK} , both U_{JK} and $U_{J'K'}$ agree with $U_{\tilde{J}\tilde{K}}$ when restricted to this intersection.

Lemma 4: If $\{J, K\}$ and $\{J', K'\}$ are both of type 3, then U_{JK} and $U_{J'K'}$ agree.

Proof: Let $J = \tilde{J} \oplus L_2, K = \tilde{K} \oplus L_1, J' = \tilde{J}' \oplus L_2, K' = \tilde{K}' \oplus L_1$, where $\tilde{J}, \tilde{K}, \tilde{J}', \tilde{K}'$ all contain neither L_1 nor L_2 . By Theorem 5.5 it suffices to show that U_{JK} and $U_{J'K'}$ coincide at one (nonzero) point in the intersection of their domains, namely $[(\tilde{J} \cap \tilde{J}') \oplus L_2] \otimes [(\tilde{K} \cap \tilde{K}') \oplus L_1]$. We consider two cases.

Case 1: $\dim(\tilde{J} \cap \tilde{J}') \geq 1$ and $\dim(\tilde{K} \cap \tilde{K}') \geq 1$. Let $\psi \in \tilde{J} \cap \tilde{J}'$ and $\varphi \in \tilde{K} \cap \tilde{K}'$ be nonzero vectors. Then in the following sequence each mapping agrees with the next at the point $\psi \otimes \varphi$: $U_{JK}, U_{\tilde{J}\tilde{K}}, U_{\tilde{J}'\tilde{K}'}, U_{J'K'}$.

Case 2: $\dim(\tilde{J} \cap \tilde{J}') = 0$ or $\dim(\tilde{K} \cap \tilde{K}') = 0$. By symmetry it suffices to treat the case $\dim(\tilde{J} \cap \tilde{J}') = 0$. Since \tilde{J} and \tilde{J}' are proper subspaces of L compatible with B , they each contain at least one element of B which we may without loss of generality assume to be ψ_3 and ψ_4 , respectively. Let K_3 and K_4 be the orthog-

onal complements of $L_3 \oplus L_2$ and $L_4 \oplus L_2$ in L . Then, by Case 1, each mapping in the following sequence agrees with the next:

$$U_{JK}, \quad U_{L_3 \oplus L_2, K_3}, \quad U_{L_4 \oplus L_3 \oplus L_2, L_5 \oplus L_1}, \\ U_{L_4 \oplus L_2, L_5 \oplus L_3 \oplus L_1}, \quad U_{L_4 \oplus L_2, K_4}, \quad U_{J'K'}.$$

Moreover, the domains of all these isometries contain $\psi_2 \otimes \psi_1$. Hence U_{JK} and $U_{J'K'}$ agree.

We can condense Lemmas 1–4 as follows. Let us say a pair $\{J, K\}$ is of *type 23* if it is of type 2 or of type 3. Then we have shown:

Lemma 5: If $\{J, K\}$ and $\{J', K'\}$ are both of type 23, then U_{JK} and $U_{J'K'}$ agree.

We now use this collection of U_{JK} s to define a linear mapping $U_B: H' \otimes_B H' \rightarrow H'$. For each pair $\{\psi_j, \psi_k\}$ ($j \neq k$) of elements of B , we can always choose a pair $\{J, K\}$ of type 23 such that $\psi_j \otimes \psi_k \in J \otimes K$ [e.g., if $j \neq 2$ and $k \neq 1$, take $\{L_j, L_k\}$, which is of type 2; if $j = 2$ and $k \neq 3$ (say), take $\{L_2 \oplus L_3, L_1 \oplus L_k \oplus L_4\}$, which is of type 3]. Now define $\theta_{jk} = U_{JK}(\psi_j \otimes \psi_k)$; by Lemma 5 the value of θ_{jk} is independent of the choice of J and K . Then define U_B by setting (for $j \neq k$) $U_B(\psi_j \otimes \psi_k) = \theta_{jk}$. This specifies U_B on a basis of $H' \otimes_B H'$; it is then determined by linearity and continuity on the whole of $H' \otimes_B H'$.

To complete the proof of Theorem 5.6, it only remains to show

Lemma 6: For any orthogonal pair $\{J, K\}$ of subspaces compatible with B , U_{JK} and U_B agree up to a phase factor.

Proof: By Theorem 4.2 it is sufficient to consider the case when the pair $\{J, K\}$ is *maximal*; i.e., when K is the orthogonal complement of J in L . Now, whenever $\psi_j \in J$ and $\psi_k \in K$, we have

$$U_{JK}(\psi_j \otimes \psi_k) = t_{jk} U_B(\psi_j \otimes \psi_k)$$

for some phase factor t_{jk} . We must show that t_{jk} is independent of j and k .

Now if $J \supset L_1$ and $K \supset L_2$, $\{J, K\}$ is of type 2; and if $J \supset L_2$ and $K \supset L_1$, $\{J, K\}$ is of type 3; in both cases the result follows at once from the definition of U_B . It is therefore only necessary to consider the situation when $L_1 \oplus L_2 \subset J$. (The case when $L_1 \oplus L_2 \subset K$ is similar.) This divides into two cases.

Case 1: $\dim J \geq 3$. Let $J = L_1 \oplus L_2 \oplus \tilde{J}$ where \tilde{J} is orthogonal to $L_1 \oplus L_2$. Let j, k, k' be such that

$\psi_j \in \tilde{J}, \psi_k \in K, \psi_{k'} \in K$. Then $\{L_1 \oplus L_j, L_k \oplus L_{k'}\}$ is a pair of type 2 so that U_B agrees with $U_{L_1 \oplus L_j, L_k \oplus L_{k'}}$ which (by Theorem 4.2) agrees up to a phase factor with U_{JK} . Since $\psi_1 \otimes \psi_k$ and $\psi_j \otimes \psi_{k'}$ both lie in the domains of all these mappings, $t_{1k} = t_{jk'}$. In the same way, by using the pair $\{L_2 \oplus L_j, L_k \oplus L_{k'} \oplus L_1\}$ (which is of type 3) we deduce that $t_{2k} = t_{jk'}$. It follows that all the phase factors are equal.

Case 2: $J = L_1 \oplus L_2$. Here $\dim K \geq 3$. Let $\psi_k, \psi_{k'}, \psi_{k''}$ be any 3 basis elements belonging to K . Then, by Theorem 4.2, $U_{L_1 \oplus L_2, K}$ agrees up to a phase factor with $U_{L_1 \oplus L_2 \oplus L_{k''}, L_k \oplus L_{k'}}$ and by Case 1 the latter agrees with U_B . But both $\psi_1 \otimes \psi_k$ and $\psi_2 \otimes \psi_{k'}$ lie in the domains of all these mappings so that $t_{1k} = t_{2k'}$, which (k and k' being arbitrary) establishes the equality of all the phase factors.

This completes the proof of Lemma 6 and with it of Theorem 5.6.

We can now separate the Bose and Fermi cases. Under the hypotheses of Theorem 5.6, let $E: H' \otimes H' \rightarrow H' \otimes H'$ be the "exchange operator" defined by $E(\psi \otimes \varphi) = \varphi \otimes \psi$ for all ψ and φ in H' . Let $V_B = U_B E$. Owing to the commutative property of the operator \circ , V_B has the same property as U_B : whenever $\{\psi, \varphi\}$ is compatible with B ,

$$\text{ray } V_B(\psi \otimes \varphi) = \text{ray } U_B(\varphi \otimes \psi) = \varphi \circ \psi = \psi \circ \varphi.$$

By Theorem 5.6 it follows that $V_B = \lambda_B U_B$, where λ_B is a phase factor. Since $\lambda_B^2 U_B = \lambda_B U_B E = U_B E^2 = U_B$, $\lambda_B = \pm 1$. We now prove that λ_B is independent of B .

Theorem 5.7: Let H' be a simple system of dimension ≥ 5 , B be an orthonormal set in H' of at least 3 elements, and λ_B be defined (as above) by the equation

$$U_B(\psi \otimes \varphi) = \lambda_B U_B(\varphi \otimes \psi),$$

valid for every orthogonal pair $\{\psi, \varphi\}$ of vectors in H' compatible with B . (U_B is defined in Theorem 5.6.) Then one of the following cases applies.

Bose case: $\lambda_B = 1$ for every such orthonormal set B .

Fermi case: $\lambda_B = -1$ for every such orthonormal set B .

Proof: We first note two lemmas.

Lemma 1: If $B' \subset B$, then $\lambda_{B'} = \lambda_B$.

For we can then take for $U_{B'}$ a suitable restriction of U_B . Lemma 1 gives at once.

Lemma 2: If $B \cap B'$ has at least 3 elements, then $\lambda_{B'} = \lambda_B$.

In proving the theorem it suffices, by Lemma 1, to consider only orthonormal sets consisting of 3 elements. Let B and B' be two such sets. We shall show that $\lambda_B = \lambda_{B'}$. To do this we proceed from B to an orthonormal set containing B' through a sequence of steps, each of which involves the addition or removal of a single unit vector. We start with $S = B$ and follow these instructions:

(i) If S has 3 elements or has 4 elements and is an orthonormal set, add to S an element of B' not already in S .

(ii) If S has 4 elements and is not an orthonormal set, add to S a unit vector orthogonal to each vector in S .

(iii) If S has 5 elements, delete from S any element which is not orthogonal to every member of the set $B' \cap S$.

This process terminates eventually, with $S \supset B'$ and the consequent impossibility of carrying out instruction (i). Writing down the finite sequence S_1, \dots, S_n obtained in this way and deleting from each S_r one element (or sometimes none), we get a sequence of orthonormal sets $B = B_1, \dots, B_n \supset B'$, to which we can apply Lemmas 1 and 2 and so deduce $\lambda_B = \lambda_{B'}$.

In virtue of Theorem 5.7, every simple system of dimension ≥ 5 is either a Fermi system or a Bose system in the following sense.

Definition 5.8: A Fermi (Bose) system is a simple system of dimension ≥ 5 for which the Fermi (Bose) case of Theorem 5.7 applies.

6. STRUCTURE OF FERMI AND BOSE SYSTEMS

In this section we determine, for a simple system H' of dimension ≥ 5 , the relation between H' and the Hilbert space H'' of 2-particle states. Our conclusion (Theorems 6.5 and 6.6) is that there is a natural isomorphism between H'' and (in the Fermi and Bose cases, respectively) the antisymmetric or symmetric tensor product of H' with itself. We define this isomorphism by means of the mappings U_B corresponding to various bases B . First, however, we must show that these mappings agree with each other up to a phase factor. The next theorem is an important step in this direction.

Throughout the following, H' denotes the Hilbert space of a simple system of dimension ≥ 5 .

Theorem 6.1: Let B and B' be bases of H' with at least 2 elements in common, and let $B'' = B \cap B'$. Let U_B and $U_{B'}$ be defined as in Theorem 5.6, with

the phase of $U_{B'}$ so adjusted that $U_{B'}(\xi \otimes \eta) = U_B(\xi \otimes \eta)$ for some orthogonal pair of vectors (ξ, η) compatible with B'' . Then, if $\{\psi, \varphi\}$ and $\{\psi', \varphi'\}$ are any orthogonal pairs of vectors compatible with B and B' , respectively,

$$\begin{aligned} \langle U_B(\psi \otimes \varphi), U_{B'}(\psi' \otimes \varphi') \rangle \\ = \langle \psi, \psi' \rangle \langle \varphi, \varphi' \rangle + \lambda \langle \psi, \varphi' \rangle \langle \varphi, \psi' \rangle, \quad (1) \end{aligned}$$

where $\lambda = \pm 1$ according as H' is a Bose or a Fermi system.

Before proving the theorem we establish two corollaries.

Corollary 6.2: If $B, B', U_B, U_{B'}$ satisfy the conditions stated in the theorem, and if $\{\psi, \varphi\}$ is any orthogonal pair of vectors compatible with both B and B' , then

$$U_B(\psi \otimes \varphi) = U_{B'}(\psi \otimes \varphi).$$

Proof: Assume ψ and φ are unit vectors. Then so are $U_B(\psi \otimes \varphi)$ and $U_{B'}(\psi \otimes \varphi)$, and by the theorem their scalar product is 1.

Corollary 6.3: If $\{\psi, \varphi\}$ and $\{\psi', \varphi'\}$ are any two pairs of orthogonal unit vectors in H' , then

$$\langle \Psi \circ \Phi, \Psi' \circ \Phi' \rangle = |\langle \psi, \psi' \rangle \langle \varphi, \varphi' \rangle + \lambda \langle \psi, \varphi' \rangle \langle \varphi, \psi' \rangle|.$$

Proof: Since $\dim H' \geq 5$, we can choose orthonormal vectors ψ_1, ψ_2, ψ_3 orthogonal to ψ and ψ' . Let B be a basis containing the set $\psi_1, \psi_2, \psi_3, \psi$ and let B' be a basis containing the set $\psi_1, \psi_2, \psi_3, \psi'$. Then $\{\psi, \varphi\}$ is compatible with B and $\{\psi', \varphi'\}$ is compatible with B' . Thus $\Psi \circ \Phi = \text{ray } U_B(\psi \otimes \varphi)$ and $\Psi' \circ \Phi' = \text{ray } U_{B'}(\psi' \otimes \varphi')$, and the result follows from the theorem.

Proof of Theorem 6.1: The proof is by a series of lemmas. Using the definition and linearity of U_B we easily obtain.

Lemma 1: The theorem holds in the case when $B = B'$.

Now let K be the linear space whose basis is B'' and let J be its orthogonal complement in H' .

Lemma 2: The theorem holds if $\psi, \psi', \varphi, \varphi'$ are all in K .

Proof: Let ζ be any unit vector in J and let $\tilde{B} = \{\zeta\} \cup B''$. Then $U_B|_{H'} \otimes_{\tilde{B}} H' = U_{B'}|_{H'} \otimes_{\tilde{B}} H'$ since these mappings, by Theorem 5.6, can differ only by a phase factor and they agree at $\xi \otimes \eta$. The result now follows from Lemma 1.

Lemma 3: The theorem holds whenever ψ and ψ' belong to J and φ and φ' belong to K . Moreover, in this case $U_B(\psi \otimes \varphi) = U_{B'}(\psi \otimes \varphi)$.

Proof: Putting $\zeta = \psi$, the argument used in the proof of Lemma 2 shows that $U_B(\psi \otimes \varphi) = U_{B'}(\psi \otimes \varphi)$. The result now follows from Lemma 1.

From Axiom 1.1 and Lemma 3 we easily deduce

Lemma 4: The theorem holds whenever, of the four vectors $\psi, \psi', \varphi, \varphi'$, two belong to J , and two belong to K .

Lemma 5: The theorem holds whenever $\psi, \psi', \varphi, \varphi'$ all belong to J .

Proof: Let $\langle \psi, \psi' \rangle = \alpha$, $\langle \psi, \varphi' \rangle = \beta$, $\langle \varphi, \psi' \rangle = \gamma$, $\langle \varphi, \varphi' \rangle = \delta$. Consider the four vectors $\psi + \xi, \varphi, \psi' - \alpha\xi + \beta\eta, \varphi' - \beta\xi - \bar{\alpha}\eta$. The first is orthogonal to the other three and the last two are orthogonal. Also the pair $\{\psi + \xi, \varphi\}$ is compatible with the basis B and the pair $\{\psi' - \alpha\xi + \beta\eta, \varphi' - \beta\xi - \bar{\alpha}\eta\}$ is compatible with \tilde{B} , where we denote by \tilde{B} the basis obtained by replacing in B' the vectors ξ, η by unit vectors ξ', η' parallel to $\alpha\xi - \beta\eta$ and $\beta\xi + \bar{\alpha}\eta$, respectively. Choose $U_{\tilde{B}}$ so that $U_{\tilde{B}}(\psi' \otimes \varphi') = U_{B'}(\psi' \otimes \varphi')$. We then have, by Axiom 1.1,

$$\begin{aligned} 0 &= \langle U_B[(\psi + \xi) \otimes \varphi], U_{\tilde{B}}[(\psi' - \alpha\xi + \beta\eta) \\ &\quad \otimes (\varphi' - \beta\xi - \bar{\alpha}\eta)] \rangle \\ &= \langle U_B(\psi \otimes \varphi + \xi \otimes \varphi), U_{\tilde{B}}(\psi' \otimes \varphi' - \alpha\xi \otimes \varphi' \\ &\quad + \beta\eta \otimes \varphi' - \beta\psi' \otimes \xi - \bar{\alpha}\psi' \otimes \eta \\ &\quad + [\alpha\xi - \beta\eta] \otimes [\beta\xi + \bar{\alpha}\eta]) \rangle. \end{aligned}$$

On expanding this product and using Axiom 1.1 again, all but three of the 12 terms vanish. We obtain

$$\begin{aligned} \langle U_B(\psi \otimes \varphi), U_{B'}(\psi' \otimes \varphi') \rangle \\ = \alpha \langle U_B(\xi \otimes \varphi), U_{B'}(\xi \otimes \varphi') \rangle \\ \quad + \beta \lambda \langle U_B(\xi \otimes \varphi), U_{B'}(\xi \otimes \psi') \rangle \\ = \alpha \delta + \beta \lambda \gamma, \end{aligned}$$

where Lemma 4 has been used on several occasions.

With the aid of Lemmas 2, 4, and 5, the theorem can now be proved. First note that both sides of (1) are linear in ψ' and φ' and antilinear in ψ and φ . It suffices therefore to consider only the 16 cases in which each of these four vectors lies in one or other of the spaces J and K . But in each of these cases either both sides of (1) vanish by Axiom 1.1 or the theorem holds by one of the lemmas. This completes the proof of Theorem 6.1.

Before studying the relation between H' and H'' for Fermi and Bose systems, we introduce some

notation. Let $B = \{\psi_i: i \in \mathfrak{J}\}$ be an orthonormal basis in 1-particle space H' . A basis for $H' \otimes H'$ is given by $\{\psi_i \otimes \psi_j: i \in \mathfrak{J}, j \in \mathfrak{J}\}$. The *antisymmetric tensor product* $H' \otimes_- H'$ is the (closed) subspace of $H' \otimes H'$ spanned by the orthonormal set $\{(\psi_i \otimes \psi_j - \psi_j \otimes \psi_i)/\sqrt{2}: i \in \mathfrak{J}, j \in \mathfrak{J}, i \neq j\}$ and the *symmetric tensor product* $H' \otimes_+ H'$ is that spanned by the orthonormal set

$$\{\psi_i \otimes \psi_i: i \in \mathfrak{J}\} \cup \{(\psi_i \otimes \psi_j + \psi_j \otimes \psi_i)/\sqrt{2}: i \in \mathfrak{J}, j \in \mathfrak{J}, i \neq j\}.$$

It is well known that these subspaces are independent of the basis B and are mutually orthogonal and that

$$H' \otimes H' = H' \otimes_- H' \oplus H' \otimes_+ H'.$$

We write also

$$H' \otimes_{B+} H' = H' \otimes_B H' \cap H' \otimes_+ H'.$$

We denote by E_- and E_+ the projections onto $H' \otimes_- H'$ and $H' \otimes_+ H'$, respectively.

Theorem 6.4: (i) For a Fermi system $U_B E_- = U_B$ and $U_B E_+ = 0$. Moreover, if $U_{B-} = (U_B | H' \otimes_- H')/\sqrt{2}$, then U_{B-} is an isometry of $H' \otimes_- H'$ into H'' .

(ii) For a Bose system $U_B E_+ = U_B$ and $U_B E_- = 0$. Moreover, if $U_{B+} = (U_B | H' \otimes_{B+} H')/\sqrt{2}$, then U_{B+} is an isometry of $H' \otimes_{B+} H'$ into H'' .

Proof: (i) Since U_{B-} is linear, it suffices to show that the orthonormal basis $\{(\psi_i \otimes \psi_j - \psi_j \otimes \psi_i)/\sqrt{2}: i \neq j\}$ for $H' \otimes_- H'$ is carried by U_{B-} into an orthonormal set in H'' .¹² But this fact is easily verified using the definition of U_{B-} and Theorem 6.1.

(ii) The proof is similar.

In the rest of this section we shall deal separately with the Fermi and Bose cases, the latter being considerably more complicated. In the Fermi (Bose) case our aim is to amalgamate the isometries U_{B-} (U_{B+}) into a single isometry U_- (U_+) from $H' \otimes_- H'$ onto H'' ($H' \otimes_+ H'$ onto H''). We treat first the Fermi case.

Theorem 6.5: For any Fermi system H' , there exists a linear isometry U_- , unique up to a phase factor, of $H' \otimes_- H'$ onto H'' such that, for any orthogonal unit vectors ψ, φ in H' ,

$$\text{ray } U_-[(\psi \otimes \varphi - \varphi \otimes \psi)/\sqrt{2}] = \psi \circ \varphi.$$

Proof: We first establish uniqueness. Assume U_- and \tilde{U}_- both have the stated property. Define a bounded linear operator $U'_B: H' \otimes_B H' \rightarrow H$ by setting, whenever $i \neq j$,

$$U'_B(\psi_i \otimes \psi_j) = U_-(\psi_i \otimes \psi_j - \psi_j \otimes \psi_i)/\sqrt{2},$$

and define \tilde{U}'_B similarly. Then both U'_B and \tilde{U}'_B satisfy the condition required of U_B in Theorem 5.6. Consequently, they can differ only by a phase factor. However, $U_- = U'_B | H' \otimes_- H')/\sqrt{2}$ and similarly for \tilde{U}_- . So U_- and \tilde{U}_- differ only by a phase factor.

We now establish the existence of U_- . Choose arbitrarily a basis B for H' and let U_B and U_{B-} be defined as in Theorems 5.6 and 6.4. Put $U_- = U_{B-}$. We must show that for any orthogonal unit vectors ψ, φ in H' ,

$$\text{ray } U_-[(\psi \otimes \varphi - \varphi \otimes \psi)/\sqrt{2}] = \psi \circ \varphi; \quad (2)$$

i.e., that every pair of orthogonal unit vectors in H' belongs to Ψ' , where Ψ' denotes the set of all orthogonal pairs of unit vectors $\{\psi, \varphi\}$ in H' such that (2) holds.

By Theorems 5.6 and 6.4 we have at once

Lemma 1: Any pair of orthogonal unit vectors compatible with B belongs to Ψ' .

Now let $\psi = \sum_i a_i \psi_i$ and $\varphi = \sum_i b_i \psi_i$ be any pair of orthogonal unit vectors in H' . We shall show that $\{\psi, \varphi\} \in \Psi'$.

Choose on the ray $\psi \circ \varphi$ a unit vector $\sum_{i \neq j} c_{ij} \theta_{ij} + \bar{\theta}$ where $\theta_{ij} = -\theta_{ji} = U_B(\psi_i \otimes \psi_j)$. We may without loss of generality assume $c_{ij} = -c_{ji}$ and $\bar{\theta} \perp \theta_{ij}$ for all i, j with $i \neq j$. By Corollary 6.4, for any orthogonal unit vectors ψ', φ' in H'

$$\begin{aligned} \langle \psi' \circ \varphi', \psi \circ \varphi \rangle &= |\langle \psi', \psi \rangle \langle \varphi', \varphi \rangle - \langle \psi', \varphi \rangle \langle \varphi', \psi \rangle|. \end{aligned} \quad (3)$$

Putting $\psi' = \psi_i$ and $\varphi' = \psi_j$, we obtain

$$|c_{ij} - c_{ji}| = |a_i b_j - a_j b_i|.$$

Thus

$$2c_{ij} = t_{ij}(a_i b_j - a_j b_i),$$

where $t_{ij} = t_{ji}$ is a number of modulus 1 defined for each i, j with $i \neq j$ and $c_{ij} \neq 0$. A short calculation shows that $\sum_{i \neq j} c_{ij} \theta_{ij}$ is a unit vector, so that $\bar{\theta} = 0$.

We next show that t_{ij} is independent of i and j . First,

Lemma 2: If i, j, k are distinct, $c_{ij} \neq 0$, and $c_{ik} \neq 0$, then $t_{ij} = t_{ik}$.

Proof: In (3) put $\psi' = \psi_i$ and $\varphi' = \alpha \psi_j + \beta \psi_k$, where $|\alpha|^2 + |\beta|^2 = 1$, $\alpha \neq 0$, $\beta \neq 0$. We obtain

$$|\bar{\alpha} c_{ij} + \bar{\beta} c_{ik}| = |\bar{\alpha} c_{ij}/t_{ij} + \bar{\beta} c_{ik}/t_{ik}|$$

so that

$$|1 + \bar{\beta} c_{ik}/\bar{\alpha} c_{ij}| = |1 + \bar{\beta} c_{ik} t_{ij}/\bar{\alpha} c_{ij} t_{ik}|.$$

Now α and β may be chosen so that $\bar{\beta} c_{ik}/\bar{\alpha} c_{ij}$ is real. It follows that $t_{ij}/t_{ik} = 1$.

Lemma 3: If $c_{ij} \neq 0$ and $c_{mn} \neq 0$, then $t_{ij} = t_{mn}$.

Proof: Assume $c_{ij}c_{mn} \neq 0$. The identity

$$(a_i b_j - a_j b_i)(a_m b_n - a_n b_m) + (a_i b_m - a_m b_i) \times (a_n b_j - a_j b_n) + (a_i b_n - a_n b_i)(a_j b_m - a_m b_j) = 0$$

then shows that either $c_{im} \neq 0$ or $c_{in} \neq 0$. Assume $c_{im} \neq 0$; then by two applications of Lemma 2 we have $t_{ij} = t_{im} = t_{mi} = t_{mn}$.

From Lemma 3 it follows that all the phase factors t_{ij} are equal so that, by means of a phase change in θ , we may obtain

$$\theta = \sum_{i \neq j} (a_i b_j - a_j b_i) \theta_{ij} / 2.$$

But this means $\theta = U_-(\psi \otimes \varphi - \varphi \otimes \psi) \sqrt{2}$, which shows that the pair $\{\psi, \varphi\}$ satisfies (2), completing the proof of the theorem.

The proof of the corresponding theorem for the Bose case is rather more complicated:

Theorem 6.6: For any Bose system there exists a linear isometry U_+ , unique up to a phase factor, of $H' \otimes_+ H'$ onto H'' such that, for any orthogonal unit vectors ψ, φ in H' ,

$$\text{ray } U_+[(\psi \otimes \varphi + \varphi \otimes \psi) / \sqrt{2}] = \psi \circ \varphi.$$

Proof: We first establish uniqueness. Let U_+ be an arbitrary isometry satisfying the condition of the theorem. Define a linear mapping $U: H' \otimes H' \rightarrow H''$ by setting

$$U(\psi \otimes \varphi) = U_+[(\psi \otimes \varphi + \varphi \otimes \psi) / \sqrt{2}].$$

Let $B = \{\psi_i : i \in \mathbb{J}\}$ be a fixed basis in H' and let $U_B: H' \otimes_B H' \rightarrow H''$ be a fixed mapping satisfying the condition of Theorem 5.6. Since $U|_{H' \otimes_B H'}$ also satisfies this condition, it differs from U_B only by a phase factor; thus by altering U_+ by a phase factor we can obtain $U|_{H' \otimes_B H'} = U_B$. We assume this done. We then have, whenever $i \neq j$,

$$U_+[(\psi_i \otimes \psi_j + \psi_j \otimes \psi_i) / \sqrt{2}] = U_B(\psi_i \otimes \psi_j).$$

We have to show that U_+ is now uniquely fixed. This will be accomplished if we can show that the value of $U_+(\psi_i \otimes \psi_i)$ is independent of the original choice of U_+ .

To do this we introduce some more orthonormal bases in H' . For each i, j with $i \neq j$, let B_{ij} be the basis obtained from B by replacing ψ_i and ψ_j by $(\psi_i + \psi_j) / \sqrt{2}$ and $(\psi_i - \psi_j) / \sqrt{2}$, and let B'_{ij} be that obtained from B by replacing ψ_i and ψ_j by $(\psi_i + i\psi_j) / \sqrt{2}$ and $(\psi_i - i\psi_j) / \sqrt{2}$. Let the phase of $U_{B_{ij}}$ and $U_{B'_{ij}}$ be chosen so that these operators agree with U_B .

Then $U_{B_{ij}}$ and $U_{B'_{ij}}$ are uniquely determined. Moreover, $U|_{H' \otimes_{B_{ij}} H'} = U_{B_{ij}}$, since these can differ only by a phase factor and both agree with U_B at, for instance, the point $\psi_m \otimes \psi_n$, where i, j, m, n are all distinct. Similarly, $U|_{H' \otimes_{B'_{ij}} H'} = U_{B'_{ij}}$. Now

$$\begin{aligned} 2\sqrt{2} U_+(\psi_i \otimes \psi_i) &= 2U(\psi_i \otimes \psi_i) \\ &= U[(\psi_i + \psi_j) \otimes (\psi_i - \psi_j) \\ &\quad + (\psi_i + i\psi_j) \otimes (\psi_i - i\psi_j)] \\ &= U_{B_{ij}}[(\psi_i + \psi_j) \otimes (\psi_i - \psi_j)] \\ &\quad + U_{B'_{ij}}[(\psi_i + i\psi_j) \otimes (\psi_i - i\psi_j)] \end{aligned}$$

which is independent of the choice of U_+ . This completes the proof of the uniqueness up to a phase factor of U_+ .

To establish the existence of a mapping U_+ with the desired property, we define it as follows. First choose a basis B in H' and a mapping $U_B: H' \otimes_B H' \rightarrow H''$; then define bases B_{ij} and B'_{ij} and fix the phases of $U_{B_{ij}}$ and $U_{B'_{ij}}$ all as above. Now, for $i \neq j$, let

$$\theta_{ij} = U_B(\psi_i \otimes \psi_j)$$

and let

$$\begin{aligned} 2\sqrt{2} \theta_{iij} &= U_{B_{ij}}[(\psi_i + \psi_j) \otimes (\psi_i - \psi_j)] \\ &\quad + U_{B'_{ij}}[(\psi_i + i\psi_j) \otimes (\psi_i - i\psi_j)]. \end{aligned}$$

Lemma 1: θ_{iij} is a unit vector which is independent of j .

Proof: Let i, j, k be indices with $i \neq j$ and $i \neq k$. Since $\dim H' \geq 5$, any pair of the bases $B, B_{ij}, B'_{ij}, B_{ik}, B'_{ik}$ has at least 2 elements in common. We can therefore apply Theorem 6.1 in determining $\langle \theta_{iij}, \theta_{iik} \rangle$. A short calculation gives $\langle \theta_{iij}, \theta_{iik} \rangle = 1$.

Let us write $\theta_{iij} = \theta_{ii}$. In the same way we can prove

Lemma 2: For any indices i, j, k ,

$$\begin{aligned} \langle \theta_{ii}, \theta_{jj} \rangle &= \delta_{ij}, \\ \langle \theta_{ii}, \theta_{jk} \rangle &= 0. \end{aligned}$$

These equations show that if we define a linear mapping $U_+: H' \otimes_+ H' \rightarrow H''$ by

$$\begin{aligned} U_+(\psi_i \otimes \psi_j + \psi_j \otimes \psi_i) / \sqrt{2} &= \theta_{ij} \quad \text{for } i \neq j, \\ U_+(\psi_i \otimes \psi_i) &= \theta_{ii}, \end{aligned}$$

then U_+ is an isometry.

We shall show that U_+ has the required property. Let Ψ denote the set of all pairs $\{\psi, \varphi\}$ of orthogonal

unit vectors in H' for which

$$\text{ray } U_+(\psi \otimes \varphi + \varphi \otimes \psi)/\sqrt{2} = \psi \circ \varphi.$$

Lemma 3: Ψ is a closed set in the set product $H' \times H'$ with the product topology.

Proof: Ray $U_+(\psi \otimes \varphi + \varphi \otimes \psi)/\sqrt{2}$ is a continuous function of the pair $\{\psi, \varphi\}$ and so, by Corollary 6.3, is $\psi \circ \varphi$.

It is convenient to define a linear mapping $U: H' \otimes H' \rightarrow H''$ by

$$U(\psi \otimes \varphi) = U_+(\psi \otimes \varphi + \varphi \otimes \psi)/\sqrt{2}.$$

Lemma 4: Let i, j be distinct indices. If $\{\psi, \varphi\}$ is a pair of orthogonal unit vectors compatible with any of the bases B, B_{ij}, B'_{ij} , then $\{\psi, \varphi\} \in \Psi'$.

Proof: Let i, j, m, n be distinct elements of J . Then

$$(a) \ U(\psi_m \otimes \varphi_n) = \theta_{mn} = U_B(\psi_m \otimes \varphi_n),$$

$$(b) \ U[(\psi_i \pm \psi_j) \otimes \varphi_m] = \theta_{im} \pm \theta_{jm} \\ = U_B[(\psi_i \pm \psi_j) \otimes \varphi_m] \\ = U_{B_{ij}}[(\psi_i \pm \psi_j) \otimes \varphi_m],$$

and similarly

$$(c) \ U[(\psi_i \pm i\psi_j) \otimes \varphi_m] = U_{B'_{ij}}[(\psi_i \pm i\psi_j) \otimes \varphi_m].$$

Moreover, since $U(\psi_i \otimes \varphi_i) = \sqrt{2}\theta_{ii}$,

$$(d) \ U[(\psi_i + \psi_j) \otimes (\psi_i - \psi_j)] = \sqrt{2}\theta_{ii} - \sqrt{2}\theta_{jj} \\ = \sqrt{2}(\theta_{iii} - \theta_{jjj}) \\ = U_{B_{ij}}[(\psi_i + \psi_j) \\ \otimes (\psi_i - \psi_j)],$$

and similarly

$$(e) \ U[(\psi_i + i\psi_j) \otimes (\psi_i + i\psi_j)] \\ = U_{B'_{ij}}[(\psi_i + i\psi_j) \otimes (\psi_i - i\psi_j)].$$

By (a) U agrees with U_B on the elements of a basis of the domain of U_B and consequently everywhere in that domain. But, whenever $\{\psi, \varphi\}$ is compatible with B , ray $U_B(\psi \otimes \varphi) = \psi \circ \varphi$ so that $\{\psi, \varphi\} \in \Psi'$.

In the same way it follows from (b), (c), (d), and (e) that $\{\psi, \varphi\} \in \Psi'$ whenever $\{\psi, \varphi\}$ is compatible with B_{ij} or with B'_{ij} for any distinct indices i, j . This completes the proof of Lemma 4.

Lemma 5: Let i, j be distinct indices and let $\psi = a_i\psi_i + a_j\psi_j$ and $\varphi = b_i\psi_i + b_j\psi_j$ be orthogonal unit vectors. Then $\{\psi, \varphi\} \in \Psi'$.

Proof: In view of the definition of U_+ , we have to show that

$$\sqrt{2} a_i b_i \theta_{ii} + (a_i b_j + a_j b_i) \theta_{ij} + \sqrt{2} a_j b_j \theta_{jj} = \theta,$$

say, is a vector on the ray $\psi \circ \varphi$.

Let B'' be the base obtained from B by replacing ψ_i and ψ_j by ψ and φ , and choose the phase of $U_{B''}$ so that $U_{B''}$ agrees with U_B . We shall show that $U_{B''}(\psi \otimes \varphi) = \theta$. Let

$$U_{B''}(\psi \otimes \varphi) = c_{ii}\theta_{ii} + c_{ij}\theta_{ij} + c_{jj}\theta_{jj} + \bar{\theta}$$

where $\bar{\theta}$ is orthogonal to θ_{ii} , θ_{ij} , and θ_{jj} . Since $\dim H' \geq 5$, each of the sets $B \cap B''$, $B_{ij} \cap B''$, and $B'_{ij} \cap B''$ has at least three elements. Hence, by Theorem 6.1,

$$\langle U_{B''}(\psi \otimes \varphi), U_B(\psi_i \otimes \psi_j) \rangle = \langle \psi, \psi_i \rangle \langle \varphi, \psi_j \rangle \\ + \langle \psi, \psi_j \rangle \langle \varphi, \psi_i \rangle,$$

i.e.,

$$c_{ij} = a_i b_j + a_j b_i.$$

Similarly, by using in turn $U_{B_{ij}}[(\psi_i + \psi_j) \otimes (\psi_i - \psi_j)]$ and $U_{B'_{ij}}[(\psi_i + i\psi_j) \otimes (\psi_i - i\psi_j)]$ instead of $U_B(\psi_i \otimes \psi_j)$, we obtain

$$c_{ii} - c_{jj} = \sqrt{2}(a_i b_i - a_j b_j)$$

and

$$c_{ii} + c_{jj} = \sqrt{2}(a_i b_i + a_j b_j).$$

Thus $U_{B''}(\psi \otimes \varphi) = \theta + \bar{\theta}$. Since both θ and $\theta + \bar{\theta}$ are unit vectors, $\bar{\theta} = 0$.

Lemma 6: Let $\psi = \sum_i a_i \psi_i$ and $\varphi = \sum_i b_i \psi_i$ be any orthogonal unit vectors in H' such that the numbers $a_i b_i$ are all different and nonzero. Then $\{\psi, \varphi\} \in \Psi'$.

Proof: Let a unit vector on the ray $\psi \circ \varphi$ be

$$\theta = \sum_i c_{ii} \theta_{ii} + \sum_{i \neq j} c_{ij} \theta_{ij} + \bar{\theta},$$

where $c_{ij} = c_{ji}$ and $\bar{\theta}$ is orthogonal to the vectors θ_{ij} , for all i and j . We shall prove the lemma by showing that, with a suitable choice of phase, $\theta = U_+(\psi \otimes \varphi)/\sqrt{2}$.

For any indices m, n with $m \neq n$, let

$$\psi' = \psi_m + \epsilon \psi_n,$$

$$\varphi' = \bar{\epsilon} \psi_m - \psi_n,$$

where ϵ is an arbitrary complex number. By Lemmas 2 and 5,

$$\theta' = U_+(\psi' \otimes \varphi' + \varphi' \otimes \psi')/\sqrt{2} \\ = \sqrt{2}(\epsilon \theta_{mm} - \epsilon \theta_{nn}) + (|\epsilon|^2 - 1)\theta_{mn}$$

is a vector of length $1 + |\epsilon|^2$ on $\psi' \circ \varphi'$. Hence, by Corollary 6.3,

$$|\langle \theta', \theta \rangle| = |\langle \psi', \psi \rangle \langle \varphi', \varphi \rangle + \langle \psi', \varphi \rangle \langle \varphi', \psi \rangle|,$$

whence

$$\begin{aligned} & |\sqrt{2}(\bar{\epsilon}c_{mm} - \epsilon c_{nn}) + 2(|\epsilon|^2 - 1)c_{mn}| \\ &= |2(\bar{\epsilon}a_m b_m - \epsilon a_n b_n) + (|\epsilon|^2 - 1)(a_m b_n + a_n b_m)|. \end{aligned} \quad (4)$$

From this equation (which holds for every ϵ and every pair m, n of distinct indices), we deduce the value of the constants c_{ij} . First put $|\epsilon| = 1$ and let $\mu = \bar{\epsilon}/\epsilon$, an arbitrary complex number of modulus unity. We have, for every such μ ,

$$|\mu c_{mm} - c_{nn}| = \sqrt{2} |\mu a_m b_m - a_n b_n|. \quad (5)$$

Allowing μ to vary and taking the maximum of each side of (5) gives

$$|c_{mm}| + |c_{nn}| = |\sqrt{2} a_m b_m| + |\sqrt{2} a_n b_n|.$$

Let p be distinct from m and n . Writing the previous equation also for the pairs m, p and n, p , we obtain $|c_{mm}| = \sqrt{2} |a_m b_m|$ so that

$$c_{mm} = t_m \sqrt{2} a_m b_m,$$

where t_m is a phase factor which may depend on m . Substituting in (5) and observing that both sides of (5) must attain their maximum values for the same *single* value of μ (since $0 \neq c_{mm} \neq c_{nn} \neq 0$), we obtain $t_m = t_n$. Thus all these phase factors are equal. Let us assume that the phase of θ is so chosen that for every index i , $c_{ii} = \sqrt{2} a_i b_i$.

Now put $\epsilon = 0$. Then (4) gives

$$2c_{mn} = t_{mn}(a_m b_n + a_n b_m),$$

where t_{mn} is a phase factor. We shall show that, for every distinct pair of indices m, n , $t_{mn} = 1$. To do this let $\epsilon = \sqrt{2} \mu$, where μ is an arbitrary phase factor. Equation (4) then yields

$$\begin{aligned} & |(\bar{\mu}c_{mm} - \mu c_{nn}) + c_{mn}| \\ &= |(\bar{\mu}c_{mm} - \mu c_{nn}) + t_{mn}c_{mn}|. \end{aligned} \quad (6)$$

Now, as μ varies, the point $\bar{\mu}c_{mm} - \mu c_{nn}$ describes a (nondegenerate) ellipse in the complex plane. By (6) the points $-c_{mn}$ and $-t_{mn}c_{mn}$ are equidistant from every point on this ellipse. They must therefore coincide. Unless $c_{mn} = 0$ (in which case there is nothing to prove), this means $t_{mn} = 1$.

Substituting the values thus found for the constants c_{ii} and c_{ij} , we have

$$\theta = \sum_i \sqrt{2} a_i b_i \theta_{ii} + \sum_{i \neq j} [(a_i b_j + a_j b_i)/\sqrt{2}] \theta_{ij} + \bar{\theta};$$

i.e.,

$$\theta = U_+(\psi \otimes \varphi + \varphi \otimes \psi)/\sqrt{2} + \bar{\theta}.$$

Now θ is a unit vector and $\bar{\theta}$ is orthogonal to the unit vector $U_+(\psi \otimes \varphi + \varphi \otimes \psi)/\sqrt{2}$; so $\bar{\theta} = 0$. This completes the proof of Lemma 6.

From Lemmas 6 and 3 it follows that every orthogonal pair of unit vectors in H' belongs to Ψ , which proves Theorem 6.6.

Transmission of Electromagnetic Waves through a Conducting Slab. I. The Two-Sided Wiener-Hopf Solution

G. A. BARAFF

Bell Telephone Laboratories, Incorporated, Murray Hill, New Jersey

(Received 16 July 1967)

This is the first of a series of papers dealing with propagation of electromagnetic waves through a metallic slab of finite thickness. In this first paper, we present a method for solving the integrodifferential equation governing the electric field in the interior of the metal when the electrons in the interior of the metal suffer diffuse reflection at each surface. The method is potentially of use in a wide class of problems, namely, the finite-slab generalization of all those semi-infinite-medium problems which are conventionally studied by the Wiener-Hopf technique. The solution given here is an iterative one with successive terms converging as $e^{-L/l}$, where L is the thickness of the slab and l is the range of the kernel of the integral term in the equation.

1. INTRODUCTION

Electromagnetic waves propagating through metals, semimetals, and semiconductors have provided an exceptionally useful tool for the study of solid state plasmas,¹ of electron-electron interactions,² and of the various excitations which the solid can support, such as phonons³ and magnons.⁴ In most cases, the analysis of an experiment or the predictions of a theory are carried out using the infinite-medium dispersion relationship. This is as it should be, for the next stage of complexity, solving the boundary-value problem, is often orders of magnitude more complex and often provides answers which are really no more informative than those the dispersion relation could have provided. Nonetheless, there are cases, such as the anomalous skin effect,⁵ field normal cyclotron resonance⁶ and size effect,⁷ and helicon propagation near doppler-shifted cyclotron resonance,⁸ where the presence of the boundaries plays a critical role in the phenomena of interest. In the first three of these, there is no infinite-medium dispersion relationship to speak of, while in the last, we are in a regime where the dispersion relationship is losing validity.

Some of the boundary-value problems associated with the above phenomena have already been solved. In order of increasing complexity, solutions have been given for the fields in a semi-infinite medium, assuming specular reflection of the electrons from the

boundary,^{5,7-9} for a finite slab assuming specular reflection¹⁰; and for the semi-infinite medium, assuming diffuse reflection of the electrons from the boundary.^{5,6,11a} The purpose of the present effort is to determine the field in a finite slab, assuming diffuse reflection of the electrons, and to compare the solution with that for the finite slab, assuming specular reflection of the electrons.

This first paper deals exclusively with how the fields are calculated and not at all with the results of that calculation, which is the subject of the second and third papers.^{11b} In this paper, we present a two-sided Wiener-Hopf method for solving the integrodifferential equation which governs the electric field in the interior of the metal. The equation is of the form

$$\frac{d^2}{dz^2} e(z) + Ae(z) = \int_0^L K(|z - z'|)e(z') dz',$$

$$0 \leq z \leq L,$$

where $e(z)$, the unknown electric field, satisfies boundary conditions at $z = 0$ and at $z = L$, A is a given constant, L is the width of the slab, and K is a given kernel. In the electromagnetic-wave problem, $A = \omega^2/c^2$ and $K = -i\omega\mu_0\sigma(|z - z'|)$, where σ is the nonlocal conductivity kernel, but the method presented here makes use only of one special property of K , namely that $|K|$ decays exponentially at infinity like $\exp -|z|/l$, where l is some positive length. (l is the electron mean free path in the electromagnetic-wave problem.) Thus, the method presented here is a general one and is potentially of use in a wide variety of problems, namely, the finite slab generalizations of all those semi-infinite-medium problems which are

¹ *Proceedings of the Symposium on Plasma Effects in Solids, Paris, 1964* (Dunod Cie., Paris, 1965) and references cited therein.

² W. M. Walsh and P. M. Platzman, *Phys. Rev. Letters* **9**, 514 (1967); S. Shultz and G. Dunifer, *Phys. Rev. Letters* **18**, 283 (1967).

³ C. C. Grimes and S. J. Buchsbaum, *Phys. Rev. Letters* **12**, 357 (1964); V. G. Skobov and E. A. Kaner, *Soviet Phys. JETP* **19**, 189 (1964).

⁴ F. A. Stern and E. R. Callen, *Phys. Rev.* **131**, 512 (1963); C. C. Grimes, *Bull. Am. Phys. Soc.* **10**, 471 (1965).

⁵ G. E. H. Reuter and E. H. Sondheimer, *Proc. Roy. Soc. (London)* **A195**, 336 (1948).

⁶ R. G. Chambers, *Phil. Mag.* **1**, 459 (1956).

⁷ M. Ya Azbel and E. A. Kaner, *Soviet Phys.—JETP* **5**, 730 (1957).

⁸ P. B. Miller and R. R. Haering, *Phys. Rev.* **128**, 126 (1962).

⁹ J. C. McGroddy, J. L. Stanford, and E. A. Stern, *Phys. Rev.* **141**, 437 (1966).

¹⁰ P. M. Platzman and S. J. Buchsbaum, *Phys. Rev.* **132**, 2 (1963).

^{11a} R. B. Dingle, *Physica* **9**, 311 (1953).

^{11b} G. A. Baraff, *Phys. Rev.* **167**, 625 (1968); the third will also appear in *Phys. Rev.*

conventionally studied by the Wiener–Hopf technique.¹² The reader interested *solely* in the physics is invited to read through the remainder of this section and the next, and then to await publication of the second paper: The physical content of this first paper is purposely restricted to shorten the presentation of the method.

The method presented here is a Wiener–Hopf method in that the basic unknown, the Fourier transform of the field, can be found only after certain analytic properties of the transform are utilized; without exploiting these analytic properties, there is not sufficient information to determine the transform. The essence of the present approach is that the field in the slab can be expressed as a sum of two functions, one decaying to the right and one decaying to the left. These functions must be defined outside the slab in such a way that (even though they do not represent the field exterior to the slab) their transforms have singularities in the finite k plane. Without this, there would be no useful analytic properties to exploit. In particular, the functions cannot vanish identically outside the slab. We therefore devise an extension of the original equation to determine these functions over all space. Having done this, we take the transforms of the extended equation and find that the two transforms will have the correct analytic properties only if an auxiliary condition is satisfied. This auxiliary condition is identical to the standard inhomogeneous Hilbert problem on an arc, and the machinery for satisfying it exists—the book by Muskhelishvili¹³ is the source of our knowledge here. Applying the Muskhelishvili techniques (which we describe without rigor) leads to an iterative solution. The zeroth term contains the standard Wiener–Hopf solution and also contains the Fabry–Perot fringes arising from multiple internal reflection of the wave. Higher iterates converge at least like $e^{-L/l}$. The higher iterates describe multiple reflection of the single-particle excitations, and so the convergence parameter for the series is really the amplitude of the single-particle excitations at the far side of the slab. This can be considerably smaller than $e^{-L/l}$, and the zeroth iterate can provide an accurate solution. The discussion of these matters is, however, part of the physics which will be found in the second and third papers.^{11b}

2. THE INTEGRODIFFERENTIAL EQUATION

For the purposes of this study of electromagnetic waves propagating along a magnetic field normal to

¹² E. Hopf, *Mathematical Problems of Radiative Equilibrium* (Cambridge University Press, New York, 1934).

¹³ N. I. Muskhelishvili, *Singular Integral Equations* (P. Noordhoff Ltd., Groningen, The Netherlands, 1953).

the faces of a metallic slab, it is sufficient to characterize the metal as a uniform free electron gas immersed in a uniform fixed background of neutralizing positive charge extending from $z = 0$ to $z = L$.¹⁴ One can, with no loss of generality, take the electric currents and electromagnetic fields to be transverse, circularly polarized, and monochromatic with a time dependence $e^{-i\omega t}$. As a result of the nearly random velocities of the electrons, the current j at plane z , time t , will be carried by electrons which are influenced by the electric field e at other planes z' at earlier times t' . The relation between the current and field is thus a nonlocal one:

$$\begin{aligned} j_{\pm}(z) &= \int_0^L \sigma_{\pm}(z, z', \omega) e_{\pm}(z') dz', \quad 0 \leq z \leq L, \\ j_{\pm}(z) &= J_x(z) \pm iJ_y(z), \\ e_{\pm}(z) &= E_x(z) \pm iE_y(z). \end{aligned} \quad (2.1)$$

This relationship, used to eliminate the current from Maxwell's equation, leads in the usual way to an integrodifferential equation for the field:

$$\left(\frac{d^2}{dz^2} + k_0^2 \right) e_{\pm}(z) = -i\omega\mu_0 \int_0^L dz' \sigma_{\pm}(z, z', \omega) e_{\pm}(z'), \quad (2.2a)$$

$$k_0^2 = \omega^2/c^2. \quad (2.2b)$$

In order to calculate the conductivity kernel σ_{\pm} , we use the Boltzmann equation in the simplified collision-time approximation commonly used in transport studies:

$$\begin{aligned} \left(\frac{\partial}{\partial t} + V_z \frac{\partial}{\partial z} \right) f(z, \mathbf{p}, t) + q(\mathbf{e}_{\pm} + \mathbf{V} \times \mathbf{B}) \cdot \nabla_{\mathbf{p}} f \\ = -(f - f_0)/\tau, \end{aligned} \quad (2.3)$$

where $f_0(p)$ is the uniform time-independent equilibrium distribution function. The equation is to be linearized by putting

$$\begin{aligned} f(z, \mathbf{p}, t) &= f_0(\mathbf{p}) + f_1(z, \mathbf{p})e^{-i\omega t}, \\ \mathbf{B}(z, t) &= \mathbf{B}_0 + \mathbf{b}(z)e^{-i\omega t}, \end{aligned}$$

and solved for f_1 .

The boundary conditions on f of interest here are those corresponding to diffuse scattering: An electron at the surface $z = 0$ traveling towards the interior of the metal must have just previously collided with that surface. Its distribution function must be f_0 , the distribution in the absence of fields, because there has not been sufficient time for the field to alter the motion of the electron. Similar considerations apply

¹⁴ The range of validity of this model is delineated in Ref. 10.

to the boundary conditions at $z = L$. Thus, the boundary condition to be applied to (2.3) is

$$\begin{aligned} f_1(z = 0, \mathbf{p}) &= 0 \quad \text{for } V_z > 0, \\ f_1(z = L, \mathbf{p}) &= 0 \quad \text{for } V_z < 0. \end{aligned} \quad (2.4)$$

The method of solving (2.3) and calculating the conductivity kernel may be inferred from a study of the Reuter-Sondheimer paper.⁵ The magnetic field B_0 present in (2.3) causes no complication and the result for a free-electron model, after all the integrals are done, is exactly the Reuter-Sondheimer result with ω replaced by $\omega \pm \omega_c$, where ω_c is the cyclotron frequency.¹⁵

The important feature of this result is that *in spite of the boundaries at $z = 0$ and $z = L$* , the boundary condition (2.4) results in a conductivity tensor which depends on the two spatial arguments *only through their difference*. That is, in the metal,

$$\sigma_{\pm}(z, z', \omega) = \sigma_{\pm}(|z - z'|, \omega), \quad \begin{aligned} 0 \leq z \leq L \\ 0 \leq z' \leq L, \end{aligned} \quad (2.5)$$

where the kernel on the right is the ordinary infinite-medium conductivity tensor. Thus, the integrodifferential equation to be solved (with the polarization index \pm and the argument ω suppressed) is of the form

$$\left(\frac{d^2}{dz^2} + k_0^2\right)e(z) = -i\omega\mu_0 \int_0^L dz' \sigma(|z - z'|)e(z'), \quad 0 \leq z \leq L. \quad (2.6a)$$

The boundary condition is that the field be incident on the slab with unit amplitude from the left:

$$\left(1 + \frac{1}{ik_0} \frac{d}{dz}\right)e(z) = 2 \quad \text{at } z = 0, \quad (2.6b)$$

$$\left(1 - \frac{1}{ik_0} \frac{d}{dz}\right)e(z) = 0 \quad \text{at } z = L. \quad (2.6c)$$

At this point, we may note that the specular reflection boundary condition used by Platzman and Buchsbaum¹⁰ (electrons striking the boundaries suffer a reversal of the z component of velocity but no change in x or y component) *does not lead* to a displacement kernel of the form (2.5). It leads instead to a form which can be *reexpressed* as (2.5), provided that:

(a) The limits 0 and L of the integral are made infinite.

(b) The field $e(z)$ is made periodic with period $2L$, and even.

¹⁵ Ref. 5, Eq. (14).

Proviso (b) demands the introduction of periodic current sheets. Proviso (a) then renders the problem amenable to straightforward solution by Fourier transforms which, because of proviso (b), means by Fourier cosine series. Platzman and Buchsbaum arrived correctly at these conclusions and were able to evaluate their series solution rather simply. In our case, the finite limits 0 and L are responsible for the complexity of the problem which does not yield directly to Fourier analysis.

3. THE TWO-SIDED EXTENSION OF THE ORIGINAL EQUATION

If the slab were infinitely thick, that is, if L were infinite, then Eqs. (2.6) would be exactly the semi-infinite-medium problem solved by Reuter and Sondheimer using the Wiener-Hopf technique. Recall that the starting point is the observation that the integral equation (2.6a) determines $e(z)$ only for $z > 0$, so that there is complete freedom in the definition of $e(z)$ for $z < 0$. In the Wiener-Hopf method, the choice is made to define $e(z)$ as zero for negative z . Having made that choice, the integral equation is no longer satisfied for $z < 0$ because, although the left side of (2.6a) has been defined as zero, the right side does not vanish. Hence, the first step is to define a compensating function, call it $h(z)$, which, when added to the integral equation, results in an equation which is true for all z . The conventional choice would be

$$\begin{aligned} h(z) &= +i\omega\mu_0 \int_0^{\infty} dz' \sigma(|z - z'|)e(z'), \quad z < 0, \\ &= 0, \quad z \geq 0, \end{aligned}$$

because this, added to the right side of (2.6a) (with $L = \infty$), would extend the validity of (2.6a) to all z . Then, the Fourier transform of the augmented equation

$$\left(\frac{d^2}{dz^2} + k_0^2\right)e(z) = -i\omega\mu_0 \int_{-\infty}^{\infty} dz' \sigma(|z - z'|)e(z') + h(z)$$

can be taken and solved for the two transforms $E(k)$ and $H(k)$. The reason that one equation can be solved for two unknowns, E and H , is that by virtue of $e(z)$ being defined as zero for $z < 0$, all singularities of $E(k)$ will be located in one half of the complex k plane, while all singularities of $H(k)$, the transform of an $h(z)$ defined as zero for $z > 0$, will be located in the other half plane. This information about the analyticity of E and H underlies the Wiener-Hopf method.

Now, we should like to apply the same sort of ideas to (2.6a) with L finite. First we observe that, since (2.6) defines $e(z)$ only for $0 \leq z \leq L$, we might

define $e(z)$ in some convenient way outside this range. If we were to do this, we might add to (2.6a) whatever is needed to make the equation valid for all z . Having done that, we could then take Fourier transforms and, hopefully, solve for the transforms, using such information about analyticity as we have available.

The first impulse, by analogy with Wiener-Hopf, is to define $e(z)$ as zero outside the range $0 \leq z \leq L$. This leads nowhere, for any function which is defined as zero for $z < 0$ and $z > L$ has a transform $E(k)$ which is analytic everywhere in the complex k plane. Then there are no analytic properties to exploit.

A second approach is motivated by the knowledge that the surface often produces boundary transients which die off with distance into the bulk; that is, that die off from the surface $z = 0$ towards larger z and that die off from the surface $z = L$ towards smaller z . Thus, the field $e(z)$ can be expressed in the slab as a sum of two parts, one of which decays to the left and the other of which decays to the right. We denote the leftwards-growing function as $f(z)$ and, purely for convenience in later manipulations, we denote the rightwards-growing function as $g(L - z)$. That is, in the region $0 \leq z \leq L$ where $e(z)$ is defined by (2.6), we write

$$e(z) = f(z) + g(L - z). \tag{3.1}$$

The choice of defining $e(z)$ outside the slab is now the choice of defining f and g outside the slab. We want f and g to extend to infinity so that their transforms will have singularities, and yet we want them to decay in the correct direction. This can be satisfied by demanding that

$$f(z) \sim e^{-\mu z}; \quad g(z) \sim e^{-\mu z}, \quad z \rightarrow +\infty, \tag{3.2a}$$

where $\mu > 0$. Then, we also want the singularities of $F(k)$ and $G(k)$, the Fourier transforms of f and g , to be confined to one half the complex k plane. This can be done by choosing

$$f(z) = 0, \quad g(z) = 0, \quad z < 0. \tag{3.2b}$$

At this point we can substitute (3.1) into the integral equation (2.6a), to get

$$\begin{aligned} & \left[\frac{d^2}{dz^2} + k_0^2 \right] [f(z) + g(L - z)] \\ & + i\omega\mu_0 \int_0^L \sigma(|z - z'|) [f(z') + g(L - z')] dz' = 0, \\ & \qquad \qquad \qquad 0 \leq z \leq L. \end{aligned} \tag{3.3}$$

[There will be no need to use $e(z)$ again; the program will be to determine f and g and, only at the end, to evaluate $e(z)$ from (3.1).]

It is unlikely that (3.3) will be valid outside the range $0 \leq z \leq L$, so now, again by analogy with Wiener-Hopf, we add to (3.3) just what is needed to make it valid everywhere. In this case we write the compensating function as $h(z) + i(L - z)$, which we add to the right side of (3.3) to obtain

$$\begin{aligned} & \left[\frac{d^2}{dz^2} + k_0^2 \right] [f(z) + g(L - z)] \\ & + i\omega\mu_0 \int_0^L \sigma(|z - z'|) [f(z') + g(L - z')] dz' \\ & \qquad \qquad \qquad = h(z) + i(L - z). \end{aligned} \tag{3.4}$$

The two functions must be chosen so that (3.4) is valid for all z , and also so that $h(z) + i(L - z)$ will vanish in the range $0 \leq z \leq L$. This last is necessary in order that (3.3) should be a consequence of (3.4). One way of guaranteeing this vanishing is to demand that

$$h(z) = 0, \quad i(z) = 0, \quad z \geq 0.$$

This still leaves considerable freedom in the choice of h and i separately for $z < 0$. Before making this choice, however, let us note that because f and g are defined to vanish at $z < 0$, the following two statements are true:

$$\begin{aligned} & \int_0^L \sigma(|z - z'|) f(z') dz' \\ & = \int_{-\infty}^{\infty} \sigma(|z - z'|) f(z') dz' - \int_L^{\infty} \sigma(|z - z'|) f(z') dz', \\ & \int_0^L \sigma(|z - z'|) g(L - z') dz' \\ & = \int_{-\infty}^{\infty} \sigma(|z - z'|) g(L - z') dz' \\ & \qquad \qquad \qquad - \int_{-\infty}^0 \sigma(|z - z'|) g(L - z') dz'. \end{aligned}$$

This we insert into Eq. (3.4). A slight rearrangement of the terms yields

$$\begin{aligned} & \left(\frac{d^2}{dz^2} + k_0^2 \right) f(z) + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(z - z') f(z') dz' \\ & - i\omega\mu_0 \int_{-\infty}^0 \sigma(z - z') g(L - z') dz' - h(z) \\ & = - \left[\left(\frac{d^2}{dz^2} + k_0^2 \right) g(L - z) \right. \\ & \quad + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(z - z') g(L - z') dz' \\ & \quad \left. - i\omega\mu_0 \int_L^{\infty} \sigma(z - z') f(z') dz' - i(L - z) \right]. \end{aligned} \tag{3.5}$$

We have at this stage far more unknowns than equations. In particular, $h(z)$ and $i(z)$ are not separately defined. Hence, we can choose to satisfy Eq.

(3.5) by demanding that each side separately vanishes. This will provide an extra equation for determining the two functions f and g . In the equations resulting from the vanishing of the right side, we replace z by $L - z$. The equation resulting from the vanishing of the left side is left as it stands. The two equations, vanishing of the left side and of the right side, are

$$\left(\frac{d^2}{dz^2} + k_0^2\right)f(z) + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(|z - z'|)f(z') dz' - i\omega\mu_0 \int_{-\infty}^0 \sigma(|z - z'|)g(L - z') dz' = h(z), \quad (3.6a)$$

$$\left(\frac{d^2}{dz^2} + k_0^2\right)g(z) + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(|z - z'|)g(z') dz' - i\omega\mu_0 \int_{-\infty}^0 \sigma(|z - z'|)f(L - z') dz' = i(z). \quad (3.6b)$$

It is clear from (3.6) that the demand that f and g vanish at negative z demands that $h(z)$ and $i(z)$, the compensating functions, should satisfy

$$h(z) = -i\omega\mu_0 \int_{-\infty}^0 \sigma(|z - z'|)g(L - z') dz' + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(|z - z'|)f(z') dz', \quad z < 0, \quad (3.7a)$$

$$i(z) = -i\omega\mu_0 \int_{-\infty}^0 \sigma(|z - z'|)f(L - z') dz' + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(|z - z'|)g(z') dz', \quad z < 0. \quad (3.7b)$$

Recall also that we have defined

$$h(z) = 0, \quad i(z) = 0, \quad z \geq 0. \quad (3.7c)$$

In addition to (3.6) we have to supply boundary conditions for f and g . Two boundary conditions have already been given—exponential decay at infinity. The other two may be obtained by substituting (3.1) into (2.6b, c)

$$\left[f(z) - \frac{1}{ik_0} \frac{df}{dz}\right]_{z=L} + \left[g(z) + \frac{1}{ik_0} \frac{dg}{dz}\right]_{z=0} = 0, \quad (3.8a)$$

$$\left[f(z) + \frac{1}{ik_0} \frac{df}{dz}\right]_{z=0} + \left[g(z) - \frac{1}{ik_0} \frac{dg}{dz}\right]_{z=L} = 2. \quad (3.8b)$$

In Eqs. (3.6) and (3.8), we have a pair of coupled integral equations and a pair of coupled boundary conditions which, for the two-sided slab, are the exact analogs of the single integral equation and boundary condition for the one-sided semi-infinite slab.

Equations (3.6), (3.7), and (3.8) possess considerable symmetry with respect to interchange of f and g , and of h and i . Only the boundary condition (3.8)

lacks this interchange symmetry. To take advantage of this, it is convenient to introduce symmetric and antisymmetric boundary conditions which give rise to functions $f^\pm(z)$, $g^\pm(z)$, etc. These replace $f(z)$, $g(z)$, etc., in (3.6) and (3.7), but instead of (3.8), they satisfy

$$\left[f^\pm(z) - \frac{1}{ik_0} \frac{d}{dz} f^\pm(z)\right]_{z=L} + \left[g^\pm(z) + \frac{1}{ik_0} \frac{d}{dz} g^\pm(z)\right]_{z=0} = \pm 1, \quad (3.9a)$$

$$\left[f^\pm(z) + \frac{1}{ik_0} \frac{d}{dz} f^\pm(z)\right]_{z=0} + \left[g^\pm(z) - \frac{1}{ik_0} \frac{d}{dz} g^\pm(z)\right]_{z=L} = 1. \quad (3.9b)$$

Then, if $f^\pm(z)$, etc., satisfy (3.6), (3.7), and (3.9), the functions

$$f(z) = f^+(z) + f^-(z), \quad (3.10a)$$

$$g(z) = g^+(z) + g^-(z) \quad (3.10b)$$

will satisfy (3.6), (3.7), and (3.8). Furthermore, on interchanging f^\pm and g^\pm in the equations, one finds that

$$f^\pm(z) = \pm g^\pm(z), \quad (3.11a)$$

$$h^\pm(z) = \pm i^\pm(z). \quad (3.11b)$$

Inserting (3.11) into (3.6), (3.9) gives

$$\left(\frac{d^2}{dz^2} + k_0^2\right)f^\pm(z) + i\omega\mu_0 \int_{-\infty}^{\infty} \sigma(|z - z'|)f^\pm(z') dz' \mp i\omega\mu_0 \int_{-\infty}^0 \sigma(|z - z'|)f^\pm(L - z') dz' = h^\pm(z), \quad (3.12a)$$

$$\left[f^\pm(z) + \frac{1}{ik_0} \frac{d}{dz} f^\pm(z)\right]_{z=0} \pm \left[f^\pm(z) - \frac{1}{ik_0} \frac{d}{dz} f^\pm(z)\right]_{z=L} = 1. \quad (3.12b)$$

This form is the most suitable starting point for the use of Fourier transforms we have been able to devise. The Fourier techniques to be used will make extensive use of the analytic properties of the transforms. For this reason, the next essential task is to study the analytic properties of the transforms of the quantities in (3.12a).

4. ANALYTIC PROPERTIES OF THE TRANSFORMS

Consider first the transform of $f(z)$. (We suppress superscripts \pm .)

$$F(k) \equiv \int_{-\infty}^{\infty} f(z)e^{-ikz} dz. \quad (4.1)$$

The requirements of Eq. (3.2) guarantee analyticity

of $F(k)$ for $\text{Im } k < \mu$, where μ is the positive rate of exponential decay. Secondly, like almost all other kernels arising in physical problems, the conductivity kernel has a range; call it l . That is, $|\sigma(z)| \sim e^{-|z|/l}$ as $z \rightarrow \pm\infty$. Hence, its transform

$$V(k) \equiv -i\omega\mu_0 \int \sigma(z)e^{-ikz} dz \quad (4.2)$$

will be analytic in the strip $-1/l < \text{Im } k < 1/l$. Thirdly, the definition provided by (3.7) guarantees that $H(k)$, the transform of $h(z)$, will be analytic for $-1/l < \text{Im } k < \infty$. Finally, the transform of the last term on the left of (3.12a) is analytic in the strip $-1/l < \text{Im } k < 1/l$. It turns out that $\mu \leq 1/l$, so that all transforms are analytic in the range $-1/l < \text{Im } k < \mu$. Therefore, we can take the transform of Eq. (3.12a) for k in the strip $-1/l < \text{Im } k < \mu$ in the usual way.

The transformed equation is

$$[-k^2 + k_0^2 - V(k)]F(k) \mp i\omega\mu_0 \int_{-\infty}^{\infty} dz e^{-ikz} \int_{-\infty}^0 \sigma(|z-z'|)f(L-z') dz' = H(k) + f'(0) + ikf(0). \quad (4.3)$$

In the integral here, σ and f may be expressed in terms of their transforms:

$$\pm \left(\frac{1}{2\pi}\right)^2 \int dz e^{-ikz} \int_{-\infty}^0 dz' \int_{-\infty}^{\infty} V(q_1)e^{iq_1(z-z')} dq_1 \times \int_{-\infty}^{\infty} F(q)e^{iq(L-z')} dq.$$

The z integration gives a delta function which removes the q_1 integration, leaving

$$\pm \frac{1}{2\pi} \int_{-\infty}^0 dz' V(k)e^{-ikz'} \int_{-\infty}^{\infty} dq F(q)e^{iq(L-z')}.$$

The z' integral will converge if $\text{Im}(k+q) > 0$. We choose q along the real axis, which means that we regard k as lying in the upper half plane. Performing the z' integration, we insert the result into (4.3) and find

$$[-k^2 + k_0^2 - V(k)]F(k) \mp V(k) \frac{1}{2\pi i} \int_{-\infty}^{\infty} \frac{F(q)e^{iqL}}{q+k} dq = H(k) + f'(0) + ikf(0). \quad (4.4)$$

This equation is valid over the entire k plane, not just the upper half k plane, if the integral is extended by analytic continuation.

It is convenient to define a function $\psi(k)$ and an integral operator $K(k)$ by

$$k^2 - k_0^2 + V(k) \equiv \psi(k), \quad (4.5a)$$

$$\frac{1}{2\pi i} \int_{-\infty}^{\infty} \frac{dq e^{iqL}}{q+k} F(q) \equiv K(k)F. \quad (4.5b)$$

The function $\psi(k)$ is equal to $k^2 - k_0^2\epsilon(k, \omega)$, the function whose zeros give the dispersion relation for the infinite-medium problem. This suggests the central role this function will play in the discussion to follow. The operator $K(k)$ is unique to the finite-slab problem: Note that it vanishes exponentially as the thickness L of the slab increases. The important property of K is that $K(k)F$ is analytic in the upper half plane for any function $F(k)$ which vanishes as k goes to infinity.

Using the notation (4.5), we rewrite (4.4) as

$$\psi(k)[F(k) \pm K(k)F] = -S(k), \quad (4.6a)$$

$$S(k) = H(k) + f'(0) + ikf(0) \pm (k_0^2 - k^2)K(k)F. \quad (4.6b)$$

The previous discussion of analyticity has indicated that every single term on the right of (4.6b) is analytic in the upper half plane. This means that $S(k)$ is analytic in the upper half plane. Yet, the function $\psi(k)$ appearing on the left of (4.6a) has singularities in the upper half plane. Clearly, this situation imposes some constraints on the contents of the square brackets in order that the singularities of $\psi(k)$ on the left do not appear in $S(k)$ on the right. These *restrictions*, which are completely unrelated to the analytic considerations involved in the standard Wiener-Hopf method, provide the equations we must ultimately solve. It is these restrictions which we must now consider.

Before doing so, let us note that only the unknown transform $F(k)$, and not the unknown $H(k)$, appears in these square brackets. The restrictions on the contents of the brackets will not involve $H(k)$ and, since these restrictions will determine $F(k)$, there will be no further need for considering the unknown $H(k)$. $H(k)$ can of course be calculated from (3.7) once $F(k)$ is known, but there is no reason to do so. This situation is analogous to that in the standard Wiener-Hopf, wherein the compensating function, which is necessary at the outset, becomes superfluous *before* the calculation is completed.¹⁶

5. RESTRICTIONS RESULTING FROM ANALYTICITY REQUIREMENTS

Consider the form of Eq. (4.6a). In those regions of the upper half k plane where $\psi(k)$ is analytic and nonzero, the upper half plane analyticity of K and of S demand that $F(k)$ be analytic also. However, at those points where $\psi(k)$ vanishes, $F(k)$ can have poles without destroying the analyticity of S . If there

¹⁶ See, for example, P. M. Morse and H. Feshbach, *Methods of Theoretical Physics* (McGraw-Hill Book Co., Inc., New York, 1953), p. 978.

are N such points, then $F(k)$ will have the form

$$F(k) = \sum_1^N \frac{\varphi_n}{k - k_n} + \Phi(k), \quad (5.1)$$

where

$$\psi(k_n) = 0, \quad \text{Im } k_n > 0, \quad n = 1, \dots, N, \quad (5.2)$$

and where $\Phi(k)$ is analytic in the upper half plane at all points that $\psi(k)$ is analytic. [The φ_n are numbers which are yet to be determined. However, we turn our attention first to $\Phi(k)$.] Since $F(k)$ is analytic below $\text{Im } k = \mu > 0$, it follows that the only singularities of $\Phi(k)$ are in the upper half plane and that these coincide with the singularities of ψ . To proceed further, one must describe the singularities of ψ . Although the discussion here will be couched in terms of the dispersion function for the electromagnetic problem, it will be obvious how an arbitrary kernel whose Fourier transform is given will have to be treated.

If the electron gas representing the metal is treated as a zero-temperature noninteracting system of free fermions and if the scattering time is treated as constant, then the conductivity tensor has a Fourier transform such that¹⁰

$$\psi(k) = k^2 - k_0^2 - \frac{3\mu_0 n q^2 \omega}{2 k p_F} \left[\frac{1}{2} \left(1 - \frac{\beta^2}{k^2} \right) \ln \frac{\beta - k}{\beta + k} - \frac{\beta}{k} \right], \quad (5.3)$$

where

$$\beta = (\pm 1 + \omega/\omega_c)/R + (i/l). \quad (5.4)$$

(The \pm here refers to the sense of circular polarization, as in 2.1 to 2.5.) Here, n is the number density of the electron gas, p_F its Fermi momentum, ω_c the cyclotron frequency, $R = V_F/\omega_c$ is the cyclotron radius and $l = V_F\tau$ is the mean free path. The function $\psi(k)$ has branch points at $k = \pm\beta$. If we trace through the details of the derivation of ψ from $\sigma(|z - z'|)$, it appears that a natural choice of branch cuts is along the two lines $k = \pm\beta u$, where $1 \leq u < \infty$. Thus, the only upper half plane singularity of ψ is a branch cut running along the line $k = \beta u$, and $\Phi(k)$ is analytic in the entire cut plane, that is, at all k excluding the points $k = \beta u$.

At infinity, $\Phi(k)$ must decrease at least as rapidly as $1/k$, for otherwise $F(k)$ would not have the $1/k$ behavior demanded of Fourier transforms in general and by Eq. (4.4) in particular. The fact that $F(k)$ is a Fourier transform imposes conditions on the behavior of $\Phi(k)$ near the branch cut. For instance, imagine the inversion of $F(k)$ as being carried out along a contour which has been swept upward from the real axis to encircle any poles of F and the cut of Φ . The

contribution to the contour from the neighborhood of the branch point will diverge if Φ grows faster than a pole at $k = \beta$, that is, near $k = \beta$, $\Phi(k)$ must satisfy $|\Phi(k)| \leq \text{const}/|k - \beta|^\alpha$, where $\alpha \leq 1$. Excluding the case of a pole at the branch point (which could be treated separately), we have $\alpha < 1$.

The main condition determining $\Phi(k)$ arises when we try to make $S(k)$ analytic on the cut. The functions in (4.6a) seem to take on different values as k approaches the cut from one side or the other, but we must ensure that $S(k)$ has the same value, no matter from which side k approaches the cut. If it does, then $S(k)$ will be analytic across the cut. Accordingly, let a plus or minus superscript denote the value of ψ , F , KF , and S as k approaches the point βu in the cut from one side or the other; that is, let $\psi^\pm(\beta u)$, $F^\pm(\beta u)$, $K^\pm(\beta u)F$, and $S^\pm(\beta u)$ denote the limiting values of $\psi(k)$, $F(k)$, etc., where $k = \beta(u \pm i\epsilon)$ and where ϵ (a real positive number) goes to zero in the limit. If $S(k)$ is to be analytic across the cut, it is necessary that $S^+(\beta u) = S^-(\beta u)$. Using (4.6a) to enforce this equality,

$$\begin{aligned} \psi^+(\beta u)[F^+(\beta u) \pm K^+(\beta u)F] \\ = \psi^-(\beta u)[F^-(\beta u) \pm K^-(\beta u)F]. \end{aligned} \quad (5.5)$$

However, K is analytic on the cut (as it is throughout the upper half plane) and so we can drop the \pm superscript and write $K^\pm(\beta u)F = K(\beta u)F$. We can evaluate KF by using (5.1) in (4.5b) as

$$K(\beta u)F = \sum \frac{\varphi_n e^{ik_n L}}{k_n + \beta u} + K(\beta u)\Phi \quad (5.6a)$$

with

$$K(k)\Phi \equiv \frac{1}{2\pi i} \int_{-\infty}^{\infty} \frac{dq e^{iqL}}{q + k} \Phi(q). \quad (5.6b)$$

Using (5.1) and (5.6) in (5.5) gives

$$\begin{aligned} \psi^+(\beta u)\Phi^+(\beta u) - \psi^-(\beta u)\Phi^-(\beta u) &= [\psi^-(\beta u) - \psi^+(\beta u)] \\ &\times \left[\sum \frac{\varphi_n}{\beta u - k_n} \pm \sum \frac{\varphi_n e^{ik_n L}}{\beta u + k_n} \pm K(\beta u)\Phi \right] \end{aligned} \quad (5.7)$$

as the condition the discontinuity in Φ must satisfy if $S(k)$ is to be analytic. This is the fundamental equation to be solved. The method of solving it draws heavily on the theory of singular integral equations as developed by Muskhelishvili and has points of similarity with the method of elementary solutions developed by Case¹⁷ and used by Zelazny *et al.*^{18,19} for neutron transport problems. The form of this equation seems to suggest that the unknown function $\Phi(k)$ can

¹⁷ K. M. Case, Ann. Phys. 9, 1 (1960).

¹⁸ R. Zelazny, A. Kuzell, and J. Mika, Ann. Phys. 16, 69 (1961).

¹⁹ R. Zelazny and A. Kuzell, Physica 27, 797 (1961).

be regarded as being proportional to the unknown coefficients φ_n so that, as an ansatz, one could write

$$\Phi(k) = \sum_n \varphi_n A_n(k).$$

A separate equation for each of the $A_n(k)$ functions would result. This would still leave the coefficients φ_n to be determined. The spatial boundary condition (3.12b) can be utilized to determine one of the coefficients, but it is not yet clear what information will fix the others. Moreover, in typical dispersion functions, even the number of roots available, and hence the number of coefficients to be determined, changes with changes in the physical parameters. Thus, the number of added conditions needed to fix the φ_n will somehow have to depend on the physical parameters, which at first sight seems unusual.

The resolution of this paradoxical situation appears in the next section, for there we shall see that a set of mathematical side conditions, which must be imposed to solve for the discontinuity in $\Phi(k)$ using Muskhelishvili methods, can be satisfied only if certain relations among the φ_n are true. The number of these relations, happily, is equal to the number of coefficients to be determined and thus, the same techniques which determine the $A_n(k)$ functions will also determine the φ_n coefficients.

6. SOLUTION OF THE DISCONTINUITY CONDITION

Were it not for the presence of $\pm K\Phi$ on the right of (5.7), we could rewrite (5.7) as

$$\Phi^+(\beta u) = G(\beta u)\Phi^-(\beta u) + g(\beta u), \quad (6.1a)$$

where

$$G(\beta u) \equiv \psi^-(\beta u)/\psi^+(\beta u) \quad (6.1b)$$

and

$$g(\beta u) \equiv \{[\psi^-(\beta u) - \psi^+(\beta u)]/\psi^+(\beta u)\}p(\beta u) \quad (6.1c)$$

are given functions along the cut. In our case, however,

$$p(\beta u) \equiv \sum_1^N \varphi_n \left[\frac{1}{\beta u - k_n} \pm \frac{e^{ik_n L}}{\beta u + k_n} \right] \pm K(\beta u)\Phi \quad (6.2)$$

contains the unknown function Φ and so is not *a priori* known. Fortunately, the troublesome term is of order $e^{-L/l}$. This is because $\Phi(k)$ is analytic for k below β ; that is, for $\text{Im } k < 1/l$, and hence the q integration in (5.6b) can be shifted upwards as far as $\text{Im } q = 1/l$. It is this fact which underlies the convergence of the iterative solution to be given.

Equation (6.1a) is the classical inhomogeneous Hilbert problem on an arc as discussed by Musk-

helishvili.²⁰ His solution, in outline form, is essentially the following: Suppose that there is available a function $X(k)$ with the following properties:

(a) $X(k)$ is analytic in the cut plane, the cut running along $k = \beta u$. (6.3a)

(b) $X(k)$ is nonvanishing in the cut plane. (6.3b)

(c) $X(k)$ does not vanish at the branch point $k = \beta$. (6.3c)

(d) $X(k)$ grows less rapidly than a pole at $k = \beta$. (6.3d)

(e) $X(k)$ satisfies the boundary condition along the cut that

$$X^+(\beta u)/X^-(\beta u) = G(\beta u). \quad (6.3e)$$

(f) At infinity, $X(k)$ behaves in one of the following three ways:

(1) $X(k) \rightarrow \text{constant}$. (6.3f1)

(2) $X(k) \rightarrow k^{-n}$ (n positive integer). (6.3f2)

(3) $X(k) \rightarrow k^m$ (m positive integer). (6.3f3)

Then using (6.3e), Eq. (6.1a) becomes

$$\frac{\Phi^+(\beta u)}{X^+(\beta u)} - \frac{\Phi^-(\beta u)}{X^-(\beta u)} = \frac{g(\beta u)}{X^+(\beta u)}, \quad (6.4)$$

which is a boundary discontinuity condition on the function $\Phi(k)/X(k)$. [In the same sense, Eq. (5.7) was a boundary condition on the function $\psi(k)\Phi(k)$, but Φ/X has certain properties which $\psi\Phi$ does not possess; namely, because of (6.3a) and (6.3b), Φ/X is analytic *everywhere except along the cut*.] As a result, one can use the Plemelj formulas

$$1/(X \pm i\epsilon) = PX^{-1} \pm i\pi\delta(X)$$

to verify that a solution to (6.4) is

$$\frac{\Phi(k)}{X(k)} = \frac{1}{2\pi i} \int_1^\infty \frac{g(\beta t)\beta dt}{(\beta t - k)X^+(\beta t)} + P_s(k) \quad (6.5)$$

where $P_s(k)$ is an arbitrary polynomial of degree s . (s is as yet unspecified.) The conditions (6.3c) and (6.3d) (as well as some on G and g we have not mentioned explicitly, but which are true here) lead to the conclusion, as Muskhelishvili shows, that (6.5) is the only solution to (6.4).

²⁰ Muskhelishvili is careful to stress that his treatment applies to arcs of finite length, whereas we are here working on an arc (the branch cut of ψ) which extends to infinity. If one wishes to retain the mathematical rigor of Muskhelishvili, the way to do it is to alter the definition of the spatial kernel $\sigma(|z - z'|)$ at $|z - z'| \rightarrow 0$ so as to cause the branch cut in $\psi(k)$ to terminate at another branch point at large k . The discussion given here will change only in slight detail. Then at the very end of the problem, the singularity in $\sigma(|z - z'|)$ at $z = z'$ can be restored, pushing the added branch point out to infinity.

The presence or absence of $P_s(k)$ and the value of s depend on the behavior of $X(k)$ at infinity. If Φ/X vanishes at infinity, then no polynomial $P_s(k)$ can be tolerated in (6.5). For instance, in case (6.3f1), where $X \rightarrow \text{constant}$, the $1/k$ behavior of Φ is enough to make Φ/X vanish at infinity. The polynomial $P_s(k)$ will be absent for similar reasons in the case (6.3f3) where $X(k)$ has algebraic growth at infinity. On the other hand, for the case (6.3f2), where $X \sim k^{-n}$, $\Phi/X \sim k^{n-1}$, a polynomial of degree $s = n - 1$ is allowed. This introduces n arbitrary constants into the solution for $\Phi(k)$.

If perchance $X(k)$ does go like k^m , then Φ/X goes as $1/k^{m+1}$ and Eq. (6.5), even with $P_s(k)$ deleted, does not in general provide a solution unless $g(\beta t)$ has just those properties which cause the integral also to have $1/k^{m+1}$ behavior as $k \rightarrow \infty$. The properties needed are obtained by expanding $(t - k)^{-1}$ in the integral as a power series in t/k and demanding that the coefficients of $1/k, 1/k^2, \dots, 1/k^m$ vanish. Thus, if $X(k)$ goes as k^m , Eq. (6.5) provides a solution only if m additional conditions are met:

$$\int_1^\infty \frac{t^l g(\beta t) dt}{X^+(\beta t)} = 0, \quad l = 0, 1, \dots, m - 1. \quad (6.6)$$

We see now that case (6.3f1) provides a unique mathematical solution, (6.3f2) provides a mathematical solution with n arbitrary constants, and case (6.3f3) provides a unique mathematical solution only under m additional constraints. It will turn out that the N constants φ_n which appear in $g(\beta u)$ provide exactly the amount of freedom needed to obtain a unique answer to the underlying physical problem in all three cases. However, to show this and to continue with the solution, it is necessary to construct $X(k)$.

Muskhelishvili also gives the prescription for constructing the function $X(k)$. Taking the logarithm of (6.3e) gives him

$$\log X^+(\beta u) - \log X^-(\beta u) = \log G(\beta u),$$

which can be regarded as a boundary condition on $\log X$. This leads him to consider the function

$$\Gamma(k) = \frac{1}{2\pi i} \int_1^\infty \frac{\log G(\beta u) \beta du}{\beta u - k} \quad (6.7a)$$

and

$$X_0(k) = \exp \Gamma(k). \quad (6.7b)$$

The integral here will converge only if $G(\beta u)$ goes to 1 at infinity (it does for any $\sigma(z)$ which has a transform), and even then only if we choose that branch of the logarithm for which $\log \psi^-(\beta u) = \log \psi^+(\beta u)$ as $u \rightarrow \infty$. This choice of branches will be made in all that follows.

Clearly $X_0(k)$ satisfies all the conditions of (6.3) except possibly (6.3c) and (6.3d), relating to growth or decay near the branch point. Muskhelishvili then, without destroying any of the other conditions, multiplies or divides $X_0(k)$ by integer powers of $(\beta - k)$ to produce a function

$$X(k) = (\beta - k)^\lambda X_0(k), \quad (6.8)$$

λ a positive or negative integer or zero, which stays finite at $k = \beta$ and which does not grow as fast as a pole there. Such a function satisfies all conditions (6.3). Its behavior at infinity, which follows from (6.7) and (6.8), is as k^λ , which puts it in one of the three cases (6.3), depending on the value of λ .

The choice of λ depends on behavior of $\Gamma(k)$ in the neighborhood of the branch point. One can write, again following Muskhelishvili,

$$\begin{aligned} \Gamma(k) &= \frac{1}{2\pi i} \int_1^\infty \left[\frac{\log G(\beta)}{\beta u - k} + \frac{\log G(\beta u) - \log G(\beta)}{\beta u - k} \right] \beta du \\ &= \frac{1}{2\pi i} [-\log(\beta - k) \log G(\beta) + \xi(k)], \end{aligned}$$

where ξ is bounded at $k = \beta$. We shall anticipate an important result which we will prove shortly and write $\log G(\beta) = 2\pi i(N - 1)$, where N is the number of upper half plane roots of $\psi(k)$, that is, the number of coefficients φ_n to be determined. Because of this we have

$$\Gamma(k) = -(N - 1) \log(\beta - k) + (2\pi i)^{-1} \xi(k)$$

which, using (6.7) and (6.8), gives

$$X(k) = (\beta - k)^{\lambda+1-N} \exp \xi(k) / 2\pi i.$$

The only choice of integer λ which causes $X(k)$ to be nonzero at $k = \beta$, but keeps its growth at β slower than a pole, is $\lambda = N - 1$. Putting this value in (6.8),

$$X(k) = (\beta - k)^{N-1} \exp \frac{1}{2\pi i} \int_1^\infty \frac{\log G(\beta u) \beta du}{\beta u - k}. \quad (6.9)$$

This choice of λ renders the physical problem determinate: Suppose first that $N = 0$. Then $X(k) \sim k^{-1}$, which means $\varphi(k)/X(k) \sim \text{constant}$ at infinity. Hence, a polynomial of degree zero (a constant) can appear in (6.5) giving

$$\Phi(k) = \frac{X(k)}{2\pi i} \int_1^\infty \frac{g(\beta t) \beta dt}{(\beta t - k) X^+(\beta t)} + P_0 X(k) \quad (6.10a)$$

and, with $N = 0$, Eq. (5.1) gives

$$F(k) = \Phi(k). \quad (6.10b)$$

The Fourier transform $F(k)$ thus contains one arbitrary constant, P_0 , whose value can be adjusted to

satisfy (3.12b), the spatial boundary condition on the field. On the other hand, if N is a positive integer, then $X(k) \sim k^{N-1}$, which means that

$$\Phi(k) = \frac{X(k)}{2\pi i} \int_1^\infty \frac{g(\beta t)\beta dt}{(\beta t - k)X^+(\beta t)} \quad (6.11)$$

is a solution, provided that $N - 1$ side conditions of the form (6.6) are satisfied. These $N - 1$ constraints imposed on the N constants φ_n , which appear in $g(\beta u)$, still leave enough freedom to satisfy the one spatial boundary condition on the field.

The idea that the side conditions (6.6) can determine the coefficients of the poles appears in the paper by Case.¹⁷ Implied in his treatment is the necessity of a condition relating the number of poles and the Muskhelishvili index κ ($= -\lambda$). An explicit demonstration that the condition is satisfied for the transport equation appears in the paper by Zelazny, Kuzell, and Mika.¹⁸ In our case, a similar demonstration can be made, also based simply on counting the zeros of $\psi(k)$. The relevant theorem is that if a function

$$\psi(k) = R(k)e^{i\theta(k)}, \quad R, \theta \text{ real}, \quad (6.12)$$

is analytic within and on the boundary of a region, then the number of zeros enclosed by the boundary is $1/2\pi$ times the change in θ around the boundary. Here, we take a contour composed of the circle at infinity deformed inwards to enclose the two branch cuts of ψ , along $k = \pm\beta u$, $1 \leq u < \infty$. Since $\psi(k)$ is even, half of its zeros are in the upper half plane and half the change in θ is contributed by that part of the contour lying in the upper half plane. Hence N , the number of upper half plane zeros of $\psi(k)$, is equal to $1/2\pi$ times the change in θ along the upper half plane contour. The contour is composed of the upper infinite semicircle and the loop

$$k = \beta(u - i\epsilon), \quad \infty \rightarrow u \rightarrow 1 \quad (\text{in below cut}),$$

$$k = \beta + i\beta\epsilon e^{i\varphi}, \quad \pi \rightarrow \varphi \rightarrow 0$$

(semicircle about branch),

$$k = \beta(u + i\epsilon), \quad 1 \rightarrow u \rightarrow \infty \quad (\text{out above cut}),$$

surrounding the upper branch cut.

At infinity, $\psi(k) = k^2 + V(k) \approx k^2$, which gives a phase change of 2π around the upper infinite semicircle. Along the cut, $\psi(k)$ takes on a perfectly definite value at $k = \beta$. Hence, there is no change in θ on the tiny semicircle around the branch point. Moreover, because $|\psi(k)|$ is the same on either side of the cut at $k = \beta$ and again at $k = \infty$, the change in θ can be expressed as the change in $i^{-1} \log \psi$. Thus

$$N = (2\pi)^{-1}(2\pi + i^{-1}\Delta \log \psi)$$

where $\Delta \log \psi$ is the change in $\log \psi$ coming in below the cut and going out above the cut:

$$\Delta \log \psi = [\psi^-(\beta) - \psi^-(\infty)] - [\psi^+(\infty) - \psi^+(\beta)].$$

Recalling the necessity of choosing that branch of $\log \psi^+$ for which $\log(\psi^-/\psi^+) \rightarrow 0$ at infinity, we are left with

$$N = 1 + (2\pi i)^{-1} \log \psi^-(\beta)/\psi^+(\beta)$$

or, from (6.1b)

$$\log G(\beta) = 2\pi i(N - 1). \quad (6.13)$$

This is the important result which we anticipated prior to Eq. (6.9) and which is analogous to the demonstration by Zelazny *et al.*

7. COMPLETING THE SOLUTION

The properties of $\Phi(k)$ discussed so far guarantee that it can be represented as a Cauchy integral:

$$\Phi(k) = \int_1^\infty \frac{\varphi(t)\beta dt}{\beta t - k}. \quad (7.1)$$

This representation is useful both in the evaluation of the Fourier inversion of Eq. (5.1),

$$i^{-1}f^\pm(z) = \sum_{n=1}^N \varphi_n e^{ik_n z} - \int_1^\infty \varphi(t) e^{i\beta t z} \beta dt \quad (7.2)$$

and in the evaluation of $K\Phi$, Eq. (5.6b),

$$K(\beta u)\Phi = - \int_1^\infty \frac{\varphi(t) e^{i\beta t L} dt}{u + t}. \quad (7.3)$$

In order to derive an equation governing $\varphi(t)$, we first apply the Plemelj formulas to Eq. (7.1) to obtain

$$2\pi i\varphi(v) = \Phi^+(\beta v) - \Phi^-(\beta v). \quad (7.4)$$

Then we insert (6.2) into (6.1c) and that into (6.10a) or (6.11) to obtain

$$\begin{aligned} \Phi(k) &= \frac{X(k)}{2\pi i} \int_1^\infty \frac{\beta dt}{(\beta t - k)X^+(\beta t)} \left[\frac{\psi^-(\beta t) - \psi^+(\beta t)}{\psi^+(\beta t)} \right] \\ &\times \left\{ \sum_{n=1}^N \varphi_n \left[\frac{1}{\beta t - k_n} \pm \frac{e^{ik_n L}}{\beta t - k_n} \right] \right. \\ &\left. \mp \int_1^\infty \frac{\varphi(u) e^{i\beta u L} du}{u + t} \right\} + \delta_{N0} P_0 X(k). \end{aligned} \quad (7.5)$$

Note, however, that because of (6.1b) and (6.3e)

$$\frac{1}{X^+(\beta t)} \left[\frac{\psi^-(\beta t) - \psi^+(\beta t)}{\psi^+(\beta t)} \right] = \frac{1}{X^-(\beta t)} - \frac{1}{X^+(\beta t)}.$$

We can therefore regard the t integration in (7.5) as being a line integral in the complex η plane, where the path of integration is inward along $\eta = \beta(t + i\epsilon)$

and outward along $\eta = \beta(t - i\epsilon)$. That is,

$$\Phi(k) = \frac{X(k)}{2\pi i} \int_C \frac{d\eta}{(\eta - k)X(\eta)} \left\{ \sum_{n=1}^N \varphi_n \left[\frac{1}{\eta - k_n} \pm \frac{e^{ik_n L}}{\eta + k_n} \right] \mp \int_1^\infty \frac{\varphi(u)e^{i\beta u L} \beta du}{\eta + \beta u} \right\} + \delta_{N0} P_0 X(k). \quad (7.6)$$

The contour C consisting of the lines above and below the cut can be augmented by adding to it the infinitesimal semicircle around the branch point at $\eta = \beta$ without changing the value of the integral, because $X(\eta)$ takes on a perfectly definite finite value at $\eta = \beta$. This converts the contour C to a continuous path coming in from infinity above the cut, circling the cut, and going off to infinity again below the cut.

Suppose now, for definiteness, that $N \neq 0$. (The electromagnetic problem of primary interest will in fact have $N = 1$ or $N = 2$. For a general problem in which $N = 0$, a trivial modification of the $N \neq 0$ procedure is needed.) Then the η integrand in (7.6) goes as $[X(\eta)\eta^2]^{-1} = \eta^{-(N+1)}$ and therefore the contour C can be closed by adding to it the circle at infinity which starts at the lower side of the cut, encircles the complex η plane in the negative sense, and ends at the upper side of the cut. Within this closed contour, the only singularities are at $\eta = k$, $\eta = \pm k_n$ and $\eta = -\beta u$ and we can evaluate (7.6) using residues as

$$\Phi(k) = -\sum_{n=1}^N \varphi_n \left[\frac{1}{k - k_n} \pm \frac{e^{ik_n L}}{k + k_n} \right] \pm \int_1^\infty \frac{\varphi(u)e^{i\beta u L}}{k + \beta u} + X(k) \left\{ \sum_{n=1}^N \varphi_n \left[\frac{1}{(k - k_n)X(k_n)} \pm \frac{e^{ik_n L}}{(k + k_n)X(-k_n)} \right] \mp \int_1^\infty \frac{\varphi(u)e^{i\beta u L} \beta du}{(k + \beta u)X(-\beta u)} \right\}. \quad (7.7)$$

Using this to evaluate the right side of (7.4) gives the equation governing $\varphi(v)$, namely,

$$\varphi(v) = \frac{X^+(\beta v) - X^-(\beta v)}{2\pi i} \left\{ \sum_{n=1}^N \varphi_n \left[\frac{1}{(\beta v - k_n)X(k_n)} \pm \frac{e^{ik_n L}}{(\beta v + k_n)X(-k_n)} \right] \mp \int_1^\infty \frac{\varphi(u)e^{i\beta u L} du}{(u + v)X(-\beta u)} \right\}. \quad (7.8)$$

The linearity of this equation allows one to write

$$\varphi(v) = \sum_{n=1}^N \varphi_n f_n(v) \quad (7.9a)$$

where $f_n(v)$ satisfies the integral equation

$$f_n(v) = \frac{X^+(\beta v) - X^-(\beta v)}{2\pi i} \left\{ \frac{1}{(\beta v - k_n)X(k_n)} \pm \frac{e^{ik_n L}}{(\beta v + k_n)X(-k_n)} \mp \int_1^\infty \frac{f_n(u)e^{i\beta u L} du}{(u + v)X(-\beta u)} \right\}. \quad (7.9b)$$

This equation can be solved iteratively. It is clear that successive terms converge as $e^{-l\beta L} = e^{-L/l}$ and, in fact, if $\beta L > 1$, an asymptotic evaluation of the integral shows that convergence is more like $e^{-L/l}/(\beta L)^2$, so that even the zeroth iterate is close to exact.

Having iterated (7.9b) to sufficient accuracy to determine $f_n(v)$, we can determine the coefficients φ_n . Inserting (6.2) into (6.1c) and that into the subsidiary condition (6.6) gives

$$\int_1^\infty \frac{t^l dt}{X^+(\beta t)} \left[\frac{\psi^-(\beta t) - \psi^+(\beta t)}{\psi^+(\beta t)} \right] \times \left\{ \sum \varphi_n \left[\frac{1}{\beta t - k_n} \pm \frac{e^{ik_n L}}{\beta t + k_n} \right] \mp \int_1^\infty \frac{\varphi_n(u)e^{i\beta u L} du}{u + t} \right\} = 0, \quad l = 0, 1, \dots, N - 2.$$

The same steps which led from (7.5) to (7.7) can be used here to evaluate the t integral, giving

$$\frac{-1}{\beta^{l+1}} \sum \varphi_n \left[\frac{k_n^l}{X(k_n)} \pm \frac{(-k_n)^l e^{ik_n L}}{X(-k_n)} \mp \int_1^\infty \frac{(-\beta u)^l \varphi_n(u)e^{i\beta u L} \beta du}{X(-\beta u)} \right] = 0.$$

Inserting (7.9a) here then gives

$$\sum_{n=1}^N C_{ln} \varphi_n = 0, \quad l = 0, 1, \dots, N - 2, \quad (7.10a)$$

$$C_{ln} = \frac{k_n^l}{X(k_n)} \pm \frac{(-k_n)^l e^{ik_n L}}{X(-k_n)} \mp \beta^{l+1} \int_1^\infty \frac{(-u)^l f_n(u)e^{i\beta u L} du}{X(-\beta u)}, \quad (7.10b)$$

while inserting (7.9a) into (7.2), and that in turn into the spatial boundary condition (3.12b) gives

$$i \sum \varphi_n \left\{ \left(1 + \frac{k_n}{k_0} \right) \pm \left(1 - \frac{k_n}{k_0} \right) e^{ik_n L} - \beta \int_1^\infty \left[\left(1 + \frac{\beta t}{k_0} \right) \pm \left(1 - \frac{\beta t}{k_0} \right) e^{i\beta t L} \right] f_n(t) dt \right\} = 1. \quad (7.10c)$$

This is a system of N equations for the N coefficients φ_n . Having determined these, one has determined the functions (restoring the even and odd superscript \pm) in (7.2):

$$f^\pm(z) = i \sum_{n=1}^N \varphi_n^\pm \left[e^{ik_n z} - \beta \int_1^\infty f_n^\pm(t)e^{i\beta t z} dt \right]. \quad (7.11)$$

From these, the field in the slab follows by combining (3.1) and (3.10a) and (3.11a):

$$e(z) = f^+(z) + f^-(z) + f^+(L - z) - f^-(L - z). \quad (7.12)$$

8. DISCUSSION

The solution of the problem of determining the field in the slab is now complete. It has been reduced to performing the following sequential calculations:

- (a) Determine the roots of the dispersion relation (5.2).
- (b) Evaluate the functions $X(k)$ defined in (6.9).
- (c) Iterate (7.9) to sufficient accuracy to determine f_n .
- (d) Solve (7.10) for φ_n .
- (e) Evaluate the expressions (7.11) and (7.12) for the field $e(z)$.

Although these steps in general call for numerical methods, there are still cases where considerable analytic progress can be made. For instance, suppose the parameters of the problem are such that there is only one root of the dispersion function. Then there is no need to solve for the coefficients φ_n ; the single coefficient φ_1 is fixed by boundary conditions. If, furthermore, the thickness L were such that $L/l \gg 1$, then the first iterate of (7.9) might be sufficiently accurate. This first iterate can be obtained analytically, and one would then have a situation in which the wave corresponding to the root of the dispersion relation suffers Fabry-Perot resonances which are modified by the continuum solutions embodied in the integral.

In order to establish contact with the conventional Wiener-Hopf technique, we consider another case in which our equations for $E(k)$ can be solved explicitly—namely, the semi-infinite slab. There will be no restrictions on the kernel or on the number of roots other than the existence of a finite range l for the kernel. For the semi-infinite slab, we set L equal to infinity. Then (7.9) and (7.10) become

$$2\pi i f_n(t) = \frac{X^-(\beta t) - X^+(\beta t)}{(k_n - \beta t)X(k_n)} \quad (8.1a)$$

and

$$\sum_{n=1}^N \frac{k_n^l \varphi_n}{X(k_n)} = 0, \quad l = 0, 1, \dots, N-2. \quad (8.1b)$$

To solve (8.1b), we make use of the identity

$$\sum_{n=1}^N k_n^l / \left[\prod_{m \neq n}^N (k_m - k_n) \right] = 0, \quad l = 0, 1, \dots, N-2. \quad (8.2)$$

This identity is easily established by evaluating

$$I \equiv \frac{1}{2\pi i} \oint k^l dk / \left[\prod_{m=1}^N (k_m - k) \right], \quad l \leq N-2,$$

where the contour encloses all the k_m . Evaluating by

contours, we get the left side of (8.2). Substituting $k = 1/z$ and evaluating by contours (this time the contour encloses $z = 0$ and no singularities), we get the right side of (8.2).

To use (8.2), we define b_n by

$$\varphi_n / X(k_n) = b_n / \prod_{m \neq n} (k_m - k_n). \quad (8.3)$$

Then (8.1b) becomes

$$\sum k_n^l b_n / \left[\prod_{m \neq n} (k_m - k_n) \right] = 0, \quad l = 0, 1, \dots, N-2,$$

which, by virtue of (8.2), implies that all b_n are equal, say to ib_0 . Thus from (8.3),

$$\varphi_n = [ib_0 X(k_n)] / \left[\prod_{m \neq n} (k_m - k_n) \right]. \quad (8.4)$$

Combining this with (8.1a) gives

$$2\pi i \sum \varphi_n f_n(t) = ib_0 [X^-(\beta t) - X^+(\beta t)] \times \sum_{n=1}^N [(k_n - \beta t)]^{-1} \left[\prod_{m \neq n} (k_m - k_n) \right]^{-1}.$$

However, using (8.2) again with $k_{N+1} \equiv \beta t$ and $l = 0$ allows us to evaluate the sum here so that

$$2\pi i \sum \varphi_n f_n(t) = ib_0 [X^-(\beta t) - X^+(\beta t)] / \left[\prod_{m=1}^N (k_m - \beta t) \right]. \quad (8.5)$$

Inserting (8.4) and (8.5) into (7.11), we have

$$f^\pm(z) = -ib_0 \left\{ \sum_{n=1}^N \frac{X(k_n) e^{ik_n z}}{\prod_{m \neq n} (k_m - k_n)} - \frac{\beta}{2\pi i} \int_1^\infty \frac{[X^-(\beta t) - X^+(\beta t)] e^{i\beta t z} dt}{\prod_{m=1}^N (k_m - \beta t)} \right\}. \quad (8.6)$$

(For $N = 0$, the solution is just the integral with the factors in the denominator deleted.)

Since there is no difference between $f^\pm(z)$ in the $L = \infty$ limit, we have $e(z) = 2f^\pm(z)$ and the solution (8.6) is complete except for b_0 , a normalization factor which is easily fixed by considerations based on the relation between $E(k \rightarrow \infty)$ and the power series expansion of $e(z)$ about the origin.

The solution (8.6) given here can also be derived by a slight modification of the standard Wiener-Hopf technique.²¹ This form turns out to be most useful

²¹ The modification of the standard Wiener-Hopf method needed is to choose the function to which the Wiener-Hopf factorization will be applied to be free of all zeros, not just those in the strip $-1/l < \text{Im } k < \mu$. When this is done, all contour integrals arising can be shifted into the upper half plane where the only singularities will be the branch cut of the dispersion function. The contour can surround this cut and, after some manipulation, will give rise to the function $X(k)$ defined in (6.9) here.

for evaluating the large z behavior in those cases (such as helicon propagation below the doppler-shifted cyclotron edge) where the integral provides the long-ranged term.

ACKNOWLEDGMENTS

I should like to express my sincere thanks to Dr. M. Lax and Dr. D. E. McCumber for their encouragement and helpful comments on the presentation.

State Labeling of the Irreducible Representations of SU_n

V. SYAMALA DEVI AND L. S. R. K. PRASAD
Andhra University, Waltair, India

(Received 25 February 1967)

State labeling of the irreducible representations of SU_3 is done by using Littlewood's rules for the analysis of products of representations of unitary groups. The method is generalized to any SU_n .

INTRODUCTION

Through a purely algebraic and infinitesimal approach, Baird and Biedenharn¹ have studied the problem of state labeling for the group SU_n . They explicitly carried out the analysis for SU_3 and have pointed out that though the generalization of the method can be done in principle to any SU_n , it is however laborious. Here we derive the same results for SU_3 in a much simpler manner by using Littlewood's rules for the analysis of product of representations of unitary groups. This method is generalized to any SU_n .

1. STATE LABELING OF THE IRREDUCIBLE REPRESENTATIONS OF SU_3

The generators of SU_3 , in Cartan's canonical form can be chosen as

$$\begin{aligned} H_1 &= C_1^1 - C_3^3, & H_2 &= C_2^2 - C_3^3, & E_\alpha &= C_1^3, \\ E_\beta &= C_1^2, & E_\gamma &= C_2^3, & E_{-\alpha} &= C_3^1, \\ E_{-\beta} &= C_2^1, & E_{-\gamma} &= C_3^2, \end{aligned}$$

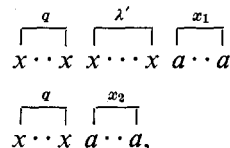
where² $C_\mu^{\mu'} = \sum_{s=1}^3 a_{\mu s} a_s^{\mu'}$, $\mu, \mu' = 1, 2, 3$ are generators of U_3 , and $a_{\mu s}$ and $a_s^{\mu'}$ are, respectively, the boson creation and annihilation operators. From the root diagram of SU_3 , the numerical coefficients that are to be taken for the generators can be found.³ The basis vectors for an IR (Irreducible Representation) of U_3 can be written as homogeneous polynomials in $a_{\mu s}$ operating on a certain vacuum ket $|0\rangle$. The inner product of the two states $P|0\rangle, P'|0\rangle$ is

$$(P, P') = \langle 0| P^+ P' |0\rangle,$$

where P^+ is obtained by replacing all $a_{\mu s}$ in P by $a_s^{\mu'}$.

The labeling problem is solved by the canonical factorization $SU_3)SU_2 \times U_1$. The generators of the SU_2 subgroup are $H_1, E_\alpha, E_{-\alpha}$ and $C_1^1 + C_3^3 - 2C_2^2 = H$ is the generator of U_1 commuting with the above SU_2 subgroup. The IR's of SU_2 and U_1 can be characterized by nonnegative integers λ' and m' , respectively. The IR's of $SU_2 \times U_1$ contained in the IR (h_1, h_2) of SU_3 are determined by Littlewood's rules⁴ for the multiplication of representations of unitary groups.

Let the symbols "x" denote the boxes in the Young diagram corresponding to (λ') of SU_2 [which is equivalent to $(\lambda' + q, q)$] and the symbols "a" denote the boxes of the diagram (m') of U_1 . Then the diagram corresponding to the IR (λ', m') of $SU_2 \times U_1$ contained in (h_1, h_2) of SU_3 is of the form



where

$$q + \lambda' + x_1 = h_1, \quad q + x_2 = h_2, \quad x_1 + x_2 = m'. \tag{1}$$

Littlewood's rules lead to the inequality

$$\lambda' \geq x_2. \tag{2}$$

Equations (1) determine x_1 and x_2 uniquely in terms of h_1, h_2, λ', m' , and hence the IR (λ', m') of $SU_2 \times U_1$ occurs, if at all, only once in (h_1, h_2) of SU_3 . The highest-weight state in (λ', m_s) is

$$P = (1)^{\lambda'+q-h_2} (13)^q (12)^{h_2-q} (2)^{h_1-\lambda'-q} \dots, \tag{3}$$

¹ G. E. Baird and L. C. Biedenharn, J. Math. Phys. 4, 1449 (1963).

² M. Moshinsky, J. Math. Phys. 4, 1128 (1963).

³ R. E. Behrends, J. Dreitlein, C. Fronsdal, and B. W. Lee, Rev. Mod. Phys. 34, 1 (1962).

⁴ D. E. Littlewood, Theory of Group Characters (Clarendon Press, Oxford, 1950), 2nd ed., p. 94.

for evaluating the large z behavior in those cases (such as helicon propagation below the doppler-shifted cyclotron edge) where the integral provides the long-ranged term.

ACKNOWLEDGMENTS

I should like to express my sincere thanks to Dr. M. Lax and Dr. D. E. McCumber for their encouragement and helpful comments on the presentation.

State Labeling of the Irreducible Representations of SU_n

V. SYAMALA DEVI AND L. S. R. K. PRASAD
Andhra University, Waltair, India

(Received 25 February 1967)

State labeling of the irreducible representations of SU_3 is done by using Littlewood's rules for the analysis of products of representations of unitary groups. The method is generalized to any SU_n .

INTRODUCTION

Through a purely algebraic and infinitesimal approach, Baird and Biedenharn¹ have studied the problem of state labeling for the group SU_n . They explicitly carried out the analysis for SU_3 and have pointed out that though the generalization of the method can be done in principle to any SU_n , it is however laborious. Here we derive the same results for SU_3 in a much simpler manner by using Littlewood's rules for the analysis of product of representations of unitary groups. This method is generalized to any SU_n .

1. STATE LABELING OF THE IRREDUCIBLE REPRESENTATIONS OF SU_3

The generators of SU_3 , in Cartan's canonical form can be chosen as

$$\begin{aligned} H_1 &= C_1^1 - C_3^3, & H_2 &= C_2^2 - C_3^3, & E_\alpha &= C_1^3, \\ E_\beta &= C_1^2, & E_\gamma &= C_2^3, & E_{-\alpha} &= C_3^1, \\ E_{-\beta} &= C_2^1, & E_{-\gamma} &= C_3^2, \end{aligned}$$

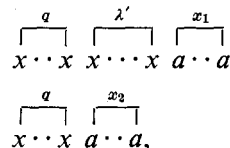
where² $C_\mu^{\mu'} = \sum_{s=1}^3 a_{\mu s} a_s^{\mu'}$, $\mu, \mu' = 1, 2, 3$ are generators of U_3 , and $a_{\mu s}$ and $a_s^{\mu'}$ are, respectively, the boson creation and annihilation operators. From the root diagram of SU_3 , the numerical coefficients that are to be taken for the generators can be found.³ The basis vectors for an IR (Irreducible Representation) of U_3 can be written as homogeneous polynomials in $a_{\mu s}$ operating on a certain vacuum ket $|0\rangle$. The inner product of the two states $P|0\rangle, P'|0\rangle$ is

$$(P, P') = \langle 0| P^+ P' |0\rangle,$$

where P^+ is obtained by replacing all $a_{\mu s}$ in P by $a_s^{\mu'}$.

The labeling problem is solved by the canonical factorization $SU_3)SU_2 \times U_1$. The generators of the SU_2 subgroup are $H_1, E_\alpha, E_{-\alpha}$ and $C_1^1 + C_3^3 - 2C_2^2 = H$ is the generator of U_1 commuting with the above SU_2 subgroup. The IR's of SU_2 and U_1 can be characterized by nonnegative integers λ' and m' , respectively. The IR's of $SU_2 \times U_1$ contained in the IR (h_1, h_2) of SU_3 are determined by Littlewood's rules⁴ for the multiplication of representations of unitary groups.

Let the symbols "x" denote the boxes in the Young diagram corresponding to (λ') of SU_2 [which is equivalent to $(\lambda' + q, q)$] and the symbols "a" denote the boxes of the diagram (m') of U_1 . Then the diagram corresponding to the IR (λ', m') of $SU_2 \times U_1$ contained in (h_1, h_2) of SU_3 is of the form



where

$$q + \lambda' + x_1 = h_1, \quad q + x_2 = h_2, \quad x_1 + x_2 = m'. \tag{1}$$

Littlewood's rules lead to the inequality

$$\lambda' \geq x_2. \tag{2}$$

Equations (1) determine x_1 and x_2 uniquely in terms of h_1, h_2, λ', m' , and hence the IR (λ', m') of $SU_2 \times U_1$ occurs, if at all, only once in (h_1, h_2) of SU_3 . The highest-weight state in (λ', m_s) is

$$P = (1)^{\lambda'+q-h_2} (13)^q (12)^{h_2-q} (2)^{h_1-\lambda'-q} \dots, \tag{3}$$

¹ G. E. Baird and L. C. Biedenharn, J. Math. Phys. 4, 1449 (1963).

² M. Moshinsky, J. Math. Phys. 4, 1128 (1963).

³ R. E. Behrends, J. Dreitlein, C. Fronsdal, and B. W. Lee, Rev. Mod. Phys. 34, 1 (1962).

⁴ D. E. Littlewood, Theory of Group Characters (Clarendon Press, Oxford, 1950), 2nd ed., p. 94.

where

$$(s_1 \cdots s_r) = \Delta_{s_1 s_2 \cdots s_r}^{12 \cdots r} = \sum_{(s_1 \cdots s_r)} \epsilon (s_1 \cdots s_r) a_{s_1 1} \cdots a_{s_r r}.$$

All the other states in the IR $(\lambda' m')$ are obtained by operating the lowering operator in SU_2 , $(E_{-\alpha})^k$ on P . All the states in the IR (h_1, h_2) may be labeled by λ', m' , and v' where v' is the eigenvalue of H_1 in SU_2 . The range of variation of λ' and m' is fixed by equations (1) and inequality (2). However, m' is not the eigenvalue of H which is the generator of U_1 .

To compare these results with those of Baird and Biedenharn, we observe that $\lambda' = 2\lambda$, where $\lambda(\lambda + 1)$ is the eigenvalue of the SU_2 Casimir invariant Λ^2 , with λ integral or half-integral and $m' = \frac{1}{3}(h_1 + h_2 - 6m)$, where $6m$ is the eigenvalue of H . This latter relation is obtained by operating H on the highest-weight state P . The nonnegative integers h_1, h_2 are the same as p, q in Ref. 1. Equations (1) and inequality (2) lead us to the following range for λ' and m' :

$$\begin{aligned} h_1 - h_2 &\leq \lambda' + m' \leq h_1 + h_2, \\ h_2 - h_1 &\leq \lambda' - m' \leq h_1 - h_2. \end{aligned} \tag{4}$$

These in turn give the range for λ and m , which is a parallelogram with vertices

$$\begin{aligned} [h_1/2, (h_1 - 2h_2)/6], & \quad [h_2/2, (h_2 - 2h_1)/6], \\ [(h_1 - h_2)/2, (h_1 + h_2)/6], & \quad [0, (2h_2 - h_1)/3]. \end{aligned}$$

These results are identical with those obtained in Ref. 1.

2. GENERALIZATION TO ANY SU_n

The generators of SU_n in Cartan's canonical form may be chosen as

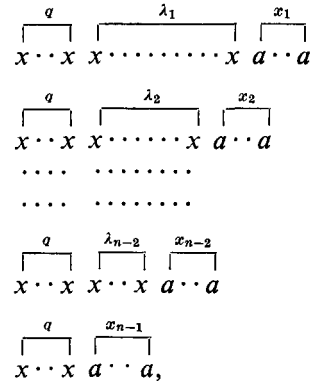
$$\begin{aligned} H_i &= C_i^i - C_{n-1}^n & i &= 1 \cdots n - 1, \\ E_\alpha &= C_i^j, \quad E_{-\alpha} = C_j^i & i < j &= 1 \cdots n. \end{aligned}$$

A canonical subgroup $SU_{n-1} \times U_1$ contained in the above SU_n has the following generators:

$$\begin{aligned} H_i &= C_i^i - C_{n-1}^{n-1}, & i &= 1 \cdots n - 2, \\ E_\beta &= C_i^j, \quad E_{-\beta} = C_j^i, & i < j &= 1 \cdots n - 1 \end{aligned}$$

which are the generators for SU_{n-1} and $H = \sum_{i=1}^{n-1} C_i^i - (n-1)C_n^n$ is the generator of U_1 commuting with the

above SU_{n-1} subgroup. The IR's (λ, m) where $(\lambda) = (\lambda_1, \lambda_2, \dots, \lambda_{n-2})$, of $SU_{n-1} \times U_1$ contained in the IR $(h) = (h_1 \cdots h_{n-1})$ of SU_n are determined according to Littlewood's rules by the following diagrams:



where

$$\begin{aligned} q + \lambda_r + x_r &= h_r, & r &= 1 \cdots n - 2, \\ q + x_{n-1} &= h_{n-1}, & x_1 + x_2 + \cdots + x_{n-1} &= m, \\ \lambda_r &\geq \lambda_{r+1} + x_{r+1}, & r &= 1 \cdots n - 3, & \lambda_{n-2} &\geq x_{n-1}. \end{aligned} \tag{5}$$

Equations (5) determine uniquely all x_r in terms of (h) , (λ) , and m . Hence, the IR (λ, m) of $SU_{n-1} \times U_1$ contained in (h) of SU_n is contained only once. The highest-weight state of the IR (λ, m) of $SU_{n-1} \times U_1$ is

$$\begin{aligned} P &= \prod_{r=1}^{n-3} (12 \cdots r)^{\lambda_r - \lambda_{r+1} - x_{r+1}} (12 \cdots n - 2)^{\lambda_{n-2} - x_{n-1}} \\ &\quad \times (12 \cdots n - 1)^q \prod_{r=1}^{n-2} (12 \cdots rn)^{x_{r+1}} (n)^{x_1} \end{aligned}$$

The highest-weight states P are labeled by the parameters (λ) and m which vary, subject to the equations and inequalities (5). The eigenvalue M of H is given by operating H on P .

After normalising the highest-weight polynomials according to the definition of the inner product mentioned, the reduced matrix elements of the generators of SU_n can be directly obtained by operating the generators on the highest-weight states P .

ACKNOWLEDGMENT

The authors are deeply indebted to Professor T. Venkatarayudu for his guidance.

Unitarity of Dynamical Propagators of Perturbed Klein-Gordon Equations

JOHN M. CHADAM

Department of Mathematics, Indiana University, Bloomington, Indiana

(Received 6 February 1967)

After discussing the basic notions of quantizations as representations of the Weyl relations, a criterion is established for a symplectic transformation on a classical linear system to be unitarily implementable in the free (zero-interaction) representation. The result is applied to the temporal propagators of $\square\varphi = m^2\varphi + K\varphi$ to obtain a condition which is sufficient to ensure that they are unitarily implementable in the free representation of the quantized Klein-Gordon field of mass m . Necessary conditions are also obtained when K commutes with $m^2I - \Delta$. Several examples are discussed, the most interesting of which is that of a mass jump (i.e., $K = m'^2I$), where the results given are fairly complete.

I. INTRODUCTION

A general, mathematically rigorous prescription for the quantization of linear systems has been available for quite some time.^{1,2} The procedure has been applied to the relativistic equations of modern physics describing the motions of free (fundamental) particles to give a precise meaning to the notion of their associated quantum fields. Using the above prescription, it is possible to quantize the perturbed Klein-Gordon equation $\square\varphi = m^2\varphi + K\varphi$ (K linear, bounded, perhaps time-dependent). However, it is not true, in general, that the transition from one time to another defined by this equation can be effected by a unitary transformation in the conventional (free, zero-interaction) representation of the quantized Klein-Gordon field associated with mass m . This would mean that not all of the dynamical matrix elements between free-particle normalizable states would be finite. From another point of view this would also imply that the total energy would not be a self-adjoint operator (observable) in this representation. In this paper, conditions on K are determined which ensure the unitary implementability of these dynamical propagators in the conventional representation of the Klein-Gordon field.

In Sec. II, the basic notions of quantizations and symplectic transformations are introduced. Conditions for symplectic transformations on general classical linear systems to be unitarily implementable in the conventional representation are discussed. A criterion equivalent to those given by Shale³ and Segal,⁴ but more functional in the present situation, is deduced. This condition is used, in Sec. III, to establish a sufficient condition for the unitary implementability of the above propagators in the conventional repre-

sentation of the Klein-Gordon field. The analysis is carried out for the case when K is a positive, bounded, linear, time-independent operator on the space of square-integrable functions defined on the spatial region over which the perturbed equation is defined.

The criterion is applied, in Sec. IV, to several more particular examples. In the commutative case (K commutes with $m^2I - \Delta$) a necessary condition is derived. It is also shown that the propagators associated with a mass jump (i.e., $K = m'^2I$) are never unitarily implementable in the free-field representation of mass m regardless of the dimension of space. However, in the cutoff theory, (i.e., for the spatial region being $[0, 2\pi]^n$) the result holds true if the number of dimensions, n , is ≤ 3 .

II. QUANTIZATIONS AND SYMPLECTIC TRANSFORMATIONS

Before discussing the problems indicated in the Introduction, a short exposition of the definitions and general theory of the transformations of field observables and states is given (following Refs. 1, 2, and 3). Given a real linear space L endowed with a non-degenerate, skew-symmetric, bilinear form $B(\cdot, \cdot)$, a Weyl system over (L, B) is defined as a mapping $x \rightarrow W(x)$ from L into unitary operators on a complex Hilbert space \mathcal{H}_W satisfying the Weyl relations

$$W(x)W(x') = \exp[\frac{1}{2}iB(x, x')]W(x + x')$$

along with the condition that $W(\cdot)$ is weakly continuous when restricted to finite-dimensional subspaces of L . Associated with any Weyl system $W(\cdot)$, the concrete Weyl algebra of bounded field observables is defined to be the uniform closure of $\cup_F \mathcal{O}_F$, where F ranges over the finite-dimensional subspaces of L , and \mathcal{O}_F is the weakly closed ring of operators generated by $\{W(x); x \in F\}$.

In the particular case when L is a real Hilbert space and $B(\cdot, \cdot)$ is $\text{Im}(\cdot, \cdot)_c$ where $(\cdot, \cdot)_c$ is the (complex) inner product of the complexification of L , such Weyl systems

¹ I. E. Segal, *Illinois J. Math.* **6**, 500 (1962).

² I. E. Segal, *Lectures in Applied Mathematics* (American Mathematical Society, Providence, R. I., 1963), Vol. 2.

³ D. Shale, *Trans. Am. Math. Soc.* **103**, 149 (1962).

⁴ I. E. Segal, *Trans. Am. Math. Soc.* **88**, 12 (1958).

are known to exist. In fact, when L is an infinite-dimensional Hilbert space⁵ there are continuously many unitarily inequivalent Weyl systems as opposed to the finite-dimensional situation where the Weyl (exponentiated) form of the Schrödinger representation is essentially the only Weyl system. Thus the concrete Weyl algebras are, in general, not unitarily equivalent. However, it is known⁶ that any two Weyl algebras arising from Weyl systems W and W' , are algebraically isomorphic by means of a unique isomorphism which takes $W(x)$ into $W'(x)$ for all x in L . Hence it is possible to define a unique abstract C^* -Weyl algebra \mathcal{O} with distinguished elements $W(x)$ satisfying the Weyl relations, as the equivalence class of all such concrete Weyl algebras. Then, in the usual manner, there is associated with each state E of \mathcal{O} a unique (to within unitary equivalence) representation of \mathcal{O} as operators on a complex Hilbert space. However, for an arbitrary state, the images of the distinguished elements $W(x)$ of \mathcal{O} under the associated representation will not in general satisfy the continuity condition mentioned previously. For this reason one usually restricts his attention to the physically more relevant representations associated with regular states; namely, those states E such that if $A, B \in \mathcal{O}$ then $E[A^*W(x)B]$ is continuous, as a function of x , on each finite-dimensional subspace of L .⁷ Representations of the Weyl relations relative to these states will be called quantizations. Our concern, in fact, will be only with pure regular states. Recalling that any pure state determines an irreducible representation, it is clear that pure regular states on \mathcal{O} give rise to irreducible quantizations.

Let $Sp(L)$ denote the multiplicative group of real bounded linear transformation on L which preserve $B(\cdot, \cdot)$. Clearly these are precisely the type of bounded linear transformations required to ensure that $x \rightarrow W(Tx)$ is a quantization if $W(\cdot)$ is. For $T \in Sp(L)$, it has been shown⁸ that $\theta(T)$ defined by

$$[\theta(T)W](x) = W(Tx), \quad x \in L$$

can be extended to a unique $*$ -automorphism of \mathcal{O} [likewise denoted by $\theta(T)$]. Furthermore any $*$ -automorphism θ of \mathcal{O} determines an automorphism θ^* of the states of \mathcal{O} through the contragredient

action

$$(\theta^*E)(A) = E(\theta^{-1}A), \quad A \in \mathcal{O}.$$

Both $\theta(\cdot)$ and $\theta^*(\cdot)$ are multiplicative. For $T \in Sp(L)$, $\theta^*(T)$ leaves the set of pure and/or regular states invariant.

The above forms the basis of the general theory of quantizations and transformations of field observables and states. Additional results will be given in the subsequent sections as they are needed. The main problem is now stated in this general context.

Suppose E is a pure regular state and $W(\cdot)$ on \mathcal{K}_W is a representative of the equivalence class of unitarily equivalent quantizations determined by E . Then $W(Tx)$, $x \in L$ is a well-defined operator on \mathcal{K}_W and for $T \in Sp(L)$ the mapping $x \rightarrow W(Tx)$ is a quantization. The problem then is to determine conditions on $T \in Sp(L)$ which will ensure that the two quantizations are unitarily equivalent; i.e., when can the induced action of T on \mathcal{K}_W be implemented by a unitary transformation Y_T satisfying $Y_T W(x) = W(Tx) Y_T$ for all $x \in L$. In the future when this situation occurs T will be said to be unitarily implementable.

The first result in this direction in the particular case of the conventional representation^{9,10} was given by Segal.⁴ Using the theory of integration over Hilbert space, it was shown that, if T is a closed, densely defined, linear operator on L with dense range, then T is unitarily implementable in the conventional representation if and only if $S = (T^*T)^{\frac{1}{2}}$ is nonsingular and has the form $I + A$ where A is Hilbert-Schmidt. Denote by (A) the set of operators which satisfy the above conditions.

Theorem 2.1. $T \in (A)$ if and only if T is bounded, linear, with bounded inverse, and $S = I + A$ where A is Hilbert-Schmidt.

Proof. Suppose $T \in (A)$. Since T is closed, densely defined, and linear it has a unique polar decomposition $T = US$ where S is self-adjoint with domain $D(S) = D(T)$ and U is a partial isometry with initial domain $U^*U = \overline{R(T^*)} = \overline{R(S)}$ ¹¹ and final domain $UU^* = \overline{R(T)} = L$ since $R(T)$ is dense. But S is nonsingular;

⁵ This will be the case of interest in the subsequent sections. In the remainder of this section the discussion will be restricted to this situation.

⁶ The proof is given in Ref. 6 in more generality.

⁷ A Weyl algebra \mathcal{O} over (L, B) , where L is a Hilbert space with B as above, not only has an abundance of states, it also has regular states.

⁸ I. E. Segal, Kgl. Danske Videnskab. Selskab, Mat.-Fys. Medd. 31, No. 12 (1959).

⁹ The conventional (zero-interaction, free) representation is characterized (see Ref. 1) as being the unique one which preserves positivity of energy. For a suitable choice of L and $B(\cdot, \cdot)$, it is precisely the Fock-Cook representation found in the usual treatment of quantum field theory. Cook (Ref. 10) was the first to give it a mathematically rigorous foundation. Henceforth the state to which this representation corresponds will be denoted by E_0 and the representation by W_0 .

¹⁰ J. M. Cook. Trans. Am. Math. Soc. 74, 222 (1953).

¹¹ $\overline{R(T^*)}$ is the closure of the range of T^* . Here we are confusing projections with their ranges. For a proof of the polar decomposition and the definition of "nonsingular," see M. A. Naimark, *Normed Rings* (P. Noordhoff Ltd., Groningen, The Netherlands, 1959), pp. 284 and 285.

hence $R(S)$ is dense in L and consequently $U^*U = L$. Thus $T = US$ where U is unitary.^{12,13} Now $S = I + A$, where A is Hilbert-Schmidt. Thus S is bounded on $D(S)$ and is closed since S is self-adjoint. Therefore $D(S) = \overline{D(S)}$ and hence $D(T) = D(S) = \overline{D(S)} = \overline{D(T)} = L$ since T is densely defined. The preceding, along with $\|Tx\|^2 = \|Sx\|^2 \leq \|S\|^2 \|x\|^2$ for all $x \in L$, and the fact that $S = I + A$ is bounded implies that T is everywhere defined on L and bounded. Since $R(S)$ is dense in L , S is one-one. For if $Sx = 0$, then $(Sx, y) = 0$ and $(x, Sy) = 0$ for all $y \in L$. But $R(S)$ is dense and so $x = 0$. Thus S possesses an inverse, S^{-1} , with $D(S^{-1}) = R(S)$ and $R(S^{-1}) = D(S) = L$. However $S = I + A$ where A is Hilbert-Schmidt and hence completely continuous and $R(I + B)$ is closed for any completely continuous operator B .¹⁴ Thus $R(S) = \overline{R(S)} = L$ since $R(S)$ is dense in L . Thus $D(S^{-1}) = R(S) = L$ and so S^{-1} is an everywhere defined, bounded, linear operator by the Hellinger-Toeplitz theorem. Consequently $T^{-1} = S^{-1}U^*$ exists as an everywhere defined bounded linear operator.

Conversely, suppose T is linear, bounded with bounded inverse on L and $S = (T^*T)^{\frac{1}{2}} = I + A$ where A is Hilbert-Schmidt. Since T is bounded, linear and everywhere defined, T is closed. Also $R(T) = D(T^{-1}) = L$. Thus T is closed, densely defined, linear with dense range and $S = I + A$ where A is Hilbert-Schmidt. Since T is bounded $S = (T^*T)^{\frac{1}{2}}$ is a (bounded) self-adjoint, and hence closed, operator with $D(S) = D(T) = L$. It only remains to show that $R(S)$ is dense in L . Since $\overline{R(S)} = \overline{R(T^*)}$ it suffices to show that $R(T^*)$ is dense. This is clearly the case because if $(y, T^*x) = 0$ for all $x \in D(T^*) = L$, then $(Ty, x) = 0$ for all $x \in L$. Thus $Ty = 0$ and hence $y = T^{-1}Ty = T^{-1}0 = 0$ which concludes the proof.

Thus it has been established that $T \in Sp(L)$ is unitarily implementable in the conventional representation if and only if it satisfies the new condition of Theorem 2.1. Indeed, precisely the same condition has been established by Shale³ employing a lengthier but more sophisticated and far-reaching analysis. However, it is desirable to modify the criterion slightly so that it is more easily applicable to the situations arising in the sections to follow. There, L is represented as the direct sum of two Hilbert spaces and hence the operators on L will be 2×2 matrices with operator valued entries. It would facilitate

matters if the difficulty of taking the square root of such transformations could be eliminated, for which purpose we give the following:

Lemma 2.2. If T is a bounded linear operator on L , then $(T^*T)^{\frac{1}{2}} - I$ is Hilbert-Schmidt if and only if $T^*T - I$ is Hilbert-Schmidt.

Proof. $0 \leq (T^*T)^{\frac{1}{2}} \leq \|T\| I < \infty$ since T is bounded. Therefore $I \leq (T^*T)^{\frac{1}{2}} + I \leq (\|T\| + 1)I < \infty$. Thus $[(T^*T)^{\frac{1}{2}} + I]^{-1}$ exists as a bounded operator on L . Now by the operational calculus $T^*T - I = [(T^*T)^{\frac{1}{2}} - I][(T^*T)^{\frac{1}{2}} + I]$. The conclusion now follows immediately from the facts that $(T^*T)^{\frac{1}{2}} + I$ is a bounded operator with a bounded inverse and that AB and BA are Hilbert-Schmidt if B is bounded and A is Hilbert-Schmidt.

Thus Theorem 2.1, Lemma 2.2, and the intervening discussion give the following:

Corollary 2.3. $T \in Sp(L)$ is unitarily implementable in the conventional representation if and only if T is a bounded linear operator with bounded inverse satisfying the property that $T^*T - I$ is Hilbert-Schmidt.

The above is the condition which is employed in the subsequent sections. This section is concluded with another easy result which will be an aid in the later discussions.

Proposition 2.4. The set of operators which satisfy the condition in Corollary 2.3 is a multiplicative subgroup of $Sp(L)$ which is closed under the taking of adjoints.

Proof. If T_1 and T_2 satisfy the hypothesis then T_1T_2 is bounded with bounded inverse $T_2^{-1}T_1^{-1}$. Also

$$\begin{aligned} (T_1T_2)^*(T_1T_2) &= T_2^*T_1^*T_1T_2 = T_2^*(I + X)T_2 \\ &= T_2^*T_2 + T_2^*XT_2 = I + Y + T_2^*XT_2 \end{aligned}$$

where $T_1^*T_1 = I + X$ and $T_2^*T_2 = I + Y$, X, Y Hilbert-Schmidt and hence $Y + T_2^*XT_2$ is Hilbert-Schmidt. If T satisfies the hypothesis then T^{-1} is bounded with bounded inverse and

$$\begin{aligned} I &= (TT^{-1})^*(TT^{-1}) = (T^{-1})^*T^*TT^{-1} \\ &= (T^{-1})^*(I + X)T^{-1} = (T^{-1})^*T^{-1} + (T^{-1})^*XT^{-1}. \end{aligned}$$

Thus $(T^{-1})^*T^{-1} = I - (T^{-1})^*XT^{-1}$ where $X = T^*T - I$ is Hilbert-Schmidt. Clearly the identity satisfies the hypothesis trivially and the associative law follows the associativity of multiplication of operators. If T satisfies the hypothesis then by Theorem 2.1,

¹² This is essentially the condition given by Seidman.¹³
¹³ T. I. Seidman, Commun. Pure Appl. Math. 17, 493 (1964).
¹⁴ C. f. J. Dieudonne, Foundations of Modern Analysis (Academic Press Inc., 1960), pp. 315 and 316.

$T = U(T^*T)^{\frac{1}{2}}$ where U is unitary. Thus,

$$(T^*)^*T^* = TT^* = UT^*TU^* = U(I + X)U^* = UU^* + UXU^* = I + UXU^*$$

and T^* is bounded with bounded inverse $(T^{-1})^*$.

III. HOMOGENEOUS TIME-INDEPENDENT PERTURBATIONS OF THE KLEIN-GORDON EQUATION

The dynamics of a free scalar classical meson field is described by the Klein-Gordon equation,

$$\square \varphi = \left(\Delta - \frac{\partial^2}{\partial t^2} \right) \varphi = m^2 \varphi, \quad m > 0, \quad (1)$$

where $\varphi(x, t)$ is a scalar function of time and space (the number of whose dimensions is not important in this section). Once one has chosen a classical solution space and imposed a skew form on it, the various associated quantum fields can be obtained by employing the prescription outlined in the previous section. Perhaps the most direct method for obtaining a suitable classical state space is by studying Eq. (1) in its abstract vector-valued form. In vector form, Eq. (1) can be written as

$$\frac{d}{dt} \begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix} = \begin{pmatrix} 0 & I \\ -B^2 & 0 \end{pmatrix} \begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix}, \quad (2)$$

where $B^2 = m^2I - \Delta$. In the following we only demand that B^2 is self-adjoint in the space of square-integrable functions over the region in space in which Eq. (1) is being studied. For example this situation occurs when Eq. (1) is studied on E^n with zero boundary conditions at infinity. In this case B^2 will be the self-adjoint extension of $m^2I - \Delta$ with domain \mathcal{S} (the Schwartz space of rapidly decreasing functions). Another typical example which also will arise later is the case of studying (1) on the interval $[0, 2\pi]$ with periodic boundary conditions. Here B^2 is the self-adjoint operator $m^2I - \Delta$ with domain $\{f \in L^2[0, 2\pi]; f'' \text{ exists, a.e., } f'' \in L^2 \text{ and } f(0) = f(2\pi), f'(0) = f'(2\pi)\}$.

Since the discussion in this section is independent of spatial dimensions and regions, only the common features of the above particular examples are used. For instance if $m \geq 0$, B^2 is positive and hence the same is true for B and $B^{\frac{1}{2}}$ on their respective domains of definition. If $m > 0$, $(B^2)^j = B^{2j}$ has bounded inverses for $j \geq 0$. In fact, for the first example, $D(B^j)$ is the space of functions $f \in L^2(E^n)$ such that $(m^2 + k^2)^{j/2}$ times the Fourier transform of f is square-integrable where k^2 is the square of the distance to the origin in the dual of E^n . Clearly $D(B^j)$ is dense in $L^2(E^n)$ for all j . For the second example it only suffices to observe that the complete orthonormal

set (CONS)

$$\left\{ \frac{e^{ikx}}{(2\pi)^{\frac{1}{2}}} \right\}_{k=-\infty}^{\infty}$$

is in $D(B^j)$ for all j and hence $D(B^j)$ is dense in $L^2[0, 2\pi]$ for all j .

We now take the classical solution space to the real Hilbert space H_B obtained by completing $D(B^{\frac{1}{2}}) \oplus D(B^{-\frac{1}{2}})$ with respect to the inner product

$$(\varphi, \psi)_{H_B} = \left[\begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix}, \begin{pmatrix} \psi_1 \\ \psi_2 \end{pmatrix} \right]_{H_B} = (B^{\frac{1}{2}}\varphi_1, B^{\frac{1}{2}}\psi_1)_{L^2} + (B^{-\frac{1}{2}}\varphi_2, B^{-\frac{1}{2}}\psi_2)_{L^2}.^{15-17}$$

Formally, the solution of (2) is given by the propagator $U_0(t, s)$ defined by

$$\begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix}(t) = U_0(t-s) \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} = \begin{pmatrix} \cos [(t-s)B], B^{-1} \sin [(t-s)B] \\ -B \sin [(t-s)B], \cos [(t-s)B] \end{pmatrix} \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix}, \quad (3)$$

where

$$\begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \in H_B$$

is the initial data given at time s . By performing the relevant calculations on a dense subspace of H_B and extending by continuity it is clear that $t \rightarrow U_0(t)$ is a one-parameter group of orthogonal transformations on H_B with infinitesimal skew-adjoint generator

$$\begin{pmatrix} 0 & I \\ -B^2 & 0 \end{pmatrix}.$$

Thus, Eq. (3) defines rigorously the unique (generalized)¹⁸ solution of Eq. (2).

For the purposes of quantization, in addition to the classical state space H_B , we must also have a skew-symmetric form on H_B . With this end in mind we (following Ref. 19) examine the operator

$$J = \begin{pmatrix} 0 & B^{-1} \\ -B & 0 \end{pmatrix}$$

on H_B .

¹⁵ This space was chosen because in the particular case of three-space it corresponds most closely to physics in that it is the unique Lorentz-invariant solution space for the Klein-Gordon equation as can be seen by the results of Strauss.¹⁷

¹⁶ $\|\cdot\|$ and (\cdot, \cdot) henceforth denotes the norm and inner product in L^2 . The summands in H_B are written as $D[B^{\frac{1}{2}}]$ and $D[B^{-\frac{1}{2}}]$ where the inner product is that which each inherits as a subspace of H_B . Since $D[B^{\frac{1}{2}}]$ and $D[B^{-\frac{1}{2}}]$ are the closures of $R(B^{\frac{1}{2}}) = D(B^{-\frac{1}{2}}) = L^2$ and $R(B^{-\frac{1}{2}}) = D(B^{\frac{1}{2}})$, respectively, in the L^2 inner product, it is clear that the closure is necessary in the second summand only.

¹⁷ W. A. Strauss, Trans. Am. Math. Soc. **108**, 12 (1963).
¹⁸ By "generalized" we mean a strongly continuous solution of the operational integral equation associated with Eq. (2). This, of course, may not necessarily be a "strict" solution of Eq. (2).

¹⁹ Roe Goodman, Ph.D. thesis, Department of Mathematics, Massachusetts Institute of Technology, Cambridge, Mass. (1963).

Lemma 3.1. J , as a transformation on H_B is an isometry which satisfies $J^* = -J$.

Proof. Suppose

$$\begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix}$$

belongs to $D(B) \oplus D(B^{-\frac{1}{2}})$ which is a dense subspace of H_B . Then

$$\left\| \begin{pmatrix} 0 & B^{-1} \\ -B & 0 \end{pmatrix} \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \right\|_{H_B}^2 = \|B^{\frac{1}{2}}(B^{-1}\varphi_2)\|^2 + \|B^{-\frac{1}{2}}(B\varphi_1)\|^2.$$

Since

$$\varphi_2 \in D(B^{-\frac{1}{2}}) = L^2, \quad B^{-1}\varphi_2 \in D(B) \subset D(B^{\frac{1}{2}}),$$

and

$$B^{\frac{1}{2}}B^{-1}\varphi_2 = B^{-\frac{1}{2}}\varphi_2.$$

Similarly, because $\varphi_1 \in D(B) \subset D(B^{\frac{1}{2}})$, $B\varphi_1 \in L^2 = D(B^{-\frac{1}{2}})$, and $B^{-\frac{1}{2}}B\varphi_1 = B^{\frac{1}{2}}\varphi_1$. Thus

$$\left\| J \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \right\|_{H_B}^2 = \|B^{-\frac{1}{2}}\varphi_2\|^2 + \|B^{\frac{1}{2}}\varphi_1\|^2 = \left\| \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \right\|_{H_B}^2.$$

Because J preserves the norm on a dense subspace, J is an isometry on H_B .

In order to compute the adjoint of J , we first notice that since J is bounded it possesses a unique adjoint given by

$$\begin{pmatrix} 0 & (-B)^+ \\ (B^{-1})^{++} & 0 \end{pmatrix}$$

where $(-B)^+$ and $(B^{-1})^{++}$ are the adjoints of the bounded operators $-B: D[B^{\frac{1}{2}}] \rightarrow D[B^{-\frac{1}{2}}]$ and $B^{-1}: D[B^{-\frac{1}{2}}] \rightarrow D[B^{\frac{1}{2}}]$, respectively, where the meaning is clear from the definition of H_B . Consider the second transformation. If $\varphi \in D(B^{-\frac{1}{2}})$ and $\psi \in D(B)$ we have $(B^{\frac{1}{2}}B^{-1}\varphi, B^{\frac{1}{2}}\psi) = (B^{-1}\varphi, B\psi) = (B^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}B\psi)$. Thus $(B^{-1})^{++} \supset B$. For the first transformation if $\varphi \in D(B)$ and $\psi \in D(B^{-\frac{1}{2}})$ $(B^{-\frac{1}{2}}B\varphi, B^{-\frac{1}{2}}\psi) = (B\varphi, B^{-1}\psi)$. But for $\psi \in D(B^{-\frac{1}{2}}) = L^2$, $B^{-1}\psi \in D(B) \subset D(B^{\frac{1}{2}})$ and hence $(B^{-\frac{1}{2}}B\varphi, B^{-\frac{1}{2}}\psi) = (B^{\frac{1}{2}}\varphi, B^{\frac{1}{2}}B^{-1}\psi)$. Thus $(-B)^+ \supset -B^{-1}$ and $J^* \supset -J$. Now $-J$ is bounded on H_B thus giving the equality $J^* = -J$.

It is clear in a similar fashion that $J^2 = -I$ and $J^{-1} = -J$ on H_B . Thus to obtain the skew-symmetric form $\Lambda_B(\cdot, \cdot)$ on H_B , the real inner product is first complexified using J and then Λ_B is taken to be the imaginary part of this complex inner product; i.e., $\Lambda_B(\cdot, \cdot) = -(J\cdot, \cdot)_{H_B} \cdot (H_B, \Lambda_B)$ as defined is the unique, relativistically invariant, classical, dynamical system associated with Eq. (1) when it is being considered in all of three space. For this reason we take it to be the classical dynamical system in the general case. It will play a distinguished role in all of the discussions to follow in that we shall attempt to define the propagators of the perturbed equation as

operators on H_B and then examine their unitary implementability in the (equivalence class of) representations obtained by quantizing (H_B, Λ_B) relative to the Fock state E_0 .

Consider now the perturbed Klein-Gordon equation

$$\square\varphi = m^2\varphi + K\varphi, \tag{4}$$

where K is a positive, bounded, linear, time-independent operator on L^2 . In vector notation, Eq. (4) can be written as

$$\frac{d}{dt} \begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix} = \begin{pmatrix} 0 & I \\ -(B^2 + K) & 0 \end{pmatrix} \begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix}, \tag{5}$$

where $B'^2 = B^2 + K$ is a positive self-adjoint operator on $D(B^2)$ because K is positive and bounded. Thus B' is well-defined and formally the solution of Eq. (5) is given by means of the propagator $U(t)$ defined by

$$\begin{aligned} \begin{pmatrix} \varphi \\ \dot{\varphi} \end{pmatrix}(t) &= U(t-s) \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \\ &= \begin{pmatrix} \cos [(t-s)B'], & B'^{-1} \sin [(t-s)B'] \\ -B' \sin [(t-s)B'], & \cos [(t-s)B'] \end{pmatrix} \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix}, \end{aligned} \tag{6}$$

where

$$\begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix} \in H_B$$

is the initial data at time s . The boundedness and linearity of K ensures²⁰ that Eq. (6) gives the unique global (generalized) solution of Eq. (5) in H_B .

We now consider the main problem: that of determining whether the transformations $U_0(t)$ and $U(t)$ are unitarily implementable in the conventional representation. As expected, $U_0(t)$ is clearly unitarily implementable because $U_0(t)$ is orthogonal and hence is bounded with bounded inverse $U_0^*(t)$ and thus $U_0^*(t)U_0(t) - I = 0$. Hence the criterion established in Corollary 2.3 is trivially satisfied. It is straightforward to check that $U_0(t)$ is symplectic on (H_B, Λ_B) .

In attempting to apply the criterion to the transformation $U(t)$, one is led very quickly to computational problems arising from the unboundedness of Δ and the noncommutativity of Δ and K . For this reason it seems more advantageous to abandon the direct approach in favor of treating the problem in a series of steps. The main result which we prove is the following:

Theorem 3.2. The dynamical propagators $U(t)$ of the perturbed Klein-Gordon equation $\square\varphi = m^2\varphi + K\varphi$, where K is a positive, linear, bounded, time-independent operator on L^2 , are unitarily implementable in the conventional representation of (H_B, Λ_B)

²⁰ I. E. Segal, Ann. Math. 78, 339 (1963).

if $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I$ is Hilbert-Schmidt as an operator from L^2 into L^2 .

First consider the transformation

$$A = \begin{pmatrix} B^{-\frac{1}{2}} & 0 \\ 0 & B^{\frac{1}{2}} \end{pmatrix} \begin{pmatrix} B'^{\frac{1}{2}} & 0 \\ 0 & B'^{-\frac{1}{2}} \end{pmatrix}$$

on H_B .

Lemma 3.3. A is a bounded linear transformation on H_B with a bounded inverse.

Proof. First we investigate some general properties of B^2 and B'^2 and their roots. In general if P and Q are (unbounded) self-adjoint operators such that $0 \leq P \leq Q$, then $0 \leq (P)^{\frac{1}{2}} \leq (Q)^{\frac{1}{2}}$ whether they commute or not. In our case $0 < B^2 \leq B'^2 \leq B^2 + \|K\| I$ on the common domain $D(B^2)$ of these operator where $\|\cdot\|$ denotes the operator norm of a transformation on L^2 . Thus $0 < B \leq B' \leq (B^2 + \|K\| I)^{\frac{1}{2}}$ on the common domain of these operators. In addition $0 < B^2 + \|K\| I \leq B^2 + 2(\|K\|)^{\frac{1}{2}}B + \|K\| I$ on $D(B^2)$. Thus $0 < B \leq B' \leq (B^2 + \|K\| I)^{\frac{1}{2}} \leq B + (\|K\|)^{\frac{1}{2}}I$ on the common domains of these operators. That their domains are the same follows directly from a result of Heinz²¹; namely, if R, S are positive operators with $D(S) \subset D(R)$ and $\|Rx\| \leq \|Sx\|$ for all $x \in D(S)$, then $\|R^\theta x\| \leq \|S^\theta x\|$ for all $x \in D(S^\theta)$, where $0 \leq \theta \leq 1$. Thus the last inequality holds on $D(B) = D(B')$. Similarly we can obtain that $0 \leq B^{\frac{1}{2}} \leq B'^{\frac{1}{2}} \leq B^{\frac{1}{2}} + \|K\|^{\frac{1}{2}} I$ on $D(B^{\frac{1}{2}}) = D(B'^{\frac{1}{2}})$.

Now A is bounded on H_B if and only if $B^{-\frac{1}{2}}B'^{\frac{1}{2}}: D[B^{\frac{1}{2}}] \rightarrow D[B'^{\frac{1}{2}}]$ and $B^{\frac{1}{2}}B'^{-\frac{1}{2}}: D[B^{-\frac{1}{2}}] \rightarrow D[B'^{-\frac{1}{2}}]$ are both bounded. Now for $\varphi \in D(B)$ a dense subset of $D(B^{\frac{1}{2}})$ we have $\|B^{\frac{1}{2}}(B^{-\frac{1}{2}}B'^{\frac{1}{2}}\varphi)\|^2 = \|B'^{\frac{1}{2}}\varphi\|^2 = (B'\varphi, \varphi) \leq ([B + (\|K\|)^{\frac{1}{2}}I]\varphi, \varphi) \leq C_K(B\varphi, \varphi) = C_K \|B^{\frac{1}{2}}\varphi\|^2$. The existence of such a finite C_K arises from the fact that $B \geq mI > 0$ and $(\|K\|)^{\frac{1}{2}}$ is finite. Thus on the dense domain, $D(B)$, the operator norm of $B^{-\frac{1}{2}}B'^{\frac{1}{2}}: D[B^{\frac{1}{2}}] \rightarrow D[B'^{\frac{1}{2}}]$ is bounded by $(C_K)^{\frac{1}{2}} < \infty$ and hence the operator is bounded. Likewise for $\varphi \in D(B^{-\frac{1}{2}})$, $\|B^{-\frac{1}{2}}B^{\frac{1}{2}}B'^{-\frac{1}{2}}\varphi\|^2 = \|B'^{-\frac{1}{2}}\varphi\|^2$ since $R(B'^{-\frac{1}{2}}) = D(B'^{-\frac{1}{2}}) = D(B^{\frac{1}{2}})$ and $B^{-\frac{1}{2}}B^{\frac{1}{2}}$ is the identity on $D(B^{\frac{1}{2}})$. Now formally $B' \geq B$ implies that $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} \geq I$ and hence $(B^{-\frac{1}{2}}B'B^{-\frac{1}{2}})^{-1} \leq I$. Thus $B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}} = (B^{-\frac{1}{2}}B'B^{-\frac{1}{2}})^{-1} \leq I$ and hence $B'^{-1} = B^{-\frac{1}{2}}(B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}})B^{-\frac{1}{2}} \leq B^{-\frac{1}{2}}IB^{-\frac{1}{2}} = B^{-1}$. By rigorizing these computations, we obtain

$$\begin{aligned} & \|B^{-\frac{1}{2}}(B^{\frac{1}{2}}B'^{-\frac{1}{2}}\varphi)\|^2 \\ &= (B'^{-1}\varphi, \varphi) \leq (B^{-1}\varphi, \varphi) = \|B^{-\frac{1}{2}}\varphi\|^2. \end{aligned}$$

Thus A is bounded on H_B and similarly it can be

shown that

$$A^{-1} = \begin{pmatrix} B'^{-\frac{1}{2}}B^{\frac{1}{2}} & 0 \\ 0 & B^{\frac{1}{2}}B'^{-\frac{1}{2}} \end{pmatrix}$$

is a bounded operator on H_B .

Lemma 3.4. A is a symplectic transformation on (H_B, Λ_B) .

Proof. The bounded transformation A is symplectic if and only if it preserves the skew form Λ_B ; i.e.,

$$\begin{aligned} \Lambda_B(Au, Av) &= -(JAU, Av)_{H_B} \\ &= -(Ju, v)_{H_B} = \Lambda_B(u, v) \end{aligned}$$

for all $u, v \in H_B$. This is equivalent to the condition that (the adjoint of A in H_B) $A^* = JA^{-1}J^{-1}$. Now since A is bounded on H_B , A^* exists as a bounded operator on H_B and

$$A^* = \begin{pmatrix} (B^{-\frac{1}{2}}B'^{\frac{1}{2}})^+ & 0 \\ 0 & (B^{\frac{1}{2}}B'^{-\frac{1}{2}})^{++} \end{pmatrix},$$

where $(B^{-\frac{1}{2}}B'^{\frac{1}{2}})^+$ and $(B^{\frac{1}{2}}B'^{-\frac{1}{2}})^{++}$ are the unique adjoints of the bounded operators $B^{-\frac{1}{2}}B'^{\frac{1}{2}}: D[B^{\frac{1}{2}}] \rightarrow D[B'^{\frac{1}{2}}]$ and $B^{\frac{1}{2}}B'^{-\frac{1}{2}}: D[B^{-\frac{1}{2}}] \rightarrow D[B'^{-\frac{1}{2}}]$, respectively. For $\varphi \in D(B^{\frac{1}{2}})$ and $\psi \in D(B'^{\frac{1}{2}}B^{\frac{1}{2}}) = D(B)$,

$$\begin{aligned} (B^{\frac{1}{2}}B^{-\frac{1}{2}}B'^{\frac{1}{2}}\varphi, B^{\frac{1}{2}}\psi) &= (B'^{\frac{1}{2}}\varphi, B^{\frac{1}{2}}\psi) = (\varphi, B'^{\frac{1}{2}}B^{\frac{1}{2}}\psi) \\ &= (\varphi, BB^{-1}B'^{\frac{1}{2}}B^{\frac{1}{2}}\psi) = (B^{\frac{1}{2}}\varphi, B^{\frac{1}{2}}B^{-1}B'^{\frac{1}{2}}B^{\frac{1}{2}}\psi). \end{aligned}$$

Thus $(B^{-\frac{1}{2}}B'^{\frac{1}{2}})^+ \supset B^{-1}B'^{\frac{1}{2}}B^{\frac{1}{2}}$. Similarly for $\varphi, \psi \in L^2$,

$$\begin{aligned} (B^{-\frac{1}{2}}B^{\frac{1}{2}}B'^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}\psi) &= (B'^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}\psi) \\ &= (\varphi, B'^{-\frac{1}{2}}B^{-\frac{1}{2}}\psi) = (\varphi, B^{-1}BB'^{-\frac{1}{2}}B^{-\frac{1}{2}}\psi) \\ &= (B^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}BB'^{-\frac{1}{2}}B^{-\frac{1}{2}}\psi). \end{aligned}$$

Thus $(B^{\frac{1}{2}}B'^{-\frac{1}{2}})^{++} \supset BB'^{-\frac{1}{2}}B^{-\frac{1}{2}}$ and hence $A^* \supset JA^{-1}J^{-1}$. But both are bounded and hence equality holds.

The methods of Lemma 3.4 can be used to show that the bounded invertible transformation $U(t)$ is likewise symplectic on (H_B, Λ_B) . Since the symplectic transformations form a group all compositions of A and $U(t)$ and their inverses which are considered below will be bounded, invertible symplectic transformations.

Before returning to the mainstream of the proof of Theorem 3.2 we introduce some notions which greatly facilitate the analytic nature of our next discussion. Following Shale³ we say that two pure regular states, E and F , are relatively normalizable, $E \sim F$, if $F(X) = E(Y^*XY)$ for all X and some fixed $Y \in \mathcal{O}$.^{22,23} It is an immediate consequence of the

²¹ Cf. L. Nirenberg, *Functional Analysis* (New York University Lecture Notes 1961), p. 123.

²² In the case of pure regular states this definition is equivalent to the more general definition for arbitrary states given by Segal²³ via a theorem of Kadison. For the details see Ref. 3, p. 163.

²³ I. E. Segal, *Can. J. Math.* **13**, 1 (1961).

definition that “ \sim ” is an equivalence relation and that for two pure regular states, $E, F, E \sim F$, with $T \in Sp(L)$ then $\theta^*(T)E \sim \theta^*(T)F$. The result of Shale, which provides the connection between these concepts and our main problem, is the following: $T \in Sp(L)$ is unitarily implementable in the representations associated with a pure regular state E if and only if $E \sim \theta^*(T)E$.²⁴ Now, using the above ideas and the previously mentioned fact that θ^* is multiplicative, we shall determine preliminary conditions which will ensure that $U(t)$ is unitarily implementable in the representations associated with the pure regular Fock–Cook state E_0 .

By the last result, since $U(t) \in Sp(H_B, \Lambda_B)$, it is unitarily implementable in the Fock representation if and only if $E_0 \sim \theta^*[U(t)]E_0$. Using the transitivity of the equivalence relation “ \sim ”, it is clear that for $E_0 \sim \theta^*[U(t)]E_0$ it is sufficient (but not necessary) that all of the following are true: (a) $E_0 \sim \theta^*(A)E_0$, (b) $\theta^*(A)E_0 \sim \theta^*[U(t)A]E_0$ and (c) $\theta^*[U(t)A]E_0 \sim \theta^*[U(t)]E_0$. Since $A \in Sp[H_B, \Lambda_B]$ condition (a) is equivalent to A being unitarily implementable in the Fock representation. By Corollary 2.3., since A is bounded with a bounded inverse on H_B , this is equivalent to $A^*A - I$ being Hilbert–Schmidt as an operator on H_B . Since E_0 is a pure regular state and $A, U(t) \in Sp(H_B, \Lambda_B)$ [and hence $AU \in Sp(H_B, \Lambda_B)$], $\theta^*(A)E_0$ and $\theta^*(U(t)A)E_0$ are pure regular states.²⁵ Thus, by the result quoted in the last paragraph, since $A^{-1} \in Sp(H_B, \Lambda_B)$ also, condition (b) is equivalent to $\theta^*(A^{-1})[\theta^*(A)E_0] \sim \theta^*(A^{-1})\{\theta^*[U(t)A]E_0\}$. But, as previously mentioned, $\theta^*(\cdot)$ is multiplicative. Thus condition (b) is equivalent to $E_0 \sim \theta^*[A^{-1}U(t)A]E_0$

which is equivalent to $[A^{-1}U(t)A]^*[A^{-1}U(t)A] - I$ being Hilbert–Schmidt on H_B . Finally, since

$$\theta^*[U(t)A]E_0 \text{ and } \theta^*[U(t)]E_0$$

are pure regular states and $U(t)^{-1} = U(-t)$ is symplectic, condition (c) is equivalent to

$$\theta^*[U(t)^{-1}]\{\theta^*[U(t)A]E_0\} \sim \theta^*[U(t)^{-1}]\{\theta^*[U(t)]E_0\}$$

or equivalently, $\theta^*(A)E_0 \sim E_0$, since $\theta^*(\cdot)$ is multiplicative. Using the symmetry of the equivalence relation “ \sim ”, we find that condition (c) is precisely the condition that $E_0 \sim \theta^*(A)E_0$ or that $A^*A - I$ is Hilbert–Schmidt on H_B .

Summarizing the above discussion, it is clear that $U(t)$ is unitarily implementable in the Fock representation if

- (i) $A^*A - I \in \Sigma(H_B)$ ²⁶
- (ii) $(A^{-1}U(t)A)^*(A^{-1}U(t)A) - I \in \Sigma(H_B)$.

These provisional conditions are now examined to obtain equivalent conditions on the perturbation K . In this respect the second condition is vacuous and hence motivates the above approach. More explicitly condition (ii) holds for all t if and only if

$$[AU(t)^{-1}A^{-1}]^*[AU(t)^{-1}A] - I \in \Sigma(H_B)$$

by proposition 2.4, or equivalently, since $U(t)^{-1} = U(-t)$, if and only if $[AU(t)A^{-1}]^*[AU(t)A^{-1}] - I \in \Sigma(H_B)$ for all t . That this is always the case (i.e., independent of K) follows from

Lemma 3.5. $[AU(t)A^{-1}]^*[AU(t)A^{-1}] - I = 0$.

Proof. On a suitable dense subspace of H_B the following formal computations will be valid. Using the results of Lemma 3.4 to compute adjoints on H_B ,

$$\begin{aligned} [AU(t)A^{-1}]^*[AU(t)A^{-1}] &= A^{-1*}U(t)^*A^*AU(t)A^{-1} \\ &= \begin{pmatrix} B^{-\frac{1}{2}}B'^{-\frac{1}{2}}B & 0 \\ 0 & B^{\frac{1}{2}}B'^{\frac{1}{2}}B^{-1} \end{pmatrix} \begin{pmatrix} B^{-1} \cos(tB')B, & -B^{-1} \sin(tB')B'B^{-1} \\ BB'^{-1} \sin(tB')B, & B \cos(tB')B^{-1} \end{pmatrix} \begin{pmatrix} B^{-1}B'^{\frac{1}{2}}B^{\frac{1}{2}} & 0 \\ 0 & BB'^{-\frac{1}{2}}B^{-\frac{1}{2}} \end{pmatrix} \\ &\quad \cdot \begin{pmatrix} B^{-\frac{1}{2}}B'^{\frac{1}{2}} & 0 \\ 0 & B^{\frac{1}{2}}B'^{-\frac{1}{2}} \end{pmatrix} \begin{pmatrix} \cos(tB') & B'^{-1} \sin(tB') \\ -B' \sin(tB') & \cos(tB') \end{pmatrix} \begin{pmatrix} B'^{-\frac{1}{2}}B^{\frac{1}{2}} & 0 \\ 0 & B'^{\frac{1}{2}}B^{-\frac{1}{2}} \end{pmatrix} \\ &= \begin{pmatrix} B^{-\frac{1}{2}} \cos^2(tB')B^{\frac{1}{2}} + B^{-\frac{1}{2}} \sin^2(tB')B^{\frac{1}{2}}, & B^{-\frac{1}{2}}B'^{-\frac{1}{2}}[\sin(tB') \cos(tB') - \sin(tB') \cos(tB')]B'^{\frac{1}{2}}B^{-\frac{1}{2}} \\ B^{\frac{1}{2}}B'^{\frac{1}{2}}[\sin(tB') \cos(tB') & \\ - \sin(tB') \cos(tB')]B'^{-\frac{1}{2}}B^{\frac{1}{2}}, & B^{\frac{1}{2}}B'^{-\frac{1}{2}}[\sin^2(tB') + \cos^2(tB')]B'^{\frac{1}{2}}B^{-\frac{1}{2}} \end{pmatrix} \\ &= \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}. \end{aligned}$$

Thus the operator agrees with the identity on a dense subspace and hence is the identity.

From Lemma 3.5 we conclude that condition (i) alone is sufficient for $U(t)$ to be unitarily implement-

able. The next two lemmas show that this condition is precisely that given in Theorem 3.2.

Lemma 3.6. $A^*A - I \in \Sigma(H_B)$ if and only if $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I \in \Sigma(L^2)$ and $B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}} - I \in \Sigma(L^2)$.

²⁴ Cf. Ref. 3, p. 163, Theorem 6.1.

²⁵ For $T \in Sp(L)$, $\theta^*(T)$ leaves the set of pure regular states invariant.

²⁶ $\Sigma(H)$ denotes the class of Hilbert–Schmidt operators on the Hilbert space H .

Proof. Using the results of Lemma 3.4,

$$A^*A = \begin{pmatrix} B^{-1}B'^{\frac{1}{2}}B^{\frac{1}{2}} & 0 \\ 0 & BB'^{-\frac{1}{2}}B^{-\frac{1}{2}} \end{pmatrix} \begin{pmatrix} B^{-\frac{1}{2}}B'^{\frac{1}{2}} & 0 \\ 0 & B^{\frac{1}{2}}B'^{-\frac{1}{2}} \end{pmatrix} \\ = \begin{pmatrix} B^{-1}B' & 0 \\ 0 & BB'^{-1} \end{pmatrix}.$$

Then, by definition,

$$A^*A - I = \begin{pmatrix} B^{-1}B' - I & 0 \\ 0 & BB'^{-1} - I \end{pmatrix} \in \Sigma(H_B),$$

if and only if

$$\sum_{\mu} \|B^{\frac{1}{2}}(B^{-1}B' - I)e_{\mu}\|^2 \\ + \sum_{\nu} \|B^{-\frac{1}{2}}(BB'^{-1} - I)f_{\nu}\|^2 < \infty,$$

where $\{e_{\mu}\}$ and $\{f_{\nu}\}$ are any CONB in $D[B^{\frac{1}{2}}]$ and $D[B^{-\frac{1}{2}}]$, respectively. Notice, in particular that one may take $\{e_{\mu}\} \subset D(B^{\frac{1}{2}})$ and $\{f_{\nu}\} \subset D(B^{-\frac{1}{2}})$ and that $B^{-\frac{1}{2}}B^{\frac{1}{2}} = I$ on $D(B^{\frac{1}{2}})$ and $B^{\frac{1}{2}}B^{-\frac{1}{2}} = I$ on $D(B^{-\frac{1}{2}})$. Hence the last statement can be written as

$$\sum_{\mu} \|B^{\frac{1}{2}}(B^{-1}B' - I)B^{-\frac{1}{2}}B^{\frac{1}{2}}e_{\mu}\|^2 \\ + \sum_{\nu} \|B^{-\frac{1}{2}}(BB'^{-1})B^{\frac{1}{2}}B^{-\frac{1}{2}}f_{\nu}\|^2 < \infty.$$

Because $\{e_{\mu}\}$ and $\{f_{\nu}\}$ are CONB in $D[B^{\frac{1}{2}}]$ and $D[B^{-\frac{1}{2}}]$ respectively $\{B^{\frac{1}{2}}e_{\mu}\}$ and $\{B^{-\frac{1}{2}}f_{\nu}\}$ are CONB in L^2 . Consider the first pair. The fact that $\{e_{\mu}\}$ is an orthonormal set in $D[B^{\frac{1}{2}}]$ if and only if $\{B^{\frac{1}{2}}e_{\mu}\}$ is orthonormal in L^2 follows from $(B^{\frac{1}{2}}e_{\mu}, B^{\frac{1}{2}}e_{\mu'}) = \delta_{\mu\mu'}$. As for the completeness, suppose $\{e_{\mu}\}$ is complete in $D[B^{\frac{1}{2}}]$ and $f \in L^2$ such that $(f, B^{\frac{1}{2}}e_{\mu}) = 0$ for all μ . Then, because $B^{\frac{1}{2}}B^{-\frac{1}{2}} = I$ on L^2 , $(f, B^{\frac{1}{2}}e_{\mu}) = (B^{\frac{1}{2}}B^{-\frac{1}{2}}f, B^{\frac{1}{2}}e_{\mu}) = 0$ which implies that $B^{-\frac{1}{2}}f = 0$ since $\{e_{\mu}\}$ is complete in $D[B^{\frac{1}{2}}]$. Thus $f = B^{\frac{1}{2}}B^{-\frac{1}{2}}f = 0$ giving the completeness of $\{B^{\frac{1}{2}}e_{\mu}\}$ in L^2 . Conversely, suppose $\{B^{\frac{1}{2}}e_{\mu}\}$ is complete in L^2 and f is an element of $D[B^{\frac{1}{2}}]$ [i.e., $f \in D(B^{\frac{1}{2}})$] such that $(B^{\frac{1}{2}}f, B^{\frac{1}{2}}e_{\mu}) = 0$ for all μ . The completeness of $\{B^{\frac{1}{2}}e_{\mu}\}$ in L^2 implies $B^{\frac{1}{2}}f = 0$ which in turn implies $f = 0$ since $B^{\frac{1}{2}}$ is positive definite. As for the second part of the claim let us define $U = B^{-\frac{1}{2}}: D(B^{-\frac{1}{2}}) \rightarrow L^2$ as a transformation from $D[B^{-\frac{1}{2}}]$ into L^2 with the usual inner product with $D(U) = D(B^{-\frac{1}{2}})$ a dense subspace of $D[B^{-\frac{1}{2}}]$ (with H_B inner product) into $R(U) = R(B^{-\frac{1}{2}}) = D(B^{\frac{1}{2}})$ a dense subspace of L^2 (with the usual inner product). Thus for

$$\varphi_1, \varphi_2 \in D(U) = D(B^{-\frac{1}{2}}), \\ (U\varphi_1, U\varphi_2) = (B^{-\frac{1}{2}}\varphi_1, B^{-\frac{1}{2}}\varphi_2) \equiv (\varphi_1, \varphi_2)_{H_B}.$$

Thus U is a transformation with dense domain and dense range and preserves inner products. Hence U can be extended to an orthogonal transformation $\tilde{U}: D[B^{-\frac{1}{2}}] \rightarrow L^2$ and $\tilde{U}(f_{\nu}) = B^{-\frac{1}{2}}f_{\nu}$. But orthogonal transformations take a CONB into a CONB. Hence $\{f_{\nu}\}$ is a CONB in $D[B^{-\frac{1}{2}}]$ if and only if $\{B^{-\frac{1}{2}}f_{\nu}\}$ is a CONB in L^2 .

Because both terms in the summation written above are positive, it is precisely the statement that

$$B^{\frac{1}{2}}(B^{-1}B' - I)B^{-\frac{1}{2}} = B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I \in \Sigma(L^2),^{27}$$

and

$$B^{-\frac{1}{2}}(BB'^{-1} - I)B^{\frac{1}{2}} = B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}} - I \in \Sigma(L^2),^{27}$$

which completes the proof.

In fact the conditions in the previous lemma can be reduced to the single requirement stated in Theorem 3.2.

Lemma 3.7. $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I \in \Sigma(L^2)$ if and only if $B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}} - I \in \Sigma(L^2)$.

Proof. We first check that $B'^{\frac{1}{2}}B^{-\frac{1}{2}}: L^2 \rightarrow L^2$ is a bounded operator with bounded inverse. Suppose $\varphi \in D(B^{\frac{1}{2}})$ then $\|B'^{\frac{1}{2}}B^{-\frac{1}{2}}\varphi\|^2 = (B'^{\frac{1}{2}}B^{-\frac{1}{2}}\varphi, B'^{\frac{1}{2}}B^{-\frac{1}{2}}\varphi) = (B'B^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}\varphi)$. Now $B' \leq B + \|K\|^{\frac{1}{2}}I$ on $D(B)$ and $B^{-\frac{1}{2}}\varphi \in D(B)$. Hence

$$\|B'^{\frac{1}{2}}B^{-\frac{1}{2}}\varphi\|^2 \leq ([B + \|K\|^{\frac{1}{2}}I]B^{-\frac{1}{2}}\varphi, B^{-\frac{1}{2}}\varphi) \\ = (B^{-\frac{1}{2}}BB^{-\frac{1}{2}}\varphi, \varphi) + \|K\|^{\frac{1}{2}}(B^{-1}\varphi, \varphi) \leq \|\varphi\|^2 \\ + (\|K\|^{\frac{1}{2}}/m)\|\varphi\|^2 = [1 + \|K\|^{\frac{1}{2}}/m]\|\varphi\|^2.$$

Thus $B'^{\frac{1}{2}}B^{-\frac{1}{2}}$ is bounded by $[1 + \|K\|^{\frac{1}{2}}/m]^{\frac{1}{2}}$ on the dense subspace $D(B^{\frac{1}{2}})$ of L^2 and hence by continuity on all of L^2 . Now since both $B'^{\frac{1}{2}}$ and $B^{-\frac{1}{2}}$ are positive definite (hence one-one), the inverse of $B'^{\frac{1}{2}}B^{-\frac{1}{2}}$ exists on $R(B'^{\frac{1}{2}}) = D(B'^{-\frac{1}{2}}) = L^2$. Clearly the inverse is $B^{\frac{1}{2}}B'^{-\frac{1}{2}}$ whose boundedness can be established by essentially the same type of argument. Then by Proposition 2.4,

$$B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I = (B'^{\frac{1}{2}}B^{-\frac{1}{2}})^*(B'^{\frac{1}{2}}B^{-\frac{1}{2}}) - I \in \Sigma(L^2)$$

and

$$B^{\frac{1}{2}}B'^{-1}B^{\frac{1}{2}} - I = ((B'^{\frac{1}{2}}B^{-\frac{1}{2}})^{-1})^*((B'^{\frac{1}{2}}B^{-\frac{1}{2}})^{-1}) \\ - I \in \Sigma(L^2)$$

are equivalent conditions.

By combining Lemmas 3.3–3.7 and the intervening discussions, it is clear that the desired conclusion of Theorem 3.2 has been established. It should be pointed out that all of the conditions on the perturbation K mentioned in the hypotheses of Theorem 3.2 have been used and, in fact, are essential for the development of the argument. Indeed, the fact that K was time-independent permits the construction of the fixed Hilbert space $H_{B'}$ and the corresponding classical state space $(H_{B'}, \Lambda_{B'})$. Then the result in Lemma 3.5 can be anticipated because the action of $A^{-1}U(t)A$ on H_B is essentially the same as that of $U(t)$ on $H_{B'}$ which is orthogonal. Hence the quantizations $x \rightarrow W_0(U(t)x)$, $x \in (H_{B'}, \Lambda_{B'})$ are unitarily equivalent to

²⁷ The equality of these operators follows from the fact that they are bounded and agree on the common dense domain $D(B^{\frac{1}{2}})$.

$x \rightarrow W_0(x)$, $x \in (H_{B'}, \Lambda_{B'})$. Thus the problem reduces to showing that the Fock-Cook representation of $(H_{B'}, \Lambda_{B'})$ and (H_B, Λ_B) are unitarily equivalent which is the condition that $A^*A - I \in \Sigma(H_B)$ since A essentially takes (H_B, Λ_B) into $(H_{B'}, \Lambda_{B'})$. Clearly the case of time-dependent homogeneous perturbations will require a much different analysis. The consideration of equations of this type as well as the more general equation $\square\varphi = m^2\varphi + f(t, \varphi)$ are reserved for a later report. However, preliminary investigations indicate that for the equation $\square\varphi = m^2\varphi + K\varphi + f(t, x)$ where K is as in Theorem 3.2 and f is real-valued, the dynamical propagator (when it exists as a bounded operator on H_B) is unitarily implementable in the Fock representation of (H_B, Λ_B) if the same is true for the propagators of $\square\varphi = m^2\varphi + K\varphi$ and $\square\varphi = m^2\varphi + f(t, x)$. The last condition is that $\|B^{-\frac{1}{2}}f(t, \cdot)\|$ is locally summable as a function of t .

IV. EXAMPLES

In the last section the main problem was reduced to the (nontrivial) problem in linear analysis of determining for which K is $B^{-\frac{1}{2}}B' B^{-\frac{1}{2}} - I \in \Sigma(L^2)$. We now examine the dependence of this condition on the type of spatial region and its dimension. The particular example of a bounded region in one dimension to be considered is the interval $[0, 2\pi]$ where, as previously mentioned, $B^2 = m^2I - \Delta$ is self-adjoint when we take as its domain $\{f \in L^2[0, 2\pi]; f'' \text{ exists, a.e., } f'' \in L^2[0, 2\pi], f(0) = f'(2\pi), f'(0) = f'(2\pi)\}$. For the unbounded region our prototypical example will be the spatial region consisting of the whole real line, with B^2 being the self-adjoint extension of $m^2I - \Delta$ with domain S (i.e., with zero boundary conditions at infinity). For higher dimensions the regions will be taken to be the Cartesian product of the above with B^2 defined analogously so that it is self-adjoint as an operator on the Hilbert space of square integrable functions over the region. The investigation is divided into two subsections, the commutative and noncommutative case, in which the perturbation K does or does not commute with B^2 , respectively.

A. Commutative Case

As anticipated, when K commutes with B^2 the criterion established in Theorem 3.2 can be modified to give a more amenable condition.

Theorem 4.1. For K as in Theorem 3.2 which commutes with B^2 , the conditions $B^{-\frac{1}{2}}B' B^{-\frac{1}{2}} - I \in \Sigma(L^2)$ and $B^{-1}KB^{-1} \in \Sigma(L^2)$ are equivalent.²⁸

²⁸ Notice that again the spatial region and its dimension are left unspecified to emphasize that Theorem 4.1, as Theorem 3.2, is independent of these features.

Proof. As in Lemma 3.7, $B'^{\frac{1}{2}}B^{-2}$ is a bounded operator on L^2 with bounded inverse $B^{\frac{1}{2}}B'^{-\frac{1}{2}}$. Thus, by Proposition 2.4,

$$B^{-\frac{1}{2}}B' B^{-\frac{1}{2}} - I = (B'^{\frac{1}{2}}B^{-\frac{1}{2}})^*(B'^{\frac{1}{2}}B^{-\frac{1}{2}}) - I \in \Sigma(L^2)$$

and

$$([B'^{\frac{1}{2}}B^{-\frac{1}{2}}]^2)^*[B'^{\frac{1}{2}}B^{-\frac{1}{2}}]^2 - I \in \Sigma(L^2)$$

are equivalent statements. But

$$([B'^{\frac{1}{2}}B^{-\frac{1}{2}}]^2)^*[B'^{\frac{1}{2}}B^{-\frac{1}{2}}]^2 - I = B^{-1}B'^2B^{-1} - I = B^{-1}KB^{-1},$$

since K commutes with B^2 (and hence with all functions of B^2).

Theorem 4.2. Suppose Ω_n is the n -fold Cartesian product of $[0, 2\pi]$ and K is a bounded linear positive, time-independent operator on $L^2(\Omega_n)$ which commutes with B^2 . Then the dynamical propagators of the equation $\square\varphi = m^2\varphi + K\varphi$ defined on Ω_n with periodic boundary conditions, are unitarily implementable in the conventional representation of (H_B, Λ_B) if $n \leq 3$.

Proof. Clearly this is precisely the situation for which Theorem 4.1 was designed. Thus $B^{-1}KB^{-1} = KB^{-2} \in \Sigma(L^2(\Omega_n))$ is a sufficient condition for the propagators to be unitarily implementable. Since K is bounded, it is sufficient that $B^{-2} \in \Sigma[L^2(\Omega_n)]$. Now

$$\left\{ \frac{e^{i(k_1x_1 + k_2x_2 + \dots + k_nx_n)}}{(2\pi)^{n/2}} \right\}_{k_1, k_2, \dots, k_n = -\infty}^{\infty}$$

is a CONB for $L^2(\Omega_n)$. Hence $B^{-2} \in \Sigma[L^2(\Omega_n)]$ if and only if

$$\sum_{k_1, k_2, \dots, k_n = -\infty}^{\infty} \left\| B^{-2} \frac{e^{i(k_1x_1 + \dots + k_nx_n)}}{(2\pi)^{n/2}} \right\|^2 = \sum_{k_1, k_2, \dots, k_n = -\infty}^{\infty} \frac{1}{(k_1^2 + k_2^2 + \dots + k_n^2 + m^2)^2} < \infty.$$

The conclusion is then an immediate consequence of the following lemma (i.e., the last sum converges for all n such that $p = 2 > \frac{1}{2}n$ or $n \leq 3$).

Lemma 4.3. If

$m \neq 0$,

$$\sum_{k_1, k_2, \dots, k_n = -\infty}^{\infty} \frac{1}{(k_1^2 + k_2^2 + \dots + k_n^2 + m^2)^p} < \infty$$

if and only if $p > \frac{1}{2}n$.

Outline of the Proof. By applying the integral test to each successive summation it is not difficult to check

that

$$2^n \int_0^\infty \dots \int_0^\infty \frac{d^n x}{(x^2 + m^2)^p}$$

$$\leq \sum_{k_1, k_2, \dots, k_n = -\infty}^\infty \frac{1}{(k_1^2 + k_2^2 + \dots + k_n^2 + m^2)^p}$$

$$\leq 2^n \left[\frac{1}{m^{2p}} + \sum_{k=1}^n \binom{n}{n-k} \int_0^\infty \dots \int_0^\infty \frac{d^k x}{(x^2 + m^2)^p} \right].$$

Thus it is clear that a necessary and sufficient condition for the sum to converge is that the n -dimensional integral

$$\int_0^\infty \dots \int_0^\infty \frac{d^n x}{(x^2 + m^2)^p}$$

be finite. But this integral can be written as

$$C_n \int_0^\infty \frac{r^{n-1} dr}{(r^2 + m^2)^p}$$

where C_n is a finite constant related to the surface area of the unit sphere in n dimensions. The last integral clearly converges if and only if $2p - (n - 1) > 1$ or equivalently $p > \frac{1}{2}n$.

Suppose $K = m'^2 I$, which commutes with B^2 . This would correspond to a cut-off theory for the subjection of a free meson to a mass jump m' at some finite time. Theorem 4.2 then says that the dynamics of this interaction can be completely carried out in the Fock representation. However, for the noncutoff theory this situation never prevails regardless of the number of dimensions. This will be a direct consequence of the next general result.

Theorem 4.4. If K commutes with B^2 as well as satisfying the hypotheses of Theorem 3.2, then a necessary and sufficient condition for the dynamical propagator $U(t)$ of $\square\varphi = m^2\varphi + K\varphi$ to be unitarily implementable in the conventional representation of (H_B, Λ_B) is that $KB^{-2} \sin(tB') \in \Sigma(L^2)$.

Proof. By Corollary 2.3, $U(t)$ is unitarily implementable if and only if $U(t)^*U(t) - I \in \Sigma(H_B)$. Now $U(t)$ is bounded with a bounded inverse $U(-t)$. Thus $U(t)^*U(t) - I \in \Sigma(H_B)$ if and only if $U^*(t) - U(t)^{-1} \in \Sigma(H_B)$. Using the facts that $U(t)^* = JU(t)^{-1}J^{-1} = JU(-t)J^{-1}$ [since $U(t)$ is symplectic] and that B^2 and B'^2 commute, we find that

$$U(t)^* = \begin{pmatrix} \cos(tB'), & -B^{-2}B'^2B'^{-1} \sin(tB') \\ B^2B'^{-2}B' \sin(tB'), & \cos(tB') \end{pmatrix}$$

and

$$U(t)^{-1} = \begin{pmatrix} \cos(tB'), & -B'^{-1} \sin(tB') \\ B' \sin(tB') & \cos(tB') \end{pmatrix}$$

Thus

$$U(t)^* - U(t)^{-1} = \begin{pmatrix} 0 & (I - B^{-2}B'^2)B'^{-1} \sin(tB') \\ (B^2B'^{-2} - I)B' \sin(tB'), & 0 \end{pmatrix}$$

As a result $U(t)^* - U(t)^{-1} \in \Sigma(H_B)$ if and only if

$$\Sigma_v \|B^{\frac{1}{2}}(I - B^{-2}B'^2)B'^{-1} \sin(tB')f_v\|^2 + \Sigma_u \|B^{-\frac{1}{2}}(B^2B'^{-2} - I)B' \sin(tB')e_u\|^2 < \infty$$

for arbitrary CONB $\{e_u\}$ and $\{f_v\}$ in $D[B^{\frac{1}{2}}]$ and $D[B^{-\frac{1}{2}}]$, respectively. As in Lemma 3.6 this is equivalent to the two conditions

$$B^{\frac{1}{2}}(I - B^{-2}B'^2)B'^{-1} \sin(tB')B^{\frac{1}{2}} \in \Sigma(L^2)$$

and

$$B^{-\frac{1}{2}}(B^2B'^{-2} - I)B' \sin(tB')B^{-\frac{1}{2}} \in \Sigma(L^2).$$

The first term,

$$B^{\frac{1}{2}}(I - B^{-2}B'^2)B'^{-1}B^{\frac{1}{2}} \sin(tB') = -BB'^{-1}KB^2 \sin(tB') = -KB^{-1}B'^{-1} \sin(tB')$$

while the second

$$B^{-1}B'(B^2B'^{-2} - I) \sin(tB') = B^{-1}B'^{-1}B'^2(B^2B'^{-2} - I) \sin(tB') = -B^{-1}B'^{-1}K \sin(tB').$$

Thus $U(t)$ is unitarily implementable if and only if $-B^{-1}B'^{-1}K \sin(tB') \in \Sigma(L^2)$ or equivalently

$$B^{-1}B'^{-1}K \sin(tB') \in \Sigma(L^2).$$

However, $B^{-1}B'$ is bounded with bounded inverse on L^2 . Thus the necessary and sufficient condition reduces to

$$B^{-1}B'B^{-1}B'^{-1}K \sin(tB') = B^{-2}K \sin(tB') \in \Sigma(L^2).$$

Corollary 4.5. The dynamical propagators arising from the perturbed Klein-Gordon equation $\square\varphi = m^2\varphi + m'^2\varphi$, $m' > 0$, defined on all of E^n with zero boundary conditions at infinity, are not unitarily implementable in the free field representation associated with the mass m regardless of the number of dimensions.

Proof. By Theorem 4.4, the necessary and sufficient condition for the unitary implementability of the propagators is that $m'^2B^{-2} \sin(tB') \in \Sigma[L^2(E^n)]$. But

$$\mathcal{F}m'^2B^{-2} \sin(tB')\mathcal{F}^{-1} = m'^2(m^2 + k_1^2 + \dots + k_n^2)^{-1} \times \sin[t(k_1^2 + k_2^2 + \dots + k_n^2 + m^2 + m'^2)^{\frac{1}{2}}]$$

where \mathcal{F} is the n -dimensional Fourier transform. Thus $m'^2B^{-2} \sin(tB')$ is unitarily equivalent to multiplication by the above function since \mathcal{F} is unitary on $L^2(E^n)$. But a multiplication operator has

continuous spectrum and unitarily equivalent operators have the same spectrum. Hence $m'^2 B^{-2} \sin(tB')$ has continuous spectrum which precludes its being Hilbert-Schmidt regardless of the dimension.

The last result points out the extreme difficulty of being able to view the dynamics of an interaction on the free representation of the noninteracting field. It also suggests the necessity of a rapid decrease to zero at infinity for the (non-commutative) potential scattering case [i.e., $K = M_v v = v(x_1 \cdots x_n)$].

In concluding this subsection, it might also be pointed out that in the commutative case the sufficient condition of Theorem 3.2 is easily deduced from the necessary and sufficient condition of Theorem 4.4 by using the boundedness of $\sin(tB')$ and Theorem 4.1. The close similarity of the necessary and the sufficient condition in this particular case seems to indicate that a better sufficient condition than that given in Theorem 3.2 is not, in general, available.

B. Noncommutative Case

In this first stage of the investigation of the condition given in Theorem 3.2, we must restrict our attention to the bounded, one-dimensional situation in order to find a noncommutative example for which the dynamical propagators are unitarily implementable. More explicitly,

Theorem 4.6. Suppose K is any bounded, linear, positive, time-independent operator on $L^2[0, 2\pi]$. Then the dynamical propagators of $\square\varphi = m^2\varphi + K\varphi$ (valid for $x \in [0, 2\pi]$, with periodic boundary conditions) are unitarily implementable in the conventional representation of (H_B, Λ_B) .

Proof. Suppose the number of spatial dimensions is left arbitrary. The argument will be based on the analog of the generalized Hölder theorem for non-commutative O_p spaces. $T \in \mathfrak{B}(H)$ is in O_p if $\text{Tr}(|T|^p) < \infty$ where $|T| = (T^*T)^{\frac{1}{2}}$ and Tr is the usual trace of an operator. The result of interest is that if $T_1 \in O_p$ and $T_2 \in O_q$ then $T_1 T_2 \in O_r$ where $1/p + 1/q = 1/r$ for $0 < r < \infty$.²⁹ Clearly $O_2(H) = \Sigma(H)$.

By the above result, to show that $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I = B^{-\frac{1}{2}}(B' - B)B^{-\frac{1}{2}} \in \Sigma[L^2(\Omega_n)]$ it suffices to show that

$$B^{-\frac{1}{2}} \in O_4[L^2(\Omega_n)]$$

and

$$B' - B \in \mathfrak{B}[L^2(\Omega_n)] = O_\infty[L^2(\Omega_n)].$$

For then $B^{-\frac{1}{2}}(B' - B)B^{-\frac{1}{2}} \in O_p$ for $1/p = \frac{1}{4} + 0 + \frac{1}{4} = \frac{1}{2}$ or $p = 2$. Now on $D(B)$, $0 < B \leq B' \leq B +$

$(\|K\|)^{\frac{1}{2}}I$ or $0 \leq B' - B \leq (\|K\|)^{\frac{1}{2}}I$. But $D(B)$ is dense in $L^2(\Omega_n)$, hence $B' - B$ is a bounded operator $L^2(\Omega_n)$. Also, since $B^{-\frac{1}{2}}$ is self-adjoint, $|B^{-\frac{1}{2}}| = B^{-\frac{1}{2}}$ and

$$\text{Tr}(|B^{-\frac{1}{2}}|^4) = \text{Tr}(B^{-2}) = \text{Tr}[(B^2)^{-1}].$$

Using the CONB

$$\left\{ \frac{e^{ik_1 x_1 + \cdots + k_n x_n}}{(2\pi)^{n/2}} \right\}_{k_1, k_2, \dots, k_n = -\infty}^{\infty}$$

it is clear that $B^{-\frac{1}{2}} \in O_4$ if and only if

$$\begin{aligned} &\text{Tr}[(B^2)^{-1}] \\ &= \sum_{k_1, k_2, \dots, k_n = -\infty}^{\infty} \frac{1}{k_1^2 + k_2^2 + \cdots + k_n^2 + m^2} < \infty. \end{aligned}$$

By Lemma 4.3 this sum converges if and only if $\frac{1}{2}n < 1$ or $n = 1$.

Thus the above result gives an example of a non-commutative situation which satisfies the conditions of Theorem 3.2. However, it also points out the difficulties which would arise in attempting to apply this straightforward approach to more natural cases such as bounded and unbounded spatial regions in three dimensions. The above proof has already demonstrated that $B^{-\frac{1}{2}} \in O_4[L^2(\Omega_n)]$ only if $n = 1$. As for the unbounded region, regardless of the dimension, B^2 and hence $B^{-\frac{1}{2}}$ has continuous spectrum and therefore $B^{-\frac{1}{2}} \in O_p$ only for $p = \infty$. But, by proposition 2.4 $B^{-\frac{1}{2}}B'B^{-\frac{1}{2}} - I = B^{-\frac{1}{2}}(B' - B)B^{-\frac{1}{2}} \in \Sigma(L^2)$ if and only if

$$B'^{-\frac{1}{2}}BB'^{-\frac{1}{2}} - I = B'^{-\frac{1}{2}}(B - B')B'^{-\frac{1}{2}} \in \Sigma(L^2).$$

Thus, the above predicament may be avoided by inquiring instead whether $B'^{-\frac{1}{2}} \in O_4(L^2)$. This tact may be profitable in the higher dimensional bounded case with $K = M_v$ (i.e., potential scattering). However, it seems unlikely that much can be gained from this approach in the unbounded case since it can be shown by a result of Kuroda³⁰ that for $K = M_v$, $v \in L^1(E^n) \cap L^2(E^n)$, $n \leq 3$, the absolutely continuous parts of the operators B'^2 and B^2 are unitarily equivalent. Thus the spectrum of B'^2 and hence $B'^{-\frac{1}{2}}$ is continuous which precludes $B'^{-\frac{1}{2}} \in O_4$. It seems inevitable then that a more careful examination of $B' - B$ is needed in order to establish conditions insuring that it is in O_p , $p < \infty$. These problems, as well as necessary conditions for the noncommutative case will be considered in later work. We also hope that the results stated here will be beneficial in indicating approaches to similar considerations for time-dependent and nonlinear perturbations.

²⁹ Cf. N. Dunford and J. T. Schwartz, *Linear Operators, Part II* (Interscience Publishers, Inc., New York, 1963), p. 1093.

³⁰ S. T. Kuroda, *J. Math. Soc. Japan* **11**, 247 (1959); *Nuovo Cimento* **12**, 431 (1959).

Relationships among the Wigner 9j Symbols

S. JANG

Centre de Recherches Nucléaires, Département de Physique Théorique, Strasbourg, France

(Received 17 April 1967)

Several identities satisfied by the 9j symbols or by the product of 6j and 9j symbols are derived by means of the symmetry properties of the Möbius strip type 15j symbol; in particular, the identity

$$\begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} \begin{pmatrix} a & j & m \\ b & k & n \\ c & l & p \end{pmatrix} = \sum_{xyz} (2x+1)(2y+1)(2z+1) \begin{pmatrix} a & d & g \\ x & j & m \end{pmatrix} \begin{pmatrix} b & e & h \\ y & k & n \end{pmatrix} \begin{pmatrix} c & f & i \\ z & l & p \end{pmatrix} \begin{pmatrix} x & d & m \\ y & e & n \\ z & f & p \end{pmatrix} \begin{pmatrix} x & j & g \\ y & k & h \\ z & l & i \end{pmatrix}.$$

The resulting recursion relations for the 9j symbols are also considered.

I. INTRODUCTION

By considering two different coupling schemes of four angular momenta, Biedenharn¹ has obtained the identity satisfied by the Racah coefficients (we henceforth use the equivalent Wigner 6j symbol). In an analogous way, but with five angular-momenta coupling, Arima *et al.*² have arrived at another relation between the 6j and 9j symbols. On the other hand, the problem concerned with the coupling of *N* angular momenta is generally related to the definition of the 3(*N* - 1)*j* symbol. This suggests that, starting from the properties of the 3(*N* - 1)*j* symbol instead of considering different sets of transformations from one specific scheme to another, one may be able to deduce the relations satisfied by the product of the 6j, the 3(*N* - 2)*j*, and consequently the 9j symbols, since the 3(*N* - 2)*j* symbol is always expressible in terms of the 6j and 9j symbols.

Indeed, when both sides of the following symmetry relation of the 9j symbol

$$\begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = (-1)^S \begin{pmatrix} a & b & c \\ g & h & i \\ d & e & f \end{pmatrix}$$

(where *S* is the sum of all the nine arguments) are multiplied by the quantity

$$(2c+1) \begin{pmatrix} a & b & c \\ f & i & k \end{pmatrix}$$

and a summation is then carried out over the argument *c*, one can readily obtain the Biedenharn identity.

Thus,

$$\begin{aligned} & (-1)^{2k} \begin{pmatrix} d & e & f \\ b & k & h \end{pmatrix} \begin{pmatrix} g & h & i \\ k & a & d \end{pmatrix} \\ &= \sum_{xc} (-1)^{S+k+c+x} (2x+1)(2c+1) \\ & \quad \times \begin{pmatrix} a & i & k \\ b & f & x \end{pmatrix} \begin{pmatrix} a & b & c \\ i & f & x \end{pmatrix} \begin{pmatrix} a & b & c \\ g & h & i \\ d & e & f \end{pmatrix} \\ &= \sum_x (-1)^{R+x+2a+2i} (2x+1) \\ & \quad \times \begin{pmatrix} a & i & k \\ b & f & x \end{pmatrix} \begin{pmatrix} g & h & i \\ b & x & e \end{pmatrix} \begin{pmatrix} d & e & f \\ x & a & g \end{pmatrix}, \quad (1) \end{aligned}$$

where *R* = *S* - *c* + *k* and use was made of the orthogonality condition and the sum rule (Racah's back-coupling rule) for the 6j symbols.

The subject of this note is, therefore, to deduce the corresponding identity and other related formulas for the 9j symbols using the symmetry properties of the 3(6 - 1)*j* symbol and also to derive the resulting recursion relations. It is noted that one can equally obtain the desired relations by means of the different coupling schemes without relying on the known relationships for the 6j and 9j symbols, but this procedure is not advantageous when the number of coupled angular momenta is large, since it is not always straightforward to express the transformation matrices of the eigenfunctions of *N* coupled angular momenta in terms of known simpler matrices.³

³ See, for example, A. P. Yutsis, I. B. Levinson, and V. V. Vanagas, *The Theory of Angular Momentum [Mathematicheskii apparat teorii momenta kolichestva dvizheniya]*, (Vilnius, USSR, (1960), translated from Russian by A. Sen and A. R. Sen (Israel Program for Scientific Translation, Jerusalem, Israel, 1962). They give an extensive account of the transformation matrices of five angular-momenta coupling with general consideration of the 3*Nj* symbols.

¹ L. C. Biedenharn, *J. Math. and Phys.* **31**, 287 (1953).
² A. Arima, H. Horie, and Y. Tanabe, *Progr. Theoret. Phys. (Kyoto)* **11**, 143 (1954).

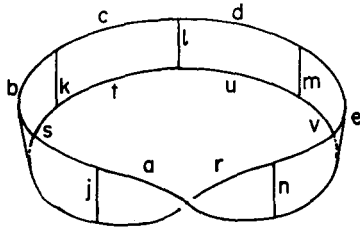


FIG. 1. Möbius strip which represents the triangular conditions of the 15j symbol.

II. IDENTITIES

A natural extension of the arguments on the 12j symbol, given by Jahn and Hope⁴ and Ord-Smith,⁵ to the 15j symbol permits us to easily write down twenty symmetry relations (see Fig. 1). It is, however, noted that the 15j symbols are not unique and can also be defined differently⁶ by another coupling manner between the fifteen arguments, and that only the 15j symbol which is represented by a Möbius strip will be considered. Among the twenty symmetry forms, five of them are given explicitly and it is remarked that the other symmetry forms reduce essentially to one of these when the 15j symbol is expressed in terms of the 9j symbols. Thus, in the notation of Ord-Smith,⁵

$$\begin{aligned} \begin{Bmatrix} a & b & c & d & e \\ j & k & l & m & n \\ r & s & t & u & v \end{Bmatrix} &= \dots \\ &= \begin{Bmatrix} b & c & d & e & r \\ k & l & m & n & j \\ s & t & u & v & a \end{Bmatrix} = \dots \\ &= \begin{Bmatrix} c & b & a & v & u \\ k & j & n & m & l \\ t & s & r & e & d \end{Bmatrix} = \dots \\ &= \begin{Bmatrix} d & c & b & a & v \\ l & k & j & n & m \\ u & t & s & r & e \end{Bmatrix} = \dots \\ &= \begin{Bmatrix} e & r & s & t & u \\ n & j & k & l & m \\ v & a & b & c & d \end{Bmatrix}. \quad (2) \end{aligned}$$

The following fundamental property of the 15j symbol can be readily shown by extending the results

⁴ H. A. Jahn and J. Hope, Phys. Rev. 93, 318 (1954).

⁵ R. J. Ord-Smith, Phys. Rev. 94, 1227 (1954).

⁶ When *N* is greater than four, there exist other less symmetric symbols besides the two 3*Nj* symbols represented by a Möbius strip and an untwisted cylindrical band, respectively (see, e.g., Ref. 3).

for the 12j symbol:

$$\begin{aligned} &\begin{Bmatrix} a & b & c & d & e \\ j & k & l & m & n \\ r & s & t & u & v \end{Bmatrix} \\ &= (-1)^S \sum_z (-1)^{2z} (2z + 1) \\ &\quad \times \begin{Bmatrix} a & b & z \\ j & k & z \\ s & r & z \end{Bmatrix} \begin{Bmatrix} b & c & z \\ k & l & z \\ t & s & z \end{Bmatrix} \begin{Bmatrix} c & d & z \\ l & m & z \\ u & t & z \end{Bmatrix} \\ &\quad \times \begin{Bmatrix} d & e & z \\ m & n & z \\ v & u & z \end{Bmatrix} \begin{Bmatrix} e & r & z \\ n & s & z \\ a & v & z \end{Bmatrix} \\ &= \sum_{xy} (-1)^{a-r+t-e-v} (2x + 1)(2y + 1) \begin{Bmatrix} a & b & c \\ j & k & x \\ r & s & t \end{Bmatrix} \\ &\quad \times \begin{Bmatrix} c & d & e \\ l & m & y \\ t & u & v \end{Bmatrix} \begin{Bmatrix} a & t & e \\ x & y & n \\ r & c & v \end{Bmatrix}, \quad (3) \end{aligned}$$

where *S* is the sum of all the fifteen arguments. It is noted that, for the purpose of symmetrical presentation, all the notations are expressed in those of Möbius strip symmetry and these notations are related to the Wigner 6j and 9j symbols by³

$$\begin{Bmatrix} a & b \\ j & k \\ r & s \end{Bmatrix} = W(absr; jk) = (-1)^{a+b+r+s} \begin{Bmatrix} a & b & j \\ r & s & k \end{Bmatrix}$$

and

$$\begin{Bmatrix} a & b & c \\ j & k & l \\ r & s & t \end{Bmatrix} = \begin{Bmatrix} j & a & b \\ r & l & c \\ s & t & k \end{Bmatrix},$$

where *W* is the Racah coefficients. Other properties⁷ of the 15j symbol follow from Eq. (3).

When two of the five 15j symbols of the symmetry relations (2) are equated in terms of the 9j symbols, it is seen that there are two types of equations: one of the 9j symbols on each side has the same form but the

⁷ The Biedenharn identity, applied to the last two 6j symbols of Eq. (3), yields

$$\sum_w (-1)^{w+z-e-m-n-v} \begin{Bmatrix} d & r \\ w & z \\ a & u \end{Bmatrix} \begin{Bmatrix} a & n \\ v & w \\ m & u \end{Bmatrix} \begin{Bmatrix} r & n \\ e & w \\ m & d \end{Bmatrix} (2w + 1).$$

Then, it is seen immediately that the summation of four 6j symbols over *z* is just the definition of the 12j symbol (see Ref. 5). Therefore, the 15j symbol of Eq. (3) is related to the 12j symbol by

$$\sum_w (2w + 1) \begin{Bmatrix} a & b & c & d \\ j & k & l & w \\ r & s & t & u \end{Bmatrix} \begin{Bmatrix} a & n \\ v & w \\ m & u \end{Bmatrix} \begin{Bmatrix} r & n \\ e & w \\ m & d \end{Bmatrix}.$$

Thus, a usefulness of the Biedenharn identity in the study of the 3*Nj* symbol is demonstrated. Yutsis *et al.* (Ref. 3) have shown in a different and general manner that this kind of relation holds between the 3*Nj* and 3(*N* - 1)*j* symbols.

other $9j$ symbols are different, and every six $9j$ symbols are different from each other. Explicitly, this can be characterized by the equation

$$\begin{aligned} & \sum_{xy} (-1)^{c-t+u-d} (2x+1)(2y+1) \\ & \quad \times \begin{Bmatrix} a & b & j \\ x & c & r \\ t & k & s \end{Bmatrix} \begin{Bmatrix} a & d & y \\ x & c & r \\ t & l & u \end{Bmatrix} \begin{Bmatrix} a & d & y \\ n & e & r \\ v & m & u \end{Bmatrix} \\ & = \sum_{xz} (-1)^{a-r+e-v} (2x+1)(2z+1) \\ & \quad \times \begin{Bmatrix} a & b & j \\ x & c & r \\ t & k & s \end{Bmatrix} \begin{Bmatrix} a & v & n \\ x & c & r \\ t & z & e \end{Bmatrix} \begin{Bmatrix} u & v & m \\ l & c & d \\ t & z & e \end{Bmatrix} \\ & = \sum_{pq} (-1)^{d-u+v-e} (2p+1)(2q+1) \\ & \quad \times \begin{Bmatrix} a & b & j \\ n & e & r \\ v & q & s \end{Bmatrix} \begin{Bmatrix} u & k & p \\ m & e & d \\ v & q & s \end{Bmatrix} \begin{Bmatrix} u & b & p \\ l & c & d \\ t & k & s \end{Bmatrix}, \quad (4) \end{aligned}$$

where the notation of the Wigner $9j$ symbol is re-introduced and the phase symmetries have not been taken into account. In the first part of the equation, the same $9j$ symbols on each side can be eliminated using the orthogonality of the $9j$ symbols. This yields

$$\begin{aligned} & \sum_x (2x+1) \begin{Bmatrix} a & b & x \\ d & e & f \\ g & h & i \end{Bmatrix} \begin{Bmatrix} a & b & x \\ j & k & f \\ l & m & i \end{Bmatrix} \\ & = \sum_y (2y+1) \begin{Bmatrix} m & d & y \\ b & e & h \\ k & f & j \end{Bmatrix} \begin{Bmatrix} m & d & y \\ i & g & h \\ l & a & j \end{Bmatrix}, \quad (5) \end{aligned}$$

which shows that the common arguments in the two $9j$ symbols on one side appear only once on the other side. When one of the common arguments on one side vanishes, this equation leads to the well-known relation Eq. (A2). Equation (5) can be rewritten in the form of the sum rule:

$$\begin{aligned} & \begin{Bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{Bmatrix} = \sum_{xyz} [a][x][y][z] \\ & \quad \times \begin{Bmatrix} a & b & c \\ j & y & f \\ l & z & i \end{Bmatrix} \begin{Bmatrix} z & b & y \\ d & e & f \\ x & h & j \end{Bmatrix} \begin{Bmatrix} a & j & l \\ d & x & z \\ g & h & i \end{Bmatrix}, \quad (6) \end{aligned}$$

where $[a] = (2a+1)$, etc. It is also possible to derive the analogous relations with $[b]$, $[f]$, or $[i]$ in place of $[a]$, but these expressions turn out to be the original one because of the symmetry character of the $9j$ symbols.

A procedure for carrying out the summation on one side of Eq. (4) is to sum over the arguments which

appear only once, namely over j and m first applying the relationship of the type (A2) on the left-hand side and then over s and e using the orthogonalities for the $6j$ symbols. This yields a simple product of two $6j$ symbols and one $9j$ symbol with fifteen distinct arguments. Now, the problem consists essentially of reducing the number of summation indices on the other side and this can be readily effected in an analogous way with the derivation of the identity (1). Thus, the following identity with twelve distinct arguments is immediately obtained:

$$\begin{aligned} & (-1)^{d+h-b-f} \begin{Bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{Bmatrix} \begin{Bmatrix} g & h & i \\ j & k & l \end{Bmatrix} \\ & = \sum_{xy} (-1)^{c+l-a-x} [x][y] \begin{Bmatrix} c & f & i \\ j & k & x \end{Bmatrix} \begin{Bmatrix} d & e & f \\ j & x & y \end{Bmatrix} \\ & \quad \times \begin{Bmatrix} b & e & h \\ j & l & y \end{Bmatrix} \begin{Bmatrix} a & b & c \\ d & y & x \\ g & l & k \end{Bmatrix}. \quad (7) \end{aligned}$$

This identity and that of (6) were originally given by Arima *et al.*² from the transformation matrices of the different coupling schemes. Another process used for reducing the multiplicity of the sum in which the fifteen distinct arguments are conserved leads to

$$\begin{aligned} & \begin{Bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{Bmatrix} \begin{Bmatrix} d & e & f \\ j & k & l \end{Bmatrix} \begin{Bmatrix} g & h & i \\ m & n & p \end{Bmatrix} \\ & = \sum_{xyz} [x][y][z] \begin{Bmatrix} a & b & c \\ x & y & z \end{Bmatrix} \begin{Bmatrix} a & y & z \\ d & k & l \\ g & n & p \end{Bmatrix} \\ & \quad \times \begin{Bmatrix} b & z & x \\ e & l & j \\ h & p & m \end{Bmatrix} \begin{Bmatrix} c & x & y \\ f & j & k \\ i & m & n \end{Bmatrix}. \quad (8) \end{aligned}$$

Finally, a product of two $9j$ symbols which have three common arguments is expressed in terms of the multiple sum of the $6j$ and $9j$ symbols:

$$\begin{aligned} & \begin{Bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{Bmatrix} \begin{Bmatrix} a & j & m \\ b & k & n \\ c & l & p \end{Bmatrix} \\ & = \sum_{xyz} [x][y][z] \begin{Bmatrix} a & d & g \\ x & j & m \end{Bmatrix} \begin{Bmatrix} b & e & h \\ y & k & n \end{Bmatrix} \begin{Bmatrix} c & f & i \\ z & l & p \end{Bmatrix} \\ & \quad \times \begin{Bmatrix} x & d & m \\ y & e & n \\ z & f & p \end{Bmatrix} \begin{Bmatrix} x & j & g \\ y & k & h \\ z & l & i \end{Bmatrix}. \quad (9) \end{aligned}$$

This identity, in a very symmetrical form, follows from Eq. (8) when the orthogonalities of the $9j$ and $6j$

symbols with respect, for example, to the arguments d, g , and to the argument a , are taken into account, respectively, but can also be proved directly by reintroducing the identity into the two $9j$ symbols of the right-hand side. When one of the twelve arguments which are not common on the left-hand side vanishes, the identity (9) reduces to that of (7) which in turn becomes Eq. (A2) if the arguments e or f are zero.

III. RECURSION FORMULAS

Contrary to the cases of the $3j$ and $6j$ symbols, almost no account has been given for the recursion

formulas of the $9j$ symbols. Because of the nine arguments involved, simple manageable algebraic expressions between $9j$ symbols are not available. In special cases, however, it is possible to derive some practical recursion formulas from relatively simple relationships for the $9j$ symbols and this is done in Appendix A.

The general recursion relations so far as two or three arguments are concerned can be derived from the identity (9). Thus, for example, when $j = k = \frac{1}{2}$ and $l = 0$, then $p = c, x = g \pm \frac{1}{2}$, and $y = h \pm \frac{1}{2}$. This leads to a recursion formula of the type

$$\begin{aligned}
 & (2g + 1)(2h + 1)[(S - 2|u - v|a + |u + v| + v + \frac{1}{2})(S - 2|u - v|b - 2|u + v|c + u + \frac{1}{2})]^{\frac{1}{2}} \begin{pmatrix} a & d & g \\ b & e & h \\ c & f & i \end{pmatrix} \\
 & = \{[d + g - u(2a + 1) + \frac{1}{2}]\{a + 2u(d - g) + u + \frac{1}{2}\}\}^{\frac{1}{2}} \\
 & \times \left((-1)^{\delta_1} [(R + 1)(R - 2i)\{e + h - v(2b + 1) + \frac{1}{2}\}\{b + 2v(e - h) + v + \frac{1}{2}\}]^{\frac{1}{2}} \begin{pmatrix} a + u & d & g - \frac{1}{2} \\ b + v & e & h - \frac{1}{2} \\ c & f & i \end{pmatrix} \right. \\
 & + (-1)^{\delta_2} [(R - 2h)(R - 2g + 1)\{e + h + v(2b + 1) + \frac{3}{2}\}\{b + 2v(h - e) + v + \frac{1}{2}\}]^{\frac{1}{2}} \left. \begin{pmatrix} a + u & d & g - \frac{1}{2} \\ b + v & e & h + \frac{1}{2} \\ c & f & i \end{pmatrix} \right) \\
 & + \{[a + d + 2u(g + 1) + 1]\{g + 2u(a - d) + 1\}\}^{\frac{1}{2}} \\
 & \times \left((-1)^{\delta_3} [(R - 2g)(R - 2h + 1)\{e + h - v(2b + 1) + \frac{1}{2}\}\{b + 2v(e - h) + v + \frac{1}{2}\}]^{\frac{1}{2}} \begin{pmatrix} a + u & d & g + \frac{1}{2} \\ b + v & e & h - \frac{1}{2} \\ c & f & i \end{pmatrix} \right. \\
 & + (-1)^{\delta_4} [(R + 2)(R - 2i + 1)\{e + h + v(2b + 1) + \frac{3}{2}\}\{b + 2v(h - e) + v + \frac{1}{2}\}]^{\frac{1}{2}} \left. \begin{pmatrix} a + u & d & g + \frac{1}{2} \\ b + v & e & h + \frac{1}{2} \\ c & f & i \end{pmatrix} \right), \tag{10}
 \end{aligned}$$

where the arguments u and v can take the values $\pm \frac{1}{2}$ independently and $\delta_1 = \frac{1}{2}(1 + u - v - |u + v|)$, $\delta_2 = \frac{1}{2}(1 - u - v - |u - v|)$, $\delta_3 = \frac{1}{2}(1 + u + v - |u - v|)$, $\delta_4 = \frac{1}{2}(1 - u + v - |u + v|)$, $S = a + b + c$, and $R = g + h + i$. It is noted that this formula is the simplest general recursion relation which can be deduced from the identity (9) and which does not contain constant arguments. Other recursion formulas are also derived from (9) for the following cases:

- $j = k = \frac{1}{2}, l = 1, m = a \pm \frac{1}{2}, n = b \pm \frac{1}{2}$, and $p = c$;
- $j = k = 1, l = 0, m = a, n = b$, and $p = c$;
- $j = k = l = 1, m = a, n = b$, and $p = c \pm 1$; etc.

In Appendix B, the explicit recursion formulas for the second case is shown.

ACKNOWLEDGMENTS

The author wishes to express his gratitude to Professor S. Gorodetzky and to Professor R. Armbruster for making it possible for him to carry out this work at Strasbourg. He would like to thank Professor J. Yoccoz and Professor G. Monsonego for their enlightening discussions and also to Dr. E. Cooperman for critically reading the manuscript.

APPENDIX A. SIMPLE RECURSION FORMULAS

As was seen in the last paragraph, the general recursion relations for the $9j$ symbols are practically of no use for the real calculation of the $9j$ values. However, the simple practical algebraic relations are obtainable when

one of the arguments is a constant. One such formula⁸ is

$$2[s(s+1)L(L+1)]^{\frac{1}{2}} \begin{Bmatrix} s & s & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L \end{Bmatrix} = [l_1(l_1+1) + j_2(j_2+1) - j_1(j_1+1) - l_2(l_2+1)] \begin{Bmatrix} s & s & 0 \\ l_1 & l_2 & L \\ j_1 & j_2 & L \end{Bmatrix}. \quad (\text{A1})$$

The proof of this formula is tedious and can be done, for example, by writing out explicitly the left-hand side $9j$ symbol in terms of the $6j$ symbols so as to contain the arguments $j_2 \pm 1$, j_2 and applying the recursion relation of the $6j$ symbol to the term with $j_2 + 1$.

Most of the recursion formulas among the $9j$ symbols which contain a constant argument can be derived from

$$(-1)^{b+a+f+h+2k} ([f][h])^{\frac{1}{2}} \begin{Bmatrix} g & h & i \\ k & a & d \end{Bmatrix} \begin{Bmatrix} k & b & f \\ d & e & f \\ h & h & 0 \end{Bmatrix} = \sum_x [x] \begin{Bmatrix} a & b & x \\ f & i & k \end{Bmatrix} \begin{Bmatrix} a & b & x \\ d & e & f \\ g & h & i \end{Bmatrix}, \quad (\text{A2})$$

which is useful when the value of the argument i is not greater than one and the sum is thus limited to a few terms. Combined with the relation (A1), a recursion formula useful in the study of nuclear spectroscopy is obtained:

$$\begin{aligned} & [L(2L+3)(S+2)(S-2j_1+1)(S-2j_2+1)(S-2L)]^{\frac{1}{2}} \begin{Bmatrix} s & s & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L+1 \end{Bmatrix} \\ & = (2L+1) \left[\bar{j}_1 - \bar{j}_2 - L + \frac{2L(\bar{j}_2 + \bar{s} - \bar{l}_2)}{\bar{l}_1 + \bar{j}_2 - \bar{j}_1 - \bar{l}_2} \right] \begin{Bmatrix} s & s & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L \end{Bmatrix} \\ & + [(L+1)(2L-1)(S+1)(S-2j_1)(S-2j_2)(S-2L+1)]^{\frac{1}{2}} \begin{Bmatrix} s & s & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L-1 \end{Bmatrix}, \quad (\text{A3}) \end{aligned}$$

where $S = j_1 + j_2 + L$ and the convention of Biedenharn *et al.*⁹ is introduced (i.e., $\bar{L} = L(L+1)$, etc.). It is observed that this equation is indeterminate if $l_1 = l_2$ and $j_1 = j_2$ or $l_1 = j_1$ and $l_2 = j_2$. In this case, the use of the formula (A1) leads to

$$\begin{aligned} & [s(s+1)(L+1)(2L+3)(2j+L+2)(2j-L)]^{\frac{1}{2}} \begin{Bmatrix} s & s & 1 \\ l & l & L \\ j & j & L+1 \end{Bmatrix} \\ & = [s(s+1)L(2L-1)(2j+L+1)(2j-L+1)]^{\frac{1}{2}} \begin{Bmatrix} s & s & 1 \\ l & l & L \\ j & j & L-1 \end{Bmatrix} \\ & + (2L+1)[j(j+1) + s(s+1) - l(l+1)] \begin{Bmatrix} s & s & 0 \\ l & l & L \\ j & j & L \end{Bmatrix}. \quad (\text{A4}) \end{aligned}$$

When $s = \frac{1}{2}$ in Eq. (A3), the first $9j$ symbol on the right-hand side can be eliminated, yielding the following

⁸ This relation was often quoted incorrectly in the literature: M. Rotenberg, R. Bivins, N. Metropolis, and J. K. Wooten, Jr., *The 3-j and 6-j Symbols* (M.I.T. Press, Cambridge, Massachusetts, 1959), p. 24; A. de-Shalit and I. Talmi, *Nuclear Shell Theory* (Academic Press Inc., New York, 1963), p. 520. The formula (A1) differs from those of these references by the factor $[(2s+1)(2L+1)]^{\frac{1}{2}}$ and also by the phase for the latter reference.

⁹ L. C. Biedenharn, J. M. Blatt, and M. E. Rose, *Rev. Mod. Phys.* **24**, 249 (1952).

expression:

$$\begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L+1 \end{pmatrix} = \left[\frac{2L+1}{L+(l_1-j_1)(2l_1+1)+(l_2-j_2)(2l_2+1)} - 1 \right] \\ \times \left[\frac{L(2L-1)(S-2L)(S-2j_1+1)(S-2j_2+1)(S+2)}{(L+1)(2L+3)(S-2L+1)(S-2j_1)(S-2j_2)(S+1)} \right]^{\frac{1}{2}} \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 1 \\ l_1 & l_2 & L \\ j_1 & j_2 & L-1 \end{pmatrix} \quad (A5)$$

with $S = j_1 + j_2 + L$. It is also possible to relate two $9j$ symbols involving the arguments $L + 1$ and L , respectively, and this can be done by introducing the expression (A5) into (A3).

APPENDIX B

The second example of the general recursion formula is the case with $j = k = 1, l = 0, m = a, n = b$, and $p = c$ in the identity (9). This yields

$$\sum_{v=0, w=0}^2 \begin{pmatrix} a & b-v & c-w \\ d & e & f \\ g & h & i \end{pmatrix} (-1)^v \cdot (b-1+\frac{1}{2}v)(c-1+\frac{1}{2}w) \cdot A(b-v, c-w) \cdot B(b-v) \cdot \Gamma(c-w) = 0, \quad (B1)$$

where

$$\begin{aligned} A(b-v, c-w) &= [(S-v)(S-v+1)(S-2a-v)(S-2a-v-1)]^{\frac{1}{2}} \text{ for } v = w = 0 \text{ or } 2 \\ &= [S(S-2a-1)(S-2b+v)(S-2c+w)]^{\frac{1}{2}} \text{ for } (v=0, w=1) \text{ or } (v=1, w=0) \\ &= -[(S-2b+v-1)(S-2b+v)(S-2c+w-1)(S-2c+w)]^{\frac{1}{2}} \\ &\hspace{15em} \text{for } (v=0, w=2) \text{ or } (v=2, w=0) \\ &= [(S-1)(S-2a-2)(S-2b+v-1)(S-2c+w-1)]^{\frac{1}{2}} \\ &\hspace{15em} \text{for } (v=1, w=2) \text{ or } (v=2, w=1) \\ &= \frac{16b(b-1)c(c-1)[e(e+1)+f(f+1)-d(d+1)]}{[B(b-1)][\Gamma(c-1)]} + a(a+1) - b(b-1) - c(c-1) \\ &\hspace{15em} \text{for } (v=1, w=1); \end{aligned}$$

$$\begin{aligned} B(b-v) &= [(R-\frac{1}{2}v+1)(R-2b+\frac{1}{2}v+1)(R-2e-\frac{1}{2}v)(R-2h-\frac{1}{2}v)]^{\frac{1}{2}} \text{ for } v = 0 \text{ or } 2 \\ &= 2[b(b-1)+e(e+1)-h(h+1)] \text{ for } v = 1; \end{aligned}$$

$$\begin{aligned} \Gamma(c-w) &= [(T-\frac{1}{2}w+1)(T-2c+\frac{1}{2}w+1)(T-2f-\frac{1}{2}w)(T-2i-\frac{1}{2}w)]^{\frac{1}{2}} \text{ for } w = 0 \text{ or } 2 \\ &= 2[c(c-1)+f(f+1)-i(i+1)] \text{ for } w = 1, \end{aligned}$$

with

$$S = a + b + c, \quad R = b + e + h \quad \text{and} \quad T = c + f + i.$$

Embeddings of the Plane-Fronted Waves and Other Space-Times

C. D. COLLINSON

Department of Applied Mathematics, The University of Hull, England

(Received 20 June 1967)

The plane-fronted waves of general relativity are embedded in a six-dimensional pseudo-Euclidean space of signature -2 . Two distinct families of embeddings are found. Embeddings of several well-known space-times are obtained. Certain results of J. Rosen are improved.

1. INTRODUCTION

In a recent paper by Rosen,¹ several well-known space-times are embedded, locally and isometrically, in pseudo-Euclidean space. One aim of the present paper is to find embedding spaces of lower dimension. In particular, the plane-fronted gravitational waves^{2,3} are embedded systematically by analysing the Gauss-Codazzi-Ricci⁴ equations. Two distinct families of embeddings of the plane-fronted gravitational waves in a six-dimensional pseudo-Euclidean space of signature -2 are obtained. These embeddings are independent of the vacuum field equations and so apply to the general plane-fronted waves.

The embedding classes of all but two of the space-times discussed by Rosen are found (the embedding class of a space-time is $p - 4$, where p is the least possible dimension of the embedding space). In order to deduce the embedding classes of several of the space-times, the following theorem is used.

Theorem 1: Solutions of the Einstein-Maxwell field equations can be embedded (locally and isometrically) in a five-dimensional pseudo-Euclidean space only if the electromagnetic field and the Weyl tensor are *both* null. The proof of this theorem is given elsewhere.⁵

Throughout this paper, Greek letters $\alpha, \beta, \gamma, \dots$ denote *tensor* indices while Roman letters m, n, p, \dots denote *tetrad* indices. A semicolon denotes an *intrinsic* derivative and γ_{mnp} are the usual Ricci rotation coefficients.⁴ A space-time can be embedded locally and isometrically in a six-dimensional pseudo-Euclidean space if and only if there exist two symmetric tensors a_{mn}, b_{mn} , and a vector s_m satisfying the following equations.^{4,6}

Gauss equation:

$$R_{mnpq} = 2e_1 a_{m[p} a_{q]n} + 2e_2 b_{m[p} b_{q]n}. \quad (1.1)$$

¹ J. Rosen, *Rev. Mod. Phys.* **37**, 204 (1965).

² H. W. Brinkman, *Math. Ann.* **94**, 119 (1925).

³ W. Kundt, *Z. Physik* **163**, 77 (1961).

⁴ L. P. Eisenhart, *Riemannian Geometry* (Princeton University Press, Princeton, N.J., 1925), p. 97.

⁵ C. D. Collinson, *Commun. Math. Phys.* (to be published).

⁶ C. D. Collinson, *J. Math. Phys.* **7**, 608 (1966).

Codazzi equations:

$$a_{m[n;p]} - \gamma_{[n^p]} a_{m\alpha} + \gamma_{m^q[n} a_{p]\alpha} = e_2 s_{[n} b_{p]m}, \quad (1.2)$$

$$b_{m[n;p]} - \gamma_{[n^p]} b_{m\alpha} + \gamma_{m^q[n} b_{p]\alpha} = -e_1 s_{[n} a_{p]m}. \quad (1.3)$$

Ricci equation:

$$s_{[m;n]} - s_a \gamma_{[m^a n]} + a_{a[m} b_{n]}^a = 0. \quad (1.4)$$

In the above, R_{mnpq} is the curvature tensor of the space-time, e_1 and e_2 are real constants of unit modulus, and square brackets denote antisymmetrization. These equations are, in fact, the integrability conditions of the differential equations

$$y_{;mn}^Q - y_{;p}^Q \gamma_m^p n = e_1 a_{mn} \eta_1^Q + e_2 b_{mn} \eta_2^Q, \quad (1.5)$$

$$\eta_{1;m}^Q = -a_{pm} y_{;}^Q - e_2 s_m \eta_2^Q, \quad (1.6)$$

and

$$\eta_{2;m}^Q = -b_{pm} y_{;}^Q + e_1 s_m \eta_1^Q, \quad (1.7)$$

where capital indices Q, R, \dots range and sum from 1 to 6. Here η_1^Q and η_2^Q are two vectors (in the embedding space) normal to the space-time and y^Q are coordinates in the embedding space. The normals η_1^Q and η_2^Q are not uniquely determined. If $e_1 = e_2$, a new set of normals can be chosen satisfying

$$\bar{\eta}_1^Q = \cos \theta \eta_1^Q + \sin \theta \eta_2^Q \quad (1.8)$$

and

$$\bar{\eta}_2^Q = \sin \theta \eta_1^Q + \cos \theta \eta_2^Q.$$

This induces the transformations

$$\bar{a}_{mn} = \cos \theta a_{mn} + \sin \theta b_{mn},$$

$$\bar{b}_{mn} = -\sin \theta a_{mn} + \cos \theta b_{mn},$$

and

$$\bar{s}_m = s_m - e_2 \theta_{;m}.$$

Similarly, if $e_1 = -e_2$, a new set of normals can be chosen satisfying

$$\bar{\eta}_1^Q = \cosh \theta \eta_1^Q + \sinh \theta \eta_2^Q \quad (1.9)$$

and

$$\bar{\eta}_2^Q = \sinh \theta \eta_1^Q + \cosh \theta \eta_2^Q.$$

This induces the transformations

$$\bar{a}_{mn} = \cosh \theta a_{mn} + \sinh \theta b_{mn},$$

$$\bar{b}_{mn} = \sinh \theta a_{mn} + \cosh \theta b_{mn},$$

and

$$\bar{s}_m = s_m - e_2 \theta_{;m}.$$

2. PLANE-FRONTED WAVES

The metric of the plane-fronted waves can be written, in terms of two real coordinates ρ, σ , and a complex coordinate ζ , in the form

$$ds^2 = 2 d\rho d\sigma - 2H d\sigma^2 - 2 d\zeta d\bar{\zeta}. \quad (2.1)$$

The conditions for a purely gravitational wave are

$$\frac{\partial H}{\partial \rho} = \frac{\partial^2 H}{\partial \zeta \partial \bar{\zeta}} = 0. \quad (2.2)$$

These plane-fronted gravitational waves are characterized by any one of the following properties:

- i. the existence of a covariant constant vector field²;
- ii. type N conformal tensor with shear-free, non-diverging rays³;
- iii. type N conformal tensor whose rays are trajectories to a one-parameter group of affine collineations⁷;
- iv. Einstein field which is mapped conformally onto another Einstein field.²

The tetrad of vectors $l^\alpha, n^\alpha, m^\alpha$, and \bar{m}^α , defined by

$$l^\alpha = \delta_1^\alpha, \quad n^\alpha = \delta_2^\alpha + H\delta_1^\alpha, \quad m^\alpha = \delta_0^\alpha,$$

where $(x^1, x^2, x^0, x^{\bar{0}}) \equiv (\rho, \sigma, \zeta, \bar{\zeta})$, satisfy the orthornormality conditions

$$l^\alpha n_\alpha = -m^\alpha \bar{m}_\alpha = 1,$$

all other contractions being zero. Such a tetrad is a null tetrad in the sense of Newman and Penrose.⁸ The intrinsic derivatives are given by

$$\begin{aligned} \phi_{;1} &= l^\alpha \partial \phi / \partial x^\alpha = \partial \phi / \partial \rho, \\ \phi_{;2} &= n^\alpha \partial \phi / \partial x^\alpha = \partial \phi / \partial \sigma + H \partial \phi / \partial \rho, \end{aligned}$$

and

$$\phi_{;3} = m^\alpha \partial \phi / \partial x^\alpha = \partial \phi / \partial \zeta.$$

Substituting the coordinates into the commutation relations satisfied by these intrinsic derivatives, and using (2.2), gives all spin coefficients zero except

$$\gamma_{242} = \partial H / \partial \zeta.$$

The Newman-Penrose field equations then yield

$$\psi_0 = \psi_1 = \psi_2 = \psi_3 = 0, \quad \psi_4 = -\partial^2 H / \partial \zeta^2,$$

where ψ_0, \dots, ψ_4 are five independent complex tetrad components of the Weyl tensor $C_{\alpha\beta\gamma\delta}$, namely

$$\begin{aligned} \psi_0 &= -C_{\alpha\beta\gamma\delta} l^\alpha m^\beta l^\gamma m^\delta = -C_{1313}, \\ \psi_1 &= -C_{\alpha\beta\gamma\delta} l^\alpha n^\beta l^\gamma m^\delta = -C_{1213}, \\ \psi_2 &= C_{\alpha\beta\gamma\delta} l^\alpha m^\beta n^\gamma \bar{m}^\delta = C_{1324}, \\ \psi_3 &= C_{\alpha\beta\gamma\delta} l^\alpha n^\beta n^\gamma \bar{m}^\delta = C_{1224}, \end{aligned}$$

and

$$\psi_4 = -C_{\alpha\beta\gamma\delta} n^\alpha \bar{m}^\beta n^\gamma \bar{m}^\delta = -C_{2424}.$$

The metric, field equations, and tetrad remain invariant in form when subjected to either of the

transformations

$$\rho' = \rho + \frac{df}{d\sigma} \bar{\zeta} + \frac{d\bar{f}}{d\sigma} \zeta + \frac{df}{d\sigma} f + \frac{d\bar{f}}{d\sigma} \bar{f}, \quad \sigma' = \sigma,$$

$$\zeta' = \zeta + f(\sigma),$$

$$l^\alpha \rightarrow l'^\alpha, \quad n^\alpha \rightarrow n'^\alpha + \frac{df}{d\sigma} \frac{d\bar{f}}{d\sigma} l'^\alpha - \frac{df}{d\sigma} \bar{m}^\alpha - \frac{d\bar{f}}{d\sigma} m^\alpha,$$

$$m^\alpha \rightarrow m'^\alpha - \frac{df}{d\sigma} l'^\alpha, \quad (2.3)$$

with

$$H' = H + \frac{d^2 f}{d\sigma^2} \bar{\zeta}' + \frac{d^2 \bar{f}}{d\sigma^2} \zeta' + \frac{df}{d\sigma} \frac{d\bar{f}}{d\sigma}$$

or

$$\rho' = \rho + f(\sigma), \quad \sigma' = \sigma, \quad \zeta' = \zeta$$

with

$$H' = H + \frac{df}{d\sigma} \quad (2.4)$$

The Gauss-Codazzi-Ricci equations [(1.1)-(1.4)] were discussed in a previous paper.⁶ In particular, for empty space-times of type N (which include the plane-fronted gravitational waves), it was shown that embedding is only possible if

$$a_{11} = a_{13} = b_{11} = b_{13} = 0.$$

Equations (1.1)-(1.3) are solved for the plane-fronted gravitational waves with $a_{12} \neq 0$ (Sec. 3) and with $a_{12} = 0, a_{34} = b_{34} = 0$ (Sec. 4). The only other possibility, that is, $a_{12} = 0$ but not both a_{34} and b_{34} zero, subdivides into several cases which are not completely solved.

3. EMBEDDINGS OF THE PLANE-FRONTED GRAVITATIONAL WAVES WITH $a_{12} \neq 0$

The Gauss equations (1.1) give

$$a_{11} = a_{13} = a_{34} = b_{11} = b_{13} = b_{34} = 0,$$

$$e_1 = -e_2,$$

$$b_{23} = \epsilon a_{23},$$

$$b_{33} = \epsilon a_{33},$$

$$b_{12} = \epsilon a_{12},$$

$$b_{22} = \epsilon a_{22} + \epsilon e_1 \psi_4 / a_{44},$$

where $\epsilon = \pm 1$. The transformation of normals (1.9) can be chosen to make a_{12} equal to unity. The Codazzi equations (1.2) and the Ricci equations (1.4) can then be solved to give s_m and a_{mn} in terms of two real functions $A(\sigma)$, which is nonzero, and $B(\sigma)$. In fact,

$$s_1 = s_3 = 0, \quad s_2 = B(\sigma),$$

$$a_{33} = -A(\sigma) \frac{\partial^2 H}{\partial \bar{\zeta}^2},$$

$$a_{23} = -\frac{dA}{d\sigma} \frac{\partial H}{\partial \bar{\zeta}} - A \frac{\partial^2 H}{\partial \bar{\zeta} \partial \sigma} - e_2 \epsilon B A \frac{\partial H}{\partial \bar{\zeta}},$$

⁷ M. Trümper and E. Schucking (unpublished).

⁸ E. Newman and R. Penrose, *J. Math. Phys.* 3, 566 (1962).

and

$$a_{22} = \rho e_2 \epsilon B + H - A \frac{\partial H}{\partial \zeta} \cdot \frac{\partial H}{\partial \bar{\zeta}} - \epsilon e_1 B \left[-\frac{\partial(AH)}{\partial \sigma} + e_1 \epsilon BAH \right] + \frac{\partial[-\partial(AH)]}{\partial \sigma^2} + \frac{\partial[e_1 \epsilon BAH]}{\partial \sigma}.$$

In this calculation, which is somewhat lengthy, two functions of integrability have been set zero by the transformations (2.3) and (2.4). The field equations (2.2) have also been used and it can be shown that the Codazzi equations (1.3) are identically satisfied.

Equations (1.5)–(1.7), with a_{mn} , b_{mn} , and s_n replaced by the expressions found above, can now be solved for the normals η_1^Q , η_2^Q , and the embedding coordinates y^Q . In particular,

$$y^Q = -e_1 AHF^Q + G^Q \bar{\zeta} + \bar{G}^Q \zeta + ((e_1 A dF^Q/d\sigma) + \epsilon ABF^Q)\rho + \int J^Q d\sigma,$$

where G^Q are constants and F^Q, J^Q are real functions of σ satisfying the equations

$$e_1 F^Q = d \left[e_1 A \frac{dF^Q}{d\sigma} + \epsilon ABF^Q \right] / d\sigma \quad (3.1)$$

and

$$\frac{dF^Q}{d\sigma} + J^Q + \epsilon e_1 AB \frac{dJ^Q}{d\sigma} - A \frac{d^2 J^Q}{d\sigma^2} - \frac{dA}{d\sigma} \cdot \frac{dJ^Q}{d\sigma} = 0. \quad (3.2)$$

To obtain explicit embeddings of the plane-fronted gravitational waves, suppose that F_1 and F_2 satisfy Eq. (3.1), J_1 and J_2 satisfy the homogeneous equation corresponding to (3.2), and L_1 and L_2 satisfy Eq. (3.2) with F^Q replaced by F_1 and F_2 , respectively. Then

$$y^Q = Z_R^Q y_1^R + C^Q,$$

where Z_R^Q, C^Q are constants and

$$y_1^1 = -AHF_1 + \rho(e_1 \epsilon ABF_1 + A dF_1/d\sigma) + \int L_1 d\sigma,$$

$$y_1^2 = -AHF_2 + \rho(e_1 \epsilon ABF_2 + A dF_2/d\sigma) + \int L_2 d\sigma,$$

$$y_1^3 = \bar{\zeta}, \quad y_1^4 = \zeta, \quad y_1^5 = \int J_1 d\sigma, \quad y_1^6 = \int J_2 d\sigma.$$

Define the functions E_1, \dots, E_7 by

$$\begin{aligned} E_1 &= F_1 J_1 - F_2 J_2, \\ E_2 &= L_1 J_1 - L_2 J_2, \\ E_3 &= A \left[J_1 \frac{dF_1}{d\sigma} - J_2 \frac{dF_2}{d\sigma} \right] - 1, \\ E_4 &= A \left[F_1 \frac{dJ_1}{d\sigma} - F_2 \frac{dJ_2}{d\sigma} \right] + 1, \end{aligned}$$

$$E_5 = J_1 \frac{dL_1}{d\sigma} - J_2 \frac{dL_2}{d\sigma} + L_1 \frac{dJ_1}{d\sigma} - L_2 \frac{dJ_2}{d\sigma},$$

$$E_6 = A \left[\frac{dF_1}{d\sigma} \cdot \frac{dJ_1}{d\sigma} - \frac{dF_2}{d\sigma} \cdot \frac{dJ_2}{d\sigma} \right] - \epsilon e_1 B,$$

$$E_7 = 2A^2 \left[\frac{dL_1}{d\sigma} \cdot \frac{dJ_1}{d\sigma} - \frac{dL_2}{d\sigma} \cdot \frac{dJ_2}{d\sigma} \right] + 1.$$

Because of the definitions of F_1, \dots, L_2 , these seven functions satisfy

$$dE_1/d\sigma = A^{-1}[E_3 + E_4],$$

$$dE_2/d\sigma = E_5,$$

$$dE_3/d\sigma = E_6 - \epsilon e_1 BE_3 + \left[1 - \epsilon e_1 \frac{d(AB)}{d\sigma} \right] E_1,$$

$$dE_4/d\sigma = E_1 + E_6 + \epsilon e_1 BE_4,$$

$$dE_5/d\sigma = A^{-2}[E_3 + E_7] + A^{-1}E_2 + A^{-1} \left[\epsilon e_1 AB - \frac{dA}{d\sigma} \right] E_5,$$

$$dE_6/d\sigma = A^{-1}[E_3 + E_4] - A^{-1} \frac{dA}{d\sigma} E_6 - \epsilon e_1 \frac{dB}{d\sigma} E_4 - \epsilon e_1 A^{-1} B \frac{dA}{d\sigma} E_4,$$

$$dE_7/d\sigma = 2AE_5 + 2E_6 + 2\epsilon e_1 BE_7.$$

The values of F_1, \dots, L_2 and $dF_1/d\sigma, \dots, dL_2/d\sigma$ can be chosen arbitrarily at any point P . If they are chosen to make E_1, \dots, E_7 zero at P , then by repeated differentiation of the above, all derivatives of E_1, \dots, E_7 will be zero at P . With this choice of F_1, \dots, L_2 and $dF_1/d\sigma, \dots, dL_2/d\sigma$, the functions E_1, \dots, E_7 will vanish in a neighborhood of P and then the pseudo-Euclidean metric

$$ds^2 = 2 dy_1^1 dy_1^5 - 2 dy_1^2 dy_1^6 - 2 dy_1^3 dy_1^4 \quad (3.3)$$

is transformed into

$$\begin{aligned} ds^2 &= 2J_1 d\sigma \left[-d(AH)F_1 - AH \frac{dF_1}{d\sigma} d\sigma + L_1 d\sigma + \rho F_1 d\sigma + d\rho \left(e_1 \epsilon BAF_1 + A \frac{dF_1}{d\sigma} \right) \right] \\ &\quad - 2J_2 d\sigma \left[-d(AH)F_2 - AH \frac{dF_2}{d\sigma} d\sigma + L_2 d\sigma + \rho F_2 d\sigma + d\rho \left(e_1 \epsilon BAF_2 + A \frac{dF_2}{d\sigma} \right) \right] \\ &\quad - 2 d\zeta d\bar{\zeta} \\ &= -2 d\sigma d(AH)E_1 - 2H d\sigma^2(E_3 + 1) \\ &\quad + 2 d\sigma^2 E_2 + 2\rho d\sigma^2 E_1 \\ &\quad + 2 d\sigma d\rho[E_3 + 1 + e_1 \epsilon BAE_1] - 2 d\zeta d\bar{\zeta} \\ &= -2H d\sigma^2 + 2 d\sigma d\rho - 2 d\zeta d\bar{\zeta}. \end{aligned}$$

Notice that, although the field equations (2.2) were used in solving the Gauss–Codazzi–Ricci equations, they play no part in the above embedding. The embedding applies to the general plane-fronted waves (2.1). The embedding space (3.3) is of signature -2 and can be put in the form

$$ds^2 = (dz^1)^2 + (dz^2)^2 - (dz^3)^2 - (dz^4)^2 - (dz^5)^2 - (dz^6)^2. \quad (3.4)$$

As examples of members of the two-parameter family of embeddings, consider the cases $A = 1, B = 0$, and $A = 1, B = \epsilon e_1$. In the first case,

$$\begin{aligned} z^1 &= (-H \sinh \sigma + \rho \cosh \sigma + \frac{1}{2}\sigma \cosh \sigma - \frac{1}{2} \sinh \sigma + \sinh \sigma)/2^{\frac{1}{2}}, \\ z^2 &= (-H \cosh \sigma + \rho \sinh \sigma + \frac{1}{2}\sigma \sinh \sigma - \frac{1}{2} \cosh \sigma - \cosh \sigma)/2^{\frac{1}{2}}, \\ z^3 &= (\zeta + \bar{\zeta})/2^{\frac{1}{2}}, \\ z^4 &= -i(\zeta - \bar{\zeta})/2^{\frac{1}{2}}, \\ z^5 &= (-H \sinh \sigma + \rho \cosh \sigma + \frac{1}{2}\sigma \cosh \sigma - \frac{1}{2} \sinh \sigma - \sinh \sigma)/2^{\frac{1}{2}}, \\ z^6 &= (-H \cosh \sigma + \rho \sinh \sigma + \frac{1}{2}\sigma \sinh \sigma - \frac{1}{2} \cosh \sigma + \cosh \sigma)/2^{\frac{1}{2}}. \end{aligned}$$

In the second case,

$$\begin{aligned} z^1 &= \alpha_1 \exp \{-\frac{1}{2}\sigma(5^{\frac{1}{2}} - 1)\} + \beta_1 \exp \{\frac{1}{2}\sigma(5^{\frac{1}{2}} - 1)\}, \\ z^2 &= \alpha_2 \exp \{\frac{1}{2}\sigma(5^{\frac{1}{2}} + 1)\} + \beta_2 \exp \{-\frac{1}{2}\sigma(5^{\frac{1}{2}} + 1)\}, \\ z^3 &= (\zeta + \bar{\zeta})/2^{\frac{1}{2}}, \\ z^4 &= -i(\zeta - \bar{\zeta})/2^{\frac{1}{2}}, \\ z^5 &= \alpha_1 \exp \{-\frac{1}{2}\sigma(5^{\frac{1}{2}} - 1)\} - \beta_1 \exp \{\frac{1}{2}\sigma(5^{\frac{1}{2}} - 1)\}, \\ z^6 &= \alpha_2 \exp \{\frac{1}{2}\sigma(5^{\frac{1}{2}} + 1)\} - \beta_2 \exp \{-\frac{1}{2}\sigma(5^{\frac{1}{2}} + 1)\}, \end{aligned}$$

where

$$\alpha_1 = -[5 + 5^{\frac{1}{2}}]/(10 \cdot 2^{\frac{1}{2}}), \quad \alpha_2 = [5 - 5^{\frac{1}{2}}]/(10 \cdot 2^{\frac{1}{2}}),$$

and

$$\begin{aligned} \beta_1 &= [-2H + (\rho - \frac{1}{2}e_1)(1 + 5^{\frac{1}{2}})]/(2 \cdot 2^{\frac{1}{2}}), \\ \beta_2 &= [-2H + (\rho - \frac{1}{2}e_1)(1 - 5^{\frac{1}{2}})]/(2 \cdot 2^{\frac{1}{2}}). \end{aligned}$$

4. EMBEDDINGS OF THE PLANE-FRONTED GRAVITATIONAL WAVES WITH $a_{12} = 0$,

$$a_{34} = b_{34} = 0$$

The Gauss equations (1.1) give

$$\begin{aligned} a_{11} = a_{13} = a_{12} = a_{34} = b_{11} = b_{13} = b_{12} = b_{34} &= 0, \\ e_1 &= -e_2, \quad b_{23} = \epsilon a_{23}, \\ b_{33} = \epsilon a_{33} = \epsilon e_1 \partial^2 H / \partial \bar{\zeta}^2, \end{aligned}$$

where $\epsilon = \pm 1$ and a transformation of normals (1.9) has been used to make

$$b_{22} = \epsilon a_{22} - \epsilon.$$

The Codazzi–Ricci equations can be solved to give a_{mn} and s_m in terms of a single real function $B(\sigma)$.

The embedding coordinates are then found to be

$$y^Q = HF^Q + G^Q \bar{\zeta} + G^Q \zeta - \left(e_1 \epsilon BF^Q + \frac{dF^Q}{d\sigma} \right) \rho + \int J^Q d\sigma,$$

where G^Q are constants, F^Q are real functions of σ satisfying the equation

$$d[e_1 \epsilon BF^Q + dF^Q/d\sigma]/d\sigma = 0, \quad (4.1)$$

and J^Q are real functions of σ satisfying

$$\frac{dF^Q}{d\sigma} - \epsilon B \frac{dJ^Q}{d\sigma} + e_1 \frac{d^2 J^Q}{J^2 \sigma^2} = 0. \quad (4.2)$$

To exhibit the embedding explicitly, suppose that F_1 and F_2 satisfy Eq. (4.1), J_1 and J_2 satisfy the homogeneous equation corresponding to (4.2), and L_1 and L_2 satisfy Eq. (4.2) with F^Q replaced by F_1 and F_2 , respectively. Then,

$$y^Q = Z_R^Q y_1^R + C^Q,$$

where Z_R^Q and C^Q are constants and

$$y_1^1 = e_1 HF_1 - \rho e_1 \left(e_1 \epsilon BF_1 + \frac{dF_1}{d\sigma} \right) + \int L_1 d\sigma,$$

$$y_1^2 = e_1 HF_2 - \rho e_1 \left(e_1 \epsilon BF_2 + \frac{dF_2}{d\sigma} \right) + \int L_2 d\sigma,$$

$$y_1^3 = \bar{\zeta}, \quad y_1^4 = \zeta, \quad y_1^5 = \int J_1 d\sigma, \quad y_1^6 = \int J_2 d\sigma.$$

The seven functions

$$E_1 = F_1 J_1 - F_2 J_2,$$

$$E_2 = L_1 J_1 - L_2 J_2,$$

$$E_3 = -e_1 \left[J_1 \frac{dF_1}{d\sigma} - J_2 \frac{dF_2}{d\sigma} \right] - 1,$$

$$E_4 = -e_1 \left[F_1 \frac{dJ_1}{d\sigma} - F_2 \frac{dJ_2}{d\sigma} \right] + 1,$$

$$E_5 = J_1 \frac{dL_1}{d\sigma} - J_2 \frac{dL_2}{d\sigma} + L_1 \frac{dJ_1}{d\sigma} - L_2 \frac{dJ_2}{d\sigma},$$

$$E_6 = -e_1 \left[\frac{dF_1}{d\sigma} \cdot \frac{dJ_1}{d\sigma} - \frac{dF_2}{d\sigma} \cdot \frac{dJ_2}{d\sigma} \right] - e_1 \epsilon B,$$

$$E_7 = 2 \left[\frac{dL_1}{d\sigma} \cdot \frac{dJ_1}{d\sigma} - \frac{dL_2}{d\sigma} \cdot \frac{dJ_2}{d\sigma} \right] + 1$$

satisfy

$$dE_1/d\sigma = -e_1 [E_3 + E_4],$$

$$dE_2/d\sigma = E_5,$$

$$dE_3/d\sigma = E_6 - e_1 \epsilon BE_3 + \epsilon (dB/d\sigma) E_1,$$

$$dE_4/d\sigma = E_6 + e_1 \epsilon BE_4,$$

$$dE_5/d\sigma = E_7 + E_3 + e_1 \epsilon BE_5,$$

$$dE_6/d\sigma = -e_1 \epsilon (dB/d\sigma) E_4,$$

$$dE_7/d\sigma = 2E_6 + 2e_1 \epsilon BE_7.$$

As in the last section, F_1, \dots, L_2 and $dF_1/d\sigma, \dots, dL_2/d\sigma$ can be chosen at any point P so that E_1, \dots, E_7 vanish in a neighborhood of P . The pseudo-Euclidean metric (3, 3) is then transformed into

$$\begin{aligned} ds^2 &= 2J_1 d\sigma \left[e_1 F_1 dH + e_1 H \frac{dF_1}{d\sigma} d\sigma \right. \\ &\quad \left. - d\rho e_1 \left(\frac{dF_1}{d\sigma} + e_1 \epsilon B F_1 \right) + L_1 d\sigma \right] \\ &\quad - 2J_2 d\sigma \left[e_1 F_2 dH + e_1 H \frac{dF_2}{d\sigma} d\sigma \right. \\ &\quad \left. - d\rho e_1 \left(\frac{dF_2}{d\sigma} + e_1 \epsilon B F_2 \right) + L_2 d\sigma \right] - 2 d\zeta d\bar{\zeta} \\ &= 2e_1 d\sigma dHE_1 - 2H d\sigma^2 (E_3 + 1) \\ &\quad + 2 d\rho d\sigma (E_3 + 1) - 2\epsilon B d\rho d\sigma E_1 \\ &\quad + 2 d\sigma^2 E_2 - 2 d\zeta d\bar{\zeta} \\ &= -2H d\sigma^2 + 2 d\rho d\sigma - 2 d\zeta d\bar{\zeta}. \end{aligned}$$

As an example of the embedding, take $B = 0$. Then the metric (3.4) is transformed into the plane-fronted wave under the transformation

$$\begin{aligned} z^1 &= (He_1\sigma - e_1\rho - \frac{1}{8}\sigma^3 - \sigma)/2^{\frac{1}{2}}, \\ z^2 &= (He_1 - \frac{3}{4}\sigma^2)/2^{\frac{1}{2}}, \\ z^3 &= (\zeta + \bar{\zeta})/2^{\frac{1}{2}}, \\ z^4 &= -i(\zeta - \bar{\zeta})/2^{\frac{1}{2}}, \\ z^5 &= (-He_1\sigma + e_1\rho + \frac{1}{8}\sigma^3 - \sigma)/2^{\frac{1}{2}}, \\ z^6 &= (-He_1 - \frac{1}{4}\sigma^2)/2^{\frac{1}{2}}. \end{aligned}$$

Again the family of embeddings found in this section does not depend upon the field equations (2.2). These embeddings, and those of the last section, remain embeddings of the plane-fronted waves when the terms $\int L_1 d\sigma$ and $\int L_2 d\sigma$ are omitted from the embedding coordinates. Using this fact, the pseudo-Euclidean metric (3.16) is transformed into the plane-fronted waves under the particularly simple transformations

$$\begin{aligned} y_1^1 &= H\sigma - \rho, \quad y_1^2 = H, \quad y_1^3 = \bar{\zeta}, \quad y_1^4 = \zeta, \\ y_1^5 &= -\sigma, \quad y_1^6 = -\frac{1}{2}\sigma^2. \end{aligned}$$

In general, this will be a *global* embedding of the plane-fronted waves in six-dimensional pseudo-Euclidean space of signature -2 . Penrose⁹ has proved that the plane waves (that is plane-fronted waves with $\partial^2 H/\partial \zeta^2$ a function of σ alone) cannot be

embedded *globally* in a pseudo-Euclidean space of dimension n and signature $2 - n$.

5. EMBEDDINGS OF OTHER SPACE-TIMES

The space-times considered in this section are those which have already been embedded by Rosen.¹ They are labeled, following Rosen, as $A1, B1, \dots, B5, C1, \dots, C8, D1, \dots, D4, E1, F1, \dots, F3, G1, H1, \dots, H3, J1, J1, \dots, J10$. The labeling A, \dots, J is characteristic of the embeddings found by Rosen but is of no consequence to the present discussion. The aim of this section is to give embeddings in pseudo-Euclidean spaces of lower dimensions than those of Rosen. It is obvious that an embedding of any space-time in four dimensions is minimal and an embedding of a nonflat space-time in five dimensions is minimal. It is also well known¹⁰ that an embedding of an empty space-time in six dimensions is minimal. Therefore the embeddings, given by Rosen, of the metrics $A1, B1, \dots, B5, C1, C2, C5, \dots, C8, D1, \dots, D4, J1, \dots, J7$ are trivially minimal and these metrics will not be considered here. The results of this, and the next section are summarized in the following table.

TABLE I. Summary.

Space-time	Dimension of embedding exhibited by Rosen	Dimension of embedding exhibited here	Is dimension of embedding known to be minimal?
C3	6	5	Yes
C4	6	—	Yes
E1	7	6	Yes
F1	7	6	Yes
F2	7	6	Yes
F3	7	6	Yes
G1	8	7	No
H1	10	7	†
H2	10	7	†
H3	10	6	Yes
J1	10	7	No
J8	8	6	Yes
J9	9 and 10	8	†
J10	9 and 10	7	Yes
J11	10	6	Yes

† These space-times contain arbitrary functions and so may include space-times of different embedding classes.

The new embeddings are now listed. In general the embedding is characterized by giving the space-time, the embedding space, and the appropriate transformation of coordinates.

⁹ R. Penrose, Rev. Mod. Phys. 37, 215 (1965).

¹⁰ Reference 4, p. 200.

C3. Interior Schwarzschild solution¹¹:

$$\frac{1}{4} \left[3 \left(1 - \frac{r_1^2}{R^2} \right)^{\frac{1}{2}} - \left(1 - \frac{r^2}{R^2} \right)^{\frac{1}{2}} \right]^2 dt^2 - \left(1 - \frac{r^2}{R^2} \right)^{-1} dr^2 - r^2 (d\theta^2 + \sin^2 \theta d\phi^2);$$

$r_1, R = \text{const.}$

$$ds^2 = (dz^1)^2 - (dz^2)^2 - (dz^3)^2 - (dz^4)^2 - (dz^5)^2;$$

$$z^1 = R \left[3 \left(1 - \frac{r_1^2}{R^2} \right)^{\frac{1}{2}} - \left(1 - \frac{r^2}{R^2} \right)^{\frac{1}{2}} \right] \sinh (t/2R),$$

$$z^2 = R \left[3 \left(1 - \frac{r_1^2}{R^2} \right)^{\frac{1}{2}} - \left(1 - \frac{r^2}{R^2} \right)^{\frac{1}{2}} \right] \cosh (t/2R),$$

$$z^3 = r \sin \theta \cos \phi,$$

$$z^4 = r \sin \theta \sin \phi,$$

$$z^5 = r \cos \theta.$$

It is interesting to note that for this space-time, the Weyl tensor is identically zero and the Gauss equations are satisfied by virtue of the symmetry of the field.

E1. Petrov space T_1 , group G_4 , metric 4¹²:

$$- (kx^4 + 1)^{\frac{1}{2}} [(dx^1)^2 + (dx^2)^2] - (kx^4 + 1)^{-\frac{1}{2}} (dx^3)^2 + (dx^4)^2; \quad k = \text{const.}$$

$$ds^2 = (dz^1)^2 - (dz^2)^2 - (dz^3)^2 - (dz^4)^2 - (dz^5)^2 + \epsilon (dz^6)^2, \quad \epsilon = \pm 1;$$

$$z^1 = \frac{9(kx^4 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^4 + 1)^{-\frac{1}{2}}}{16} + (kx^4 + 1)^{\frac{1}{2}} [(x^1)^2 + (x^2)^2 + 1]/2,$$

$$z^2 = \frac{9(kx^4 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^4 + 1)^{-\frac{1}{2}}}{16} + (kx^4 + 1)^{\frac{1}{2}} [(x^1)^2 + (x^2)^2 - 1]/2,$$

$$z^3 = x^1 (kx^4 + 1)^{\frac{1}{2}}$$

$$z^4 = x^2 (kx^4 + 1)^{\frac{1}{2}},$$

$$z^5 = \sinh x^3 (kx^4 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cosh x^3 (kx^4 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = +1,$$

or

$$z^5 = \sin x^3 (kx^4 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cos x^3 (kx^4 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = -1.$$

This embedding was found by solving the appropriate Gauss-Codazzi-Ricci equations.

¹¹ R. C. Tolman, *Relativity, Thermodynamics and Cosmology* (Oxford University Press, London, 1934).

¹² A. Z. Petrov, in *Recent Developments in General Relativity* (Pergamon Press, Ltd., London, 1962).

F1. Static cylindrically symmetric magnetic or electric geon¹³:

$$\left(1 + \frac{r^2}{a^2} \right)^2 (dt^2 - dr^2 - dz^2) - r^2 d\phi^2 \left(1 + \frac{r^2}{a^2} \right)^{-2},$$

$a = \text{const.}$

$$ds^2 = (dz^1)^2 - (dz^2)^2 + (dz^3)^2 - (dz^4)^2 - (dz^5)^2 + \epsilon (dz^6)^2; \quad \epsilon = \pm 1.$$

$$z^1 = t \left(1 + \frac{r^2}{a^2} \right),$$

$$z^2 = z \left(1 + \frac{r^2}{a^2} \right),$$

$$z^3 = \left(1 + \frac{r^2}{a^2} \right) \frac{(t^2 - z^2)}{2} + f(r) - \frac{r^2}{2a^2},$$

$$z^4 = \left(1 + \frac{r^2}{a^2} \right) \frac{(t^2 - z^2)}{2} + f(r) + \frac{r^2}{2a^2},$$

$$z^5 = \sinh \phi r \left(1 + \frac{r^2}{a^2} \right)^{-1},$$

$$z^6 = \cosh \phi r \left(1 + \frac{r^2}{a^2} \right)^{-1} \quad \text{if } \epsilon = +1,$$

or

$$z^5 = \sin \phi r \left(1 + \frac{r^2}{a^2} \right)^{-1},$$

$$z^6 = \cos \phi r \left(1 + \frac{r^2}{a^2} \right)^{-1} \quad \text{if } \epsilon = -1.$$

Here the function $f(r)$ satisfies

$$\frac{4r}{a^2} \cdot \frac{df}{dr} = \epsilon \left(1 - \frac{r^2}{a^2} \right)^2 \left(1 + \frac{r^2}{a^2} \right)^{-4} + \left(1 + \frac{r^2}{a^2} \right)^2.$$

F2. Petrov space T_1 , group G_4 , metric 5¹²:

$$- (kx^3 + 1)^{\frac{1}{2}} [(dx^1)^2 + (dx^2)^2] - (dx^3)^2 + (kx^3 + 1)^{-\frac{1}{2}} (dx^4)^2; \quad k = \text{const.}$$

$$ds^2 = (dz^1)^2 - (dz^2)^2 - (dz^3)^2 - (dz^4)^2 + (dz^5)^2 + \epsilon (dz^6)^2, \quad \epsilon = \pm 1;$$

$$z^1 = - \frac{9(kx^3 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^3 + 1)^{-\frac{1}{2}}}{16} + (kx^3 + 1)^{\frac{1}{2}} [(x^1)^2 + (x^2)^2 + 1]/2,$$

$$z^2 = - \frac{9(kx^3 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^3 + 1)^{-\frac{1}{2}}}{16} + (kx^3 + 1)^{\frac{1}{2}} [(x^1)^2 + (x^2)^2 - 1]/2,$$

$$z^3 = x^1 (kx^3 + 1)^{\frac{1}{2}},$$

$$z^4 = x^2 (kx^3 + 1)^{\frac{1}{2}},$$

$$z^5 = \sin x^3 (kx^3 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cos x^3 (kx^3 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = +1,$$

¹³ M. A. Melvin, *Phys. Letters* **8**, 65 (1964).

or

$$z^5 = \sinh x^4(kx^3 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cosh x^4(kx^3 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = -1.$$

F3. Petrov space T_1 , group G_4 , metric 6¹²:

$$(kx^3 + 1)^{\frac{1}{2}} [(dx^4)^2 - (dx^1)^2] - (dx^2)^2$$

$$- (kx^3 + 1)^{-\frac{1}{2}}(dx^2)^2; \quad k = \text{const.}$$

$$ds^2 = (dz^1)^2 - (dz^2)^2 - (dz^3)^2 + (dz^4)^2 - (dz^5)^2$$

$$+ \epsilon(dz^6)^2;$$

$$z^1 = -\frac{9(kx^3 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^3 + 1)^{-\frac{1}{2}}}{16}$$

$$+ (kx^3 + 1)^{\frac{3}{2}}[(x^4)^2 - (x^1)^2 + 1]/2,$$

$$z^2 = -\frac{9(kx^3 + 1)^{\frac{1}{2}}}{16k^2} + \frac{\epsilon(kx^3 + 1)^{-\frac{1}{2}}}{16}$$

$$+ (kx^3 + 1)^{\frac{3}{2}}[(x^4)^2 - (x^1)^2 - 1]/2,$$

$$z^3 = x^1(kx^3 + 1)^{\frac{1}{2}},$$

$$z^4 = x^4(kx^3 + 1)^{\frac{1}{2}},$$

$$z^5 = \sinh x^2(kx^3 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cosh x^2(kx^3 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = +1,$$

or

$$z^5 = \sin x^2(kx^3 + 1)^{-\frac{1}{2}},$$

$$z^6 = \cos x^2(kx^3 + 1)^{-\frac{1}{2}} \quad \text{if } \epsilon = -1.$$

G1. Degenerate static vacuum field, class C¹⁴.

This space-time is a special case of Weyl's static rotationally symmetric space-time considered next.

H1. Weyl's static rotationally symmetric solution^{13,15}:

$$\exp(2\psi) dt^2 - \exp(-2\psi)$$

$$\times [(dr^2 + dz^2) \exp(2\gamma) + r^2 d\phi^2];$$

$$\psi = \psi(r, z), \quad \gamma = \gamma(r, z).$$

$$ds^2 = (dz^1)^2 + (dz^2)^2 - (dz^3)^2 - (dz^4)^2 + d\sigma^2;$$

$$z^1 = \exp(\psi) \sin t,$$

$$z^2 = \exp(\psi) \cos t,$$

$$z^3 = r \exp(-\psi) \sin \phi,$$

$$z^4 = r \exp(-\psi) \cos \phi.$$

$d\sigma^2$ is a pseudo-Euclidean space in which is embedded the two-dimensional Riemannian space with metric

$$ds^2 = -\exp(2\gamma - 2\psi)[dr^2 + dz^2]$$

$$+ \left[dr \frac{\partial(r \exp - \psi)}{\partial r} + dz \frac{\partial(r \exp - \psi)}{\partial z} \right]^2$$

$$- \left[dr \frac{\partial(\exp \psi)}{\partial r} + dz \frac{\partial(\exp \psi)}{\partial z} \right]^2.$$

Since¹⁶ an n -dimensional Riemannian space can always be embedded locally in a pseudo-Euclidean space of dimension $\frac{1}{2}n(n + 1)$, it is certainly possible to find a space $d\sigma^2$ of dimension 3. The resulting embedding space of the space-times considered here is therefore of dimension 7.

H2. Cylindrical gravitational waveline element¹⁵:

$$-\exp(2\psi) dz^2 + \exp(-2\psi)$$

$$\times [(dt^2 - d\rho^2) \exp(2\gamma) - \rho^2 d\phi^2];$$

$$\psi = \psi(\rho, t), \quad \gamma = \gamma(\rho, t).$$

$$ds^2 = -(dz^1)^2 + (dz^2)^2 - (dz^3)^2 - (dz^4)^2 + d\sigma^2;$$

$$z^1 = \exp(\psi) \sinh z,$$

$$z^2 = \exp(\psi) \cosh z,$$

$$z^3 = \rho \exp(-\psi) \sin \phi,$$

$$z^4 = \rho \exp(-\psi) \cos \phi.$$

$d\sigma^2$ is a pseudo-Euclidean space in which is embedded the two-dimensional Riemannian space with metric

$$-\exp(2\gamma - 2\psi)[d\rho^2 - dt^2]$$

$$+ \left[d\rho \frac{\partial(\rho \exp - \psi)}{\partial \rho} + dt \frac{\partial(\rho \exp - \psi)}{\partial t} \right]^2$$

$$- \left[d\rho \frac{\partial(\exp \psi)}{\partial \rho} + dt \frac{\partial(\exp \psi)}{\partial t} \right]^2.$$

H3. "Anti-Mach" model:

$$-(dx^1)^2 + 4x^4 dx^1 dx^3 - 2 dx^2 dx^3$$

$$- 2(x^4)^2(dx^3)^2 - (dx^4)^2.$$

It is well known that this space-time is a special case of the plane gravitational waves and therefore, as a result of Secs. 3 and 4, can be embedded in a six-dimensional Euclidean space.

I1. Gödel model¹⁷:

$$[dx^0 + \exp(x^1/a) dx^2]^2 - (dx^1)^2$$

$$- \frac{1}{2} \exp(2x^1/a)(dx^2)^2 - (dx^3)^2; \quad a = \text{const.}$$

The three-dimensional Riemannian space

$$[dx^0 + \exp(x^1/a) dx^2]^2 - (dx^1)^2 - \frac{1}{2} \exp(2x^1/a)(dx^2)^2$$

can always be embedded in a pseudo-Euclidean space of dimension $\frac{1}{2} 3(3 + 1) = 6$. Hence the Gödel model

¹⁴ J. Ehlers and W. Kundt, in *Recent Developments in General Relativity* (Pergamon Press, Ltd., London, 1962).

¹⁵ N. Rosen, *Bull. Res. Council Israel Sect. A*: 3, 328 (1954).

¹⁶ A. Friedman, *J. Math. Mech.* 10, 625 (1961).

¹⁷ K. Gödel, *Rev. Mod. Phys.* 21, 447 (1949).

can be embedded in 7 dimensions. Mayer and others¹⁸⁻²⁰ have shown that in certain circumstances the Codazzi equations are consequences of the Gauss equation. For the Godel model a symmetric tensor a_{mn} can be found satisfying the Gauss equation (5.1) but not the Codazzi equation (5.2).

J8. Plane-fronted gravitational waves:

$$-dx^2 - dy^2 - 2 du dv - 2H(x, y, u) du^2.$$

$$ds^2 = -(dz^1)^2 - (dz^2)^2 - (dz^3)^2 - (dz^4)^2 + (dz^5)^2 + (dz^6)^2;$$

$$z^1 = x, \quad z^2 = y, \quad z^3 = (Hu + v + u)/2^{\frac{1}{2}},$$

$$z^4 = (H - \frac{1}{2}u^2)/2^{\frac{1}{2}}, \quad z^5 = (Hu + v - u)/2^{\frac{1}{2}},$$

$$z^6 = (H + \frac{1}{2}u^2)/2^{\frac{1}{2}}.$$

This is but one of the many embeddings of the plane-fronted waves obtained in the previous sections.

J9. Robinson-Trautman metric²¹:

$$C(u, v, y, z)dv^2 + du dv - \frac{u^2(dy^2 + dz^2)}{Q^2(v, y, z)}.$$

$$ds^2 = -(dz^1)^2 - (dz^2)^2 + (dz^3)^2 + (dz^4)^2 - (dz^5)^2 - (dz^6)^2 - (dz^7)^2 + (dz^8)^2;$$

$$z^1 = uQ^{-1} \cos z, \quad z^2 = uQ^{-1} \sin z,$$

$$z^7 = uQ^{-1} \sinh y, \quad z^8 = uQ^{-1} \cosh y,$$

$$z^3 = (\frac{1}{2}Cv + u + v)/2^{\frac{1}{2}}, \quad z^4 = \frac{1}{2}(C - v^2)/2^{\frac{1}{2}},$$

$$z^5 = (\frac{1}{2}Cv + u + v)/2^{\frac{1}{2}}, \quad z^6 = \frac{1}{2}(C + v^2)/2^{\frac{1}{2}}.$$

J10. Petrov space T_3 , group G_2^{12} :

$$-\exp(x^2)[(dx^1)^2 \exp(-2x^4) + (dx^2)^2] - 2dx^3 dx^4 + x^2[x^3 + \exp(x^2)](dx^4)^2.$$

$$ds^2 = -(dz^1)^2 - (dz^2)^2 - (dz^3)^2 - (dz^4)^2 - (dz^5)^2 + (dz^6)^2 + (dz^7)^2;$$

$$z^1 = \exp(\frac{1}{2}x^2) \exp(-x^4) \cos x^1,$$

$$z^2 = \exp(\frac{1}{2}x^2) \exp(-x^4) \sin x^1,$$

$$z^3 = [4 - \exp(-2x^4)]^{\frac{1}{2}} \exp(\frac{1}{2}x^2),$$

$$z^4 = x^3 + \frac{1}{2}x^4[1 - X],$$

$$z^5 = \frac{1}{2}[X - \frac{1}{2}(x^4)^2],$$

$$z^6 = -x^3 + \frac{1}{2}x^4[1 + X],$$

$$z^7 = \frac{1}{2}[X + \frac{1}{2}(x^4)^2],$$

where

$$X = x^2[x^3 + \exp(x^2)] + 4[4 - \exp(-2x^4)]^{-1} \times \exp(x^2) \exp(-2x^4).$$

J11. Petrov space T_2 , Group G_5^{12} :

This space-time is type N with a five-parameter group of motions. It is therefore a plane gravitational wave and can consequently be embedded in a six-dimensional pseudo-Euclidean space.

6. DISCUSSION OF THE EMBEDDINGS

The important question to be asked about a local embedding is whether or not the embedding space is of minimal dimension. For the reasons given in the last section, the embeddings, exhibited here, of the following metrics are trivially minimal:

$$C3, E1, F2, F3, H3, J8, J11.$$

The space-times $C4$ and $F1$ are both Einstein-Maxwell fields with type D Weyl tensor. Therefore, using the theorem quoted in the Introduction, the embedding of $C4$ exhibited by Rosen and the embedding of $F1$ exhibited here are both minimal.

The space-times $H1$, $H2$, and $J9$ are families of space-times and contain members of different embedding classes. It is therefore not appropriate to discuss the minimality of the embeddings of these metrics.

It can be shown that the Gauss-Codazzi-Ricci equations for embedding class two cannot be satisfied for the space-time $J10$. The embedding of this space-time exhibited here is therefore minimal. The calculation is quite straightforward and is not given. The result is of some importance. The space-time $J10$ is of type III and possesses hypersurface orthogonal geodesic rays with zero shear. The fact that this space-time is not of embedding class two shows that the necessary conditions for a space-time to be of embedding class two, found by the author,⁶ are *not* sufficient conditions.

Only the embeddings of $G1$ and $I1$ remain to be discussed. In principle it should be possible to find whether or not Eqs. (1.1)-(1.4) admit solutions for these space-times. Unfortunately, the calculations are too difficult to carry through.

ACKNOWLEDGMENT

The author would like to thank Dr. F. A. E. Pirani for his helpful comments on the original manuscript.

¹⁸ W. Mayer, Trans. Am. Math. Soc. **38**, 267 (1935).

¹⁹ C. B. Allendoerfer, Am. J. Math. **61**, 633 (1939).

²⁰ A. Schwartz, J. Math. & Phys. **20**, 30 (1941).

²¹ I. Robinson and A. Trautman, Proc. Roy. Soc. (London) **A265**, 463 (1962).

Structure of the Dirac Bracket in Classical Mechanics*

N. MUKUNDA† AND E. C. G. SUDARSHAN
Department of Physics, Syracuse University, Syracuse, New York

(Received 27 April 1967)

We discuss the structure of the Dirac bracket in classical mechanics. We consider a generalization of the usual Poisson bracket and show the close connection of this generalization to the Lagrange brackets of classical mechanics. We show how the Dirac bracket appears as a particular case of the generalized Poisson bracket, thus giving a simple reason why the Jacobi identity holds for the Dirac bracket. We also discuss the nature of the transformations generated via the Dirac bracket and the relation of these to canonical transformations.

INTRODUCTION

SEVERAL years ago, Dirac developed a canonical formalism for the Hamiltonian formulation of classical mechanical systems which are subject to constraints.¹ The usual Hamiltonian formulation of classical mechanics rests on the equivalence of the Lagrangian and the Hamiltonian equations of motion; and the passage from the Lagrangian variables of generalized position and velocity, q and \dot{q} , to the Hamiltonian variables of generalized position and momentum, q and p , is possible when and only when the velocities can be expressed in terms of the positions and the momenta. This requirement can be expressed in two equivalent ways: either (i) the Lagrangian equations of motion should specify *all* the accelerations as functions of positions and velocities; or (ii) the definitions of the momenta should not lead to any identities among the positions and momenta alone. The Dirac theory of constraints was intended to handle precisely those systems that do not fulfill this requirement, namely systems whose position and momentum variables obey certain identities and are therefore not independent. These identities are the constraints referred to earlier. In such cases the lack of complete specification of the accelerations by the Lagrangian equations of motion manifests itself also in ambiguities in the passage to an equivalent Hamiltonian formulation.

As a prelude to the quantization of such systems, Dirac proposed that the usual Poisson² brackets of classical mechanics be replaced by a new algebraic structure, now known as the Dirac bracket,¹ and that these new brackets be made to correspond to com-

mutators in quantum theory. If $f(q, p)$ and $g(q, p)$ are two functions defined on a $2N$ -dimensional phase space with coordinate variables $q_1 \cdots q_N, p_1 \cdots p_N$, then the Dirac bracket of f with g , $\{f, g\}^*$ is defined by

$$\{f, g\}^* = \{f, g\} - \{f, \theta^a\} C_{ab} \{ \theta^b, g \}.$$

Here, the curly brackets without stars are ordinary Poisson brackets. The functions $\theta^a(q, p)$ are a certain subset of all those functions whose vanishing expresses the constraints. They have the important property that, if we form a matrix whose elements are the Poisson brackets of the θ^a with one another, then this matrix is nonsingular. (It follows that we have an even number of θ^a 's.) The functions $C_{ab}(q, p)$ form the matrix inverse to the matrix of Poisson brackets:

$$C_{ab} \{ \theta^b, \theta^c \} = \delta_a^c.$$

A summation over repeated indices is assumed in the equations above.

For systems involving constraints, the Hamiltonian equations of motion can be expressed in terms of Dirac brackets, in the same way in which the equations of motion of systems without constraints are expressible in terms of Poisson brackets. Before the Dirac bracket can be introduced, however, the set of all constraints has to be separated into two classes, known as first class and second class constraints. The functions θ^a are the second class constraints, and this class is characterized precisely by the existence of the matrix C_{ab} . The Dirac brackets share many of the standard properties of Poisson brackets, namely linearity, antisymmetry, and the Jacobi identity.³ The main difference lies in the fact that with respect to them, the functions θ^a behave essentially like pure numbers. In other words, the Dirac bracket of θ^a with

* Work supported in part by the U.S. Atomic Energy Commission.

† Present address: Tata Institute of Fundamental Research, Bombay, India.

¹ P. A. M. Dirac, *Can. J. Math.* **2**, 129 (1950); "Lectures on Quantum Mechanics," Belfer Graduate School of Science Monograph Series No. 2 (Yeshiva University, New York, 1964).

² S. D. Poisson, *J. de l'Ecole Polytech.* **8**, 266 (1809).

³ C. G. J. Jacobi, *Compt. Rend.* **11**, 529 (1841); *Vorlesungen über Dynamik*, A. Clebsch, Ed. (Reiner, Berlin, 1866, 2nd ed., 1884).

any other function is identically zero. It is only the set of second class constraints that can be eliminated in this way by the use of the Dirac bracket.

In this paper, we would like to study and clarify in an algebraic way the structure and properties of the Dirac bracket, by relating it to the other two algebraic structures of classical mechanics, namely Poisson and Lagrange brackets. We will also study the relationship between the transformations generated by the Dirac brackets on the one hand, and those generated by Poisson brackets on the other. The latter are, of course, the canonical transformations of classical mechanics. The motivation for this study is the following. The original proof of the Jacobi identity for the Dirac bracket consisted of a straightforward but rather lengthy *verification* of the identity,⁴ without shedding much light on the structure of the bracket or suggesting any simple reason for suspecting that the identity might hold. Subsequently, it has been shown⁵ by Bergmann and Goldberg that one can start from a certain continuous group of coordinate transformations in phase space having special properties with respect to the constraints; one then finds that the infinitesimal Lie brackets corresponding to this group are, in fact, Dirac brackets. The associativity of the group multiplication law then automatically guarantees the Jacobi identity for the Dirac bracket. Our interest is in exhibiting in a direct and algebraic way the reason why the Dirac bracket looks the way it does and the reason why it obeys the Jacobi identity, and after that examine the group of coordinate transformations generated by it.

In Sec. 1, we briefly review the properties of Poisson and Lagrange brackets and of canonical transformations in phase space. This material is completely standard and is included only for the sake of completeness. Section 2 consists of a straightforward extension of Poisson brackets to what we will call a generalized Poisson bracket. These brackets can be related in a direct way to Lagrange brackets. In Sec. 3, we show how the Dirac brackets arise as a special case of the generalized Poisson brackets. Finally, Sec. 4 contains a discussion of the coordinate transformations generated by the Dirac brackets and of the relation of these transformations to canonical transformations. In this paper, we will not be interested in any particular Lagrangians or Hamiltonians, and we will not need to make statements which are valid only when the constraint functions vanish.

⁴ Compare the remarks by Dirac, "Lectures on Quantum Mechanics," Belfer Graduate School of Science Monograph Series No. 2 (Yeshiva University, New York, 1964), p. 42.

⁵ P. G. Bergmann and I. Goldberg, Phys. Rev. **98**, 531 (1955).

1. POISSON AND LAGRANGE BRACKET: CANONICAL TRANSFORMATIONS⁶

In the $2N$ -dimensional phase space of a classical mechanical system with canonical variables $q_1 \cdots q_N, p_1 \cdots p_N$ we define the Poisson bracket (PB) of any two functions $f(q, p)$ and $g(q, p)$ to be a third function given by

$$\{f, g\}(q, p) = \sum_{k=1}^N \left(\frac{\partial f}{\partial q_k} \cdot \frac{\partial g}{\partial p_k} - \frac{\partial f}{\partial p_k} \cdot \frac{\partial g}{\partial q_k} \right). \quad (1.1)$$

Introducing the variables

$$\omega^\mu = \sum_{k=1}^N (\delta_{\mu k} q_k + \delta_{\mu, k+N} p_k) \quad (1.2)$$

and the constant matrix

$$\epsilon^{\mu\nu} = \delta_{\mu, \nu+N} - \delta_{\mu+N, \nu}, \quad (1.3)$$

we could rewrite Eq. (1.1), defining the PB in tensor notation⁷

$$\{f, g\}(\omega) = \epsilon^{\mu\nu} \frac{\partial f(\omega)}{\partial \omega^\mu} \cdot \frac{\partial g(\omega)}{\partial \omega^\nu}. \quad (1.4)$$

In either form [(1.1) or (1.4)], the PB satisfies the Jacobi identity

$$\{\{h, f\}, g\} + \{\{f, g\}, h\} + \{\{g, h\}, f\} = 0 \quad (1.5)$$

for any three functions f, g, h . Using the form (1.4), the only property of $\epsilon^{\mu\nu}$ used is its antisymmetry.

Canonical transformations can be characterized in the following way. The PB's of the basic variables ω^μ with one another have the standard values

$$\{\omega^\mu, \omega^\nu\} = \epsilon^{\mu\nu}.$$

Using the definition (1.4), we see that the PB preserving property of canonical transformations can be transcribed as follows:

$$\frac{\partial \omega'^\mu}{\partial \omega^\sigma} \cdot \frac{\partial \omega'^\nu}{\partial \omega^\rho} \epsilon^{\sigma\rho} = \epsilon^{\mu\nu} \quad (1.6)$$

Thus canonical transformations are those transformations with respect to which $\epsilon^{\mu\nu}$ behaves as an invariant second rank antisymmetric tensor of contravariant type.

The covariant tensor $\epsilon_{\mu\nu}$ is defined as the inverse matrix to $\epsilon^{\mu\nu}$ and has the elements:

$$\epsilon_{\mu\nu} = -\delta_{\mu, \nu+N} + \delta_{\mu+N, \nu}. \quad (1.7)$$

Given any set $\phi^\alpha(\omega)$ of $2N$ -independent functions we could express the ω^μ as functions of ϕ^α . Then we

⁶ See, for example, any of the standard texts, such as: H. Goldstein, *Classical Mechanics* (Addison-Wesley Publishing Company, Inc., Reading, Mass., 1965); H. C. Corben and P. Stehle, *Classical Mechanics* (John Wiley & Sons, Inc., New York, 1960).

⁷ For an introduction to the tensor notation, see C. W. Kilmister, *Hamiltonian Dynamics* (John Wiley & Sons, Inc., New York, 1964).

define the Lagrange bracket (LB)⁸ of ϕ^α and ϕ_β according to

$$L_{\alpha\beta}(\phi) = \epsilon_{\mu\nu} \frac{\partial\omega^\mu}{\partial\phi^\alpha} \cdot \frac{\partial\omega^\nu}{\partial\phi^\beta}. \quad (1.8)$$

It is then well known that

$$L_{\alpha\sigma}(\phi)\{\phi^\sigma, \phi^\beta\} = \delta_\alpha^\beta. \quad (1.9)$$

More relevant is the identity

$$\frac{\partial}{\partial\phi^\alpha} L_{\beta\gamma}(\phi) + \frac{\partial}{\partial\phi^\beta} L_{\gamma\alpha}(\phi) + \frac{\partial}{\partial\phi^\gamma} L_{\alpha\beta}(\phi) = 0 \quad (1.10)$$

which is the Lagrangian analogue of the Jacobi identity (1.5). It is seen that the left-hand side of (1.10) is totally antisymmetric in α, β, γ and that the identity holds as a consequence of the antisymmetry of $\epsilon_{\mu\nu}$ in its indices.

2. GENERALIZED POISSON BRACKETS

In this section we consider a generalization of the PB as given in (1.4).⁹ Let there be given a set of functions $\eta^{\mu\nu}(\omega)$, antisymmetric in μ and ν , and obeying the identity (2.12) which will be derived later on. Define the generalized Poisson bracket (GPB) of any two functions $f(\omega), g(\omega)$ to be a third function $h(\omega)$ given by

$$\{f, g\}^*(\omega) \equiv h(\omega) = \eta^{\mu\nu}(\omega) \frac{\partial f(\omega)}{\partial\omega^\mu} \cdot \frac{\partial g(\omega)}{\partial\omega^\nu}. \quad (2.1)$$

We first consider the behavior of $\eta^{\mu\nu}(\omega)$ under coordinate transformations. Let $\omega^\mu \rightarrow \omega'^\mu$ be a general coordinate transformation, the ω'^μ being independent functions of the ω^μ . Given any function $f(\omega)$, we can define a new function f' by the equation

$$f'(\omega') = f(\omega). \quad (2.2a)$$

This is the transformation law characteristic of a "scalar field." By means of it we are led from a function f with a certain functional form to a function f' with a (generally) different functional form. The behavior of $\eta^{\mu\nu}(\omega)$ under a general transformation of coordinates is fixed by requiring that the GPB of two scalar fields be itself a scalar field. Thus, if in the variables ω^μ we have

$$\{f, g\}^*(\omega) = h(\omega) \quad (2.1)$$

and $f(\omega), g(\omega)$ go over into functions f' and g'

according to (2.2a), then we also define $h(\omega)$ to go over into $h'(\omega')$ according to (2.2b):

$$h'(\omega') = h(\omega). \quad (2.2b)$$

Equation (2.1) expresses h in terms of f and g . Similarly, we can express h' in terms of f' and g' :

$$\begin{aligned} h'(\omega') &= \eta^{\mu\nu}(\omega) \frac{\partial f(\omega)}{\partial\omega^\mu} \cdot \frac{\partial g(\omega)}{\partial\omega^\nu} \\ &= \eta^{\mu\nu}(\omega) \frac{\partial f'(\omega')}{\partial\omega'^\rho} \cdot \frac{\partial\omega'^\rho}{\partial\omega^\mu} \cdot \frac{\partial g'(\omega')}{\partial\omega'^\sigma} \cdot \frac{\partial\omega'^\sigma}{\partial\omega^\nu} \\ &= \eta'^{\rho\sigma}(\omega') \frac{\partial f'(\omega')}{\partial\omega'^\rho} \cdot \frac{\partial g'(\omega')}{\partial\omega'^\sigma}, \end{aligned} \quad (2.3)$$

where

$$\eta'^{\rho\sigma}(\omega') = \frac{\partial\omega'^\rho}{\partial\omega^\mu} \cdot \frac{\partial\omega'^\sigma}{\partial\omega^\nu} \eta^{\mu\nu}(\omega). \quad (2.4)$$

Thus $\eta^{\mu\nu}(\omega)$ transforms as a second-rank antisymmetric tensor of contravariant type.

We may state the content of (2.4) in the following form. *The GPB is an operation whereby, given two scalar fields f and g , a third one h is determined.* In each coordinate system, a given scalar field is represented by a specific function of those coordinates. The explicit expression of the function representing h in terms of those representing f and g , depends on the particular coordinate system. Equation (2.4) shows how this explicit expression changes when one goes from one set of coordinates to another.

Next we define the analogues of canonical transformations. For this purpose, we focus attention on the change in functional form, $f \rightarrow f'$, produced by a change of coordinates, $\omega \rightarrow \omega'$, when f' is defined in terms of f by (2.2a). Given two functions $f(\omega), g(\omega)$, we consider the functions $f'(\omega), g'(\omega)$ which are obtained by using the functional forms f', g' but with arguments ω instead of ω' . For an arbitrary change of coordinates, the GPB of $f(\omega)$ with $g(\omega)$

$$\{f, g\}^*(\omega) \equiv h(\omega) = \eta^{\mu\nu}(\omega) \frac{\partial f(\omega)}{\partial\omega^\mu} \cdot \frac{\partial g(\omega)}{\partial\omega^\nu} \quad (2.5)$$

and that of $f'(\omega)$ with $g'(\omega)$

$$\{f', g'\}^*(\omega) = k(\omega) = \eta^{\mu\nu}(\omega) \frac{\partial f'(\omega)}{\partial\omega^\mu} \cdot \frac{\partial g'(\omega)}{\partial\omega^\nu} \quad (2.6)$$

will not bear any special relationship to one another. However, if we demand that the transformation be such that

$$k(\omega) = h'(\omega), \quad (2.7)$$

⁸ J. L. Lagrange, Memoires de l'Institut de France, (1808); reprinted in Oeuvres, Vol. VI, p. 713.

⁹ See, for example, C. W. Kilmister, Ref. 7, Chap. 4.

where

$$h'(\omega') = h(\omega), \tag{2.8}$$

then we have the following consequence:

$$\{f', g'\}^*(\omega) = (\{f, g\}^*)(\omega). \tag{2.9}$$

This is the statement that the operation of taking the GPB commutes with the operation of changing the functional form of a function according to the prescription given above. The requirement (2.7) is equivalent to the following one:

$$\eta^{\rho\sigma}(\omega') = \frac{\partial\omega'^{\rho}}{\partial\omega^{\mu}} \frac{\partial\omega'^{\sigma}}{\partial\omega^{\nu}} \eta^{\mu\nu}(\omega). \tag{2.10}$$

Notice that the same functions $\eta^{\mu\nu}$, but with different arguments, appear on the two sides of (2.10). Transformations $\omega \rightarrow \omega'$, which obey (2.10), may be called canonical with respect to the GPB defined by $\eta^{\mu\nu}(\omega)$. In the particular case when $\eta^{\mu\nu}$ happens to be $\epsilon^{\mu\nu}$, (2.10) coincides with (1.9), and we obtain the usual canonical transformations of classical mechanics.

We conclude this section with a discussion of the Jacobi identity.¹⁰ We will demand that $\eta^{\mu\nu}(\omega)$ be such that for any three functions f, g, h , we have

$$\begin{aligned} \{\{f, g\}^*, h\}^* + \{\{g, h\}^*, f\}^* \\ + \{\{h, f\}^*, g\}^* = 0. \end{aligned} \tag{2.11}$$

If (2.1) is substituted in (2.11), two kinds of terms appear, those without derivatives of $\eta^{\mu\nu}$ and those with derivatives. The former vanish by themselves, due to the antisymmetry of $\eta^{\mu\nu}$. The vanishing of the latter leads to

$$\eta^{\lambda\mu}(\omega) \frac{\partial\eta^{\nu\rho}(\omega)}{\partial\omega^{\mu}} + \eta^{\nu\mu}(\omega) \frac{\partial\eta^{\rho\lambda}(\omega)}{\partial\omega^{\mu}} + \eta^{\rho\mu}(\omega) \frac{\partial\eta^{\lambda\nu}(\omega)}{\partial\omega^{\mu}} = 0. \tag{2.12}$$

Thus, the Jacobi identity for the GPB is equivalent to (2.12). We will assume that the GPB is nondegenerate in the sense that the functions $\eta^{\mu\nu}(\omega)$ form a non-singular matrix, and we denote the matrix elements of the inverse matrix by $\eta_{\mu\nu}(\omega)$:

$$\eta_{\mu\nu}(\omega)\eta^{\nu\lambda}(\omega) = \delta_{\mu}^{\lambda}. \tag{2.13}$$

Then (2.12) can be written much more simply in terms of $\eta_{\mu\nu}(\omega)$:

$$\frac{\partial\eta_{\mu\nu}(\omega)}{\partial\omega^{\lambda}} + \frac{\partial\eta_{\nu\lambda}(\omega)}{\partial\omega^{\mu}} + \frac{\partial\eta_{\lambda\mu}(\omega)}{\partial\omega^{\nu}} = 0. \tag{2.14}$$

Notice the resemblance between (2.14) and (1.10). The only difference is that in place of the ω^{μ} variables in (2.14), we have the ϕ^{α} variables in (1.10), while $\eta_{\mu\nu}(\omega)$ is replaced by the LB $L_{\alpha\beta}(\phi)$. This shows the close connection between the Jacobi identity for a GPB, on the one hand, and a standard property of LB's on the other. We will make use of this connection in the next section to derive the Dirac bracket.

3. THE DIRAC BRACKET

Let there be given a set of $2(N - \gamma)$ functions

$$\theta^a(\omega), \quad a = 1, 2, \dots, 2(N - \gamma), \tag{3.1}$$

which are independent of one another. Let us choose 2γ additional functions

$$\psi^m(\omega), \quad m = 1, 2, \dots, 2\gamma, \tag{3.2}$$

so that the θ^a and ψ^m together form $2N$ independent functions. We can form two matrices, one made up of the PB's of the θ 's and ψ 's with each other, the other made up of their LB's. According to (1.9), these matrices are inverse to one another. This may be expressed as follows:

$$\{\theta^a, \theta^b\}L_{bc}(\theta, \psi) + \{\theta^a, \psi^m\}L_{mc}(\theta, \psi) = \delta_c^a, \tag{3.3a}$$

$$\{\theta^a, \theta^b\}L_{bn}(\theta, \psi) + \{\theta^a, \psi^m\}L_{mn}(\theta, \psi) = 0, \tag{3.3b}$$

$$\{\psi^m, \theta^a\}L_{ab}(\theta, \psi) + \{\psi^m, \psi^n\}L_{nb}(\theta, \psi) = 0, \tag{3.3c}$$

$$\{\psi^m, \theta^a\}L_{an}(\theta, \psi) + \{\psi^m, \psi^p\}L_{pn}(\theta, \psi) = \delta_n^m. \tag{3.3d}$$

A similar set of equations can be written down, corresponding to taking the matrix of LB's first, and that of PB's next. Here, the L 's stand for the LB's:

$$L_{ab}(\theta, \psi) = \epsilon_{\mu\nu} \frac{\partial\omega^{\mu}}{\partial\theta^a} \cdot \frac{\partial\omega^{\nu}}{\partial\theta^b},$$

$$L_{am}(\theta, \psi) = -L_{ma}(\theta, \psi) = \epsilon_{\mu\nu} \frac{\partial\omega^{\mu}}{\partial\theta^a} \cdot \frac{\partial\omega^{\nu}}{\partial\psi^m}, \tag{3.4}$$

$$L_{mn}(\theta, \psi) = \epsilon_{\mu\nu} \frac{\partial\omega^{\mu}}{\partial\psi^m} \cdot \frac{\partial\omega^{\nu}}{\partial\psi^n}.$$

The LB's by themselves obey the identities (1.10). A subset of these involves differentiation with respect to ψ^m alone. These are

$$\frac{\partial L_{mn}(\theta, \psi)}{\partial\psi^l} + \frac{\partial L_{nl}(\theta, \psi)}{\partial\psi^m} + \frac{\partial L_{lm}(\theta, \psi)}{\partial\psi^n} = 0. \tag{3.5}$$

The considerations of the previous section show that if we set

$$\eta_{mn}(\theta, \psi) \equiv L_{mn}(\theta, \psi), \tag{3.6}$$

¹⁰ J. M. Souriau, Commun. Math. Phys. 1, 374 (1966).

if $\eta_{mn}(\theta, \psi)$ possesses an inverse $\eta^{mn}(\theta, \psi)$

$$\eta_{mn}(\theta, \psi)\eta^{nl}(\theta, \psi) = \delta_m^l, \quad (3.7)$$

and if we define a bracket by

$$\{f, g\}^*(\theta, \psi) \equiv \eta^{mn}(\theta, \psi) \frac{\partial f}{\partial \psi^m} \frac{\partial g}{\partial \psi^n}, \quad (3.8)$$

then this will obey the Jacobi identity and therefore be a GPB in 2γ variables. In this definition, we have used the fact that any function of ω^μ can be written as a function of θ^a, ψ^m . Partial differentiation with respect to a ψ^m is carried out keeping the θ^a and the other ψ 's constant. Clearly,

$$\{f, \theta^a\}^* = 0. \quad (3.9)$$

We now show that the bracket (3.8) is just the Dirac bracket (DB) of f with g . For this, we must first find the matrix η^{mn} inverse to η_{mn} . Let us assume that the submatrix of PB's of the θ^a with one another possesses an inverse:

$$C_{ab}(\theta, \psi) \cdot \{\theta^b, \theta^c\} = \delta_a^c. \quad (3.10)$$

From (3.3b), we then find

$$L_{bn}(\theta, \psi) = -C_{ba}(\theta, \psi)\{\theta^a, \psi^l\}\eta_{ln}(\theta, \psi). \quad (3.11)$$

Substituting this in (3.3d) gives

$$[\{\psi^m, \psi^l\} - \{\psi^m, \theta^b\}C_{ba}(\theta, \psi)\{\theta^a, \psi^l\}]\eta_{ln}(\theta, \psi) = \delta_n^m. \quad (3.12)$$

This shows that η^{mn} exists, and is given by

$$\eta^{mn}(\theta, \psi) = \{\psi^m, \psi^n\} - \{\psi^m, \theta^a\}C_{ab}(\theta, \psi)\{\theta^b, \psi^n\}. \quad (3.13)$$

Conversely, it may be easily shown that if η^{mn} exists, then so does C_{ab} .

We now use (3.13) in (3.8) to find

$$\begin{aligned} & \{f, g\}^* \\ &= \left[\frac{\partial f}{\partial \psi^m} \{\psi^m, \psi^n\} - \frac{\partial f}{\partial \psi^m} \{\psi^m, \theta^a\}C_{ab}(\theta, \psi)\{\theta^b, \psi^n\} \right] \cdot \frac{\partial g}{\partial \psi^n}. \end{aligned} \quad (3.14)$$

It is easy to see that by the addition of terms which in fact vanish, we can rewrite this in the form

$$\{f, g\}^* = \{f, \psi^n\} \frac{\partial g}{\partial \psi^n} - \{f, \theta^a\}C_{ab}(\theta, \psi)\{\theta^b, \psi^n\} \frac{\partial g}{\partial \psi^n}. \quad (3.15)$$

Once again, by the addition of vanishing terms, we

write (3.15) in the final form

$$\{f, g\}^* = \{f, g\} - \{f, \theta^a\}C_{ab}(\theta, \psi)\{\theta^b, g\}. \quad (3.16)$$

We see at once that this GPB in 2γ variables is just the DB written down in the Introduction. We also see that it is an expression determined solely by the functions θ^a , and does not depend on the choice of the functions ψ^m which were originally used in (3.8) to define it.

4. DIRAC BRACKET TRANSFORMATIONS

In this last section, we consider the transformations generated via the DB,¹¹ or more generally, the coordinate transformations which are canonical with respect to the DB.

We first take up the question of the nature of these transformations *per se*. Written in the form (3.8), we see that the DB is a nondegenerate GPB in 2γ variables. For the moment, we ignore the presence of the variables θ^a . Now it is a well-known fact in the mathematical literature that all such symplectic structures in a given number of variables are locally isomorphic.⁹ In other words, *by proper choice of the variables ψ^m , the coefficients in (3.8) can be made constants, so that the matrix $\|\eta^{mn}\|$ has exactly the same structure as the matrix $\|\epsilon^{\mu\nu}\|$ in (1.3)*. Nevertheless, it may be useful to give here a simple proof of this statement, at least for the sake of completeness.

For this purpose, consider (3.5):

$$\frac{\partial \eta_{mn}}{\partial \psi^l} + \frac{\partial \eta_{nl}}{\partial \psi^m} + \frac{\partial \eta_{lm}}{\partial \psi^n} = 0. \quad (3.5)$$

This equation is analogous to one set of the Maxwell equations of electrodynamics, and exactly as in that case, one can show that η_{mn} may be expressed as a "curl" of a "vector"¹²:

$$\eta_{mn}(\psi) = \frac{\partial A_m(\psi)}{\partial \psi^n} - \frac{\partial A_n(\psi)}{\partial \psi^m}. \quad (4.1)$$

We have seen earlier that η^{mn} transforms as a contravariant tensor of the second rank, under an arbitrary change of coordinates $\psi \rightarrow \psi'$. Then, η_{mn} transforms as a covariant tensor, while $A_m(\psi)$ is a covariant vector field. The differentials of the coordinates,

¹¹ Given a GPB $\{f, g\}^*$, the transformation generated by a function $\phi(\omega)$ via this GPB is the transformation

$$\begin{aligned} f(\omega) &\rightarrow f'(\omega) \equiv [(\exp \tilde{\phi})f](\omega) \\ &= f(\omega) + \{\phi, f\}^*(\omega) + (1/2!)\{\phi, \{\phi, f\}^*\}^*(\omega) + \dots \end{aligned}$$

Because of the Jacobi identity obeyed by the GPB, this transformation may be shown to be canonical with respect to this GPB.

¹² See, for instance, J. L. Synge, *Relativity: The Special Theory* (North-Holland Publishing Company, Amsterdam 1965), p. 344.

$d\psi^m$, form a contravariant vector, so that the expression

$$A_m(\psi) d\psi^m \tag{4.2}$$

is invariant. From the theory of the Pfaff problem,¹³ it is known that one can always find a coordinate system, with coordinates ξ^m , in which (4.2) assumes the form

$$\xi^1 d\xi^2 + \xi^3 d\xi^4 + \dots + \xi^{2r-1} d\xi^{2r}. \tag{4.3}$$

In this coordinate system, the components $A'_m(\xi)$ of the covariant vector field are linear in the coordinates, so that the $\eta'_{mn}(\xi)$ are constants. The same is then true of the η'^{mn} and without loss of generality we may assume that the matrix $\|\eta'^{mn}\|$ has the same structure as $\|\epsilon^{\mu\nu}\|$ in (1.3). Since the GPB is assumed to be nondegenerate, we conclude that $r = \gamma$.

Thus we can find functions ξ^1, \dots, ξ^{2r} of the ψ^m such that

$$\begin{aligned} \{f, g\}^*(\psi) &\equiv \eta'^{mn}(\psi) \frac{\partial f}{\partial \psi^m} \frac{\partial g}{\partial \psi^n} \\ &= \epsilon^{mn} \frac{\partial f}{\partial \xi^m} \frac{\partial g}{\partial \xi^n}. \end{aligned} \tag{4.4}$$

Written in this form, we immediately see that the coordinate transformations which are canonical with respect to the DB are none other than the usual canonical transformation *in the variables* ξ^m . (We are ignoring, of course, arbitrary transformations that may be performed among the variables θ^a by themselves.) The group of transformations generated via the DB (3.16) is isomorphic to the group of canonical transformations in $2r$ variables generated via the PB in $2r$ variables.¹¹ (Under these transformations, the variables θ^a do not change.)

We consider next the relation of these transformations to the usual canonical transformations in the basic variables ω^μ . We started with functions of the variables ω^μ , and given the set of functions $\theta^a(\omega)$, we defined the DB (3.16). We can then ask whether or not and under what circumstances transformations canonical with respect to the DB also belong to the group of usual canonical transformations on the ω^μ . We shall answer this question by first looking at a simple example.

Let us take the case where the $2(N - r)$ functions $\theta^a(\omega)$ are just a subset of the canonical q and p variables:

$$\theta^a(\omega) = q_{r+1}, \dots, q_N; p_{r+1}, \dots, p_N. \tag{4.5}$$

¹³ E. Goursat, *Lecons sur le Probleme de Pfaff* (Hermann & Cie., Paris, 1922), p. 14.

Then the DB assumes the form

$$\{f, g\}^*(q, p) = \sum_{k=1}^r \left(\frac{\partial f}{\partial q_k} \frac{\partial g}{\partial p_k} - \frac{\partial f}{\partial p_k} \frac{\partial g}{\partial q_k} \right). \tag{4.6}$$

In this case, the DB has the effect of "freezing" $(N - r)$ canonical pairs of degrees of freedom. If we now consider a transformation generated¹¹ via the DB by means of a function ϕ , depending only on $q_1, \dots, q_r, p_1, \dots, p_r$:

$$\begin{aligned} f(q, p) &\rightarrow f'(q, p) = f(q, p) + \{\phi, f\}^*(q, p) \\ &\quad + \frac{1}{2!} \{\phi, \{\phi, f\}^*\}^*(q, p) + \dots \\ &= (\exp \phi) f(q, p), \end{aligned} \tag{4.7}$$

such a transformation is also a canonical one in terms of the basic variables ω^μ , in which $q_1 \dots q_r, p_1 \dots p_r$ transform among themselves, while $q_{r+1} \dots q_N, p_{r+1} \dots p_N$ remain unchanged. [It should be noted that if the function ϕ used in (4.7) also depends on the "frozen" variables $q_{r+1} \dots q_N, p_{r+1} \dots p_N$, the resulting transformation is generally not canonical in the basic variables ω^μ .]

Returning to the general case of arbitrary functions θ^a , one checks easily that the DB is unaltered if in place of θ^a , one uses a set of independent functions of them, θ'^a . Thus in some cases, it may be possible to replace θ^a by θ'^a in such a way that the θ'^a are in fact a subset of $2(N - r)$ canonical q, p variables. Precisely when this can be done is shown by the following:

Theorem: The necessary and sufficient condition, that the DB determined by the $2(N - r)$ functions $\theta^a(\omega)$ corresponds to "freezing" $(N - r)$ pairs of canonical variables in the ordinary PB is that the functions θ^a form a function group of rank $2(N - r)$.¹⁴

To prove this theorem, we note that the necessity is obvious, since a set of $2(N - r)$ variables made up of $(N - r)$ q 's and the corresponding p 's does form a function group of rank $2(N - r)$. On the other hand, this condition is sufficient. For, given such a function group, one can replace the θ^a by functions of themselves, θ'^a , such that the PB's of the θ'^a with each other assume the constant values corresponding to $(N - r)$ pairs of canonical variables.¹⁴

We conclude by noting the distinguishing properties of the DB, which are suggested by (4.4). Given

¹⁴ A set of functions $\theta^a(\omega)$ forms a function group if the PB's of the θ^a with each other can be expressed as functions of the θ^a alone. The rank of a function group is the number of independent functions in the function group. For further properties and details, see: L. P. Eisenhart, *Continuous Groups of Transformation* (Dover Publications, Inc., New York, 1961).

the functions $\theta^a(\omega)$, one could adjoin an arbitrary set of new functions $\zeta^m(\omega)$, making sure only that θ^a and ζ^m together form $2N$ independent functions; and, one could then define a bracket of two functions $f(\omega)$, $g(\omega)$ by

$$\{f, g\}^\# = \epsilon^{mn} \frac{\partial f}{\partial \zeta^m} \cdot \frac{\partial g}{\partial \zeta^n}. \quad (4.8)$$

This bracket certainly obeys the Jacobi identity, and treats the θ^a like pure numbers:

$$\{f, \theta^a\}^\# = 0. \quad (4.9)$$

One can then ask what is special about the DB. There are two features which are special to the DB. Firstly, while (4.4) shows that the DB is a particular case of the bracket (4.8), in general, for arbitrarily chosen functions $\zeta^m(\omega)$, (4.8) will not be a structure determined by the θ^a alone. The variables ξ^m appearing in (4.4) are, on the other hand, determined completely by the θ^a (up to a canonical transformation of the ξ^m among themselves). Secondly, the DB bears a special relationship to the PB in the following sense. If a function $f(\omega)$ is such that

$$\{f, \theta^a\} = 0, \quad \text{all } a, \quad (4.10)$$

then for all functions $g(\omega)$, we have

$$\{f, g\}^* = \{f, g\}. \quad (4.11)$$

Such a function $f(\omega)$ therefore generates the same transformation via the DB as it does via the PB.¹¹

SUMMARY

We have shown that the Dirac bracket arises as a special case of the generalized Poisson bracket. We have traced the origin of the Jacobi identity for Dirac brackets to a standard property of Lagrange brackets. Finally, we have seen that when and only when the second class constraints θ^a form a function group does the Dirac bracket correspond to "freezing" a subset of canonical variables in the ordinary Poisson bracket.

ACKNOWLEDGMENTS

Part of this work was done while one of us (N. M.) was at the Palmer Physical Laboratory of Princeton University. He would like to thank Professor W. Bleakney and Professor M. L. Goldberger for the hospitality of that institution.

Unitary Representations of the Lorentz Groups: Reduction of the Supplementary Series under a Noncompact Subgroup*

N. MUKUNDA†

Physics Department, Syracuse University, Syracuse, N. Y.

(Received 16 May 1967)

Unitary representations of $O(2, 1)$ belonging to the exceptional class are reduced with respect to the noncompact subgroup $O(1, 1)$. We recover the result that the spectrum of the generator of this subgroup covers the real line twice. Unitary representations of $O(3, 1)$ belonging to the supplementary series are reduced with respect to the noncompact subgroup $O(2, 1)$. These representations of $O(3, 1)$ may be labeled by a parameter ρ in the range $0 < \rho < 1$. Representations corresponding to $0 < \rho \leq \frac{1}{2}$ yield upon reduction only those representations of $O(2, 1)$ that belong to the continuous nonexceptional class; each of these appears twice. A representation corresponding to $\frac{1}{2} < \rho < 1$, however, yields upon reduction a single representation of $O(2, 1)$ of the exceptional class (with parameter $\sigma = \rho - \frac{1}{2}$) and, in addition, all the representations of $O(2, 1)$ of the nonexceptional continuous class. The exceptional representation appears only once, while the nonexceptional ones appear twice each.

INTRODUCTION

The purpose of this paper is to examine some properties of the supplementary series of unitary irreducible representations (UIR's) of the homoge-

neous Lorentz groups $O(2, 1)$ and $O(3, 1)$.¹ In previous papers we have shown how the UIR's of $O(2, 1)$

* Work supported in part by the U.S. Atomic Energy Commission.

† Present address: Tata Institute of Fundamental Research, Bombay, India.

¹ The unitary irreducible representations of $O(2, 1)$ were first constructed by V. Bargmann, *Ann. Math.* **48**, 568 (1947). Descriptions of the unitary irreducible representations of $O(3, 1)$ are to be found in: M. A. Naimark, *Linear Representations of the Lorentz Group* (The Macmillan Company, New York, 1964); I. M. Gel'fand, R. A. Minlos, and Z. Ya. Shapiro, *Representations of the Rotation and Lorentz Groups and their Applications* (The Macmillan Company, New York, 1963).

the functions $\theta^a(\omega)$, one could adjoin an arbitrary set of new functions $\zeta^m(\omega)$, making sure only that θ^a and ζ^m together form $2N$ independent functions; and, one could then define a bracket of two functions $f(\omega)$, $g(\omega)$ by

$$\{f, g\}^\# = \epsilon^{mn} \frac{\partial f}{\partial \zeta^m} \cdot \frac{\partial g}{\partial \zeta^n}. \quad (4.8)$$

This bracket certainly obeys the Jacobi identity, and treats the θ^a like pure numbers:

$$\{f, \theta^a\}^\# = 0. \quad (4.9)$$

One can then ask what is special about the DB. There are two features which are special to the DB. Firstly, while (4.4) shows that the DB is a particular case of the bracket (4.8), in general, for arbitrarily chosen functions $\zeta^m(\omega)$, (4.8) will not be a structure determined by the θ^a alone. The variables ξ^m appearing in (4.4) are, on the other hand, determined completely by the θ^a (up to a canonical transformation of the ξ^m among themselves). Secondly, the DB bears a special relationship to the PB in the following sense. If a function $f(\omega)$ is such that

$$\{f, \theta^a\} = 0, \quad \text{all } a, \quad (4.10)$$

then for all functions $g(\omega)$, we have

$$\{f, g\}^* = \{f, g\}. \quad (4.11)$$

Such a function $f(\omega)$ therefore generates the same transformation via the DB as it does via the PB.¹¹

SUMMARY

We have shown that the Dirac bracket arises as a special case of the generalized Poisson bracket. We have traced the origin of the Jacobi identity for Dirac brackets to a standard property of Lagrange brackets. Finally, we have seen that when and only when the second class constraints θ^a form a function group does the Dirac bracket correspond to "freezing" a subset of canonical variables in the ordinary Poisson bracket.

ACKNOWLEDGMENTS

Part of this work was done while one of us (N. M.) was at the Palmer Physical Laboratory of Princeton University. He would like to thank Professor W. Bleakney and Professor M. L. Goldberger for the hospitality of that institution.

Unitary Representations of the Lorentz Groups: Reduction of the Supplementary Series under a Noncompact Subgroup*

N. MUKUNDA†

Physics Department, Syracuse University, Syracuse, N. Y.

(Received 16 May 1967)

Unitary representations of $O(2, 1)$ belonging to the exceptional class are reduced with respect to the noncompact subgroup $O(1, 1)$. We recover the result that the spectrum of the generator of this subgroup covers the real line twice. Unitary representations of $O(3, 1)$ belonging to the supplementary series are reduced with respect to the noncompact subgroup $O(2, 1)$. These representations of $O(3, 1)$ may be labeled by a parameter ρ in the range $0 < \rho < 1$. Representations corresponding to $0 < \rho \leq \frac{1}{2}$ yield upon reduction only those representations of $O(2, 1)$ that belong to the continuous nonexceptional class; each of these appears twice. A representation corresponding to $\frac{1}{2} < \rho < 1$, however, yields upon reduction a single representation of $O(2, 1)$ of the exceptional class (with parameter $\sigma = \rho - \frac{1}{2}$) and, in addition, all the representations of $O(2, 1)$ of the nonexceptional continuous class. The exceptional representation appears only once, while the nonexceptional ones appear twice each.

INTRODUCTION

The purpose of this paper is to examine some properties of the supplementary series of unitary irreducible representations (UIR's) of the homoge-

neous Lorentz groups $O(2, 1)$ and $O(3, 1)$.¹ In previous papers we have shown how the UIR's of $O(2, 1)$

* Work supported in part by the U.S. Atomic Energy Commission.

† Present address: Tata Institute of Fundamental Research, Bombay, India.

¹ The unitary irreducible representations of $O(2, 1)$ were first constructed by V. Bargmann, *Ann. Math.* **48**, 568 (1947). Descriptions of the unitary irreducible representations of $O(3, 1)$ are to be found in: M. A. Naimark, *Linear Representations of the Lorentz Group* (The Macmillan Company, New York, 1964); I. M. Gel'fand, R. A. Minlos, and Z. Ya. Shapiro, *Representations of the Rotation and Lorentz Groups and their Applications* (The Macmillan Company, New York, 1963).

belonging to the continuous nonexceptional and to the discrete classes can be written in such a way that the reduction of these UIR's with respect to the noncompact $O(1, 1)$ subgroup can be carried out; and in a similar way we have shown how the UIR's of $O(3, 1)$ belonging to the principal series can be reduced under the noncompact $O(2, 1)$ subgroup of $O(3, 1)$.² Here we take up the cases of the exceptional UIR's of $O(2, 1)$ and the supplementary series of UIR's of $O(3, 1)$ with a view to reducing them under $O(1, 1)$ and $O(2, 1)$, respectively. The methods we use are similar to those employed in the papers quoted above.

The material of this paper is arranged as follows. In Sec. 1, we give briefly the details concerning the UIR's of $O(2, 1)$ and $O(3, 1)$ which are of interest to us. This description is in a basis in which the UIR has been completely reduced under the relevant maximal compact subgroup. Sections 2 and 3 are devoted, respectively, to the reduction of these UIR's of $O(2, 1)$ and $O(3, 1)$ under their noncompact subgroups.

In comparison with the cases already treated, the UIR's considered here are somewhat more complicated to deal with. In particular, it turns out that the supplementary series of UIR's of $O(3, 1)$ have to be divided into two subclasses with markedly different properties in their reduction under $O(2, 1)$. These two classes will be dealt with in two subsections to Sec. 3.

1. EXCEPTIONAL UIR'S OF $O(2, 1)$ AND SUPPLEMENTARY UIR'S OF $O(3, 1)$

We summarize first the relevant UIR's of $O(2, 1)$.³ The group $O(2, 1)$ consists of all real linear unimodular transformations on three real variables χ_j , $j = 1, 2, 3$, which leave invariant the quadratic form

$$\chi_1^2 + \chi_2^2 - \chi_3^2,$$

and do not change the sign of χ_3 . Closely related is the group $SU(1, 1)$ of two-dimensional unimodular pseudounitary matrices; there is a two-to-one homomorphism from $SU(1, 1)$ to $O(2, 1)$. Elements g of $SU(1, 1)$ are in a one-to-one correspondence with matrices in the following way:

$$g \rightarrow \begin{pmatrix} \alpha & \beta \\ \bar{\beta} & \bar{\alpha} \end{pmatrix}, \quad |\alpha|^2 - |\beta|^2 = 1. \quad (1.1)$$

[α, β are complex numbers, and the bar denotes

² N. Mukunda, J. Math. Phys. 8, 2210 (1967); 9, 50 (1968). These two papers are referred to as (A) and (B), respectively. The reduction of the principal series of representations of $O(3, 1)$ under $O(2, 1)$ has also been carried out by S. Ström [Arkiv Fys. 34, 215 (1967)] and A. Sciarrino and M. Toller [J. Math. Phys. 8, 1252 (1967)].

³ More details may be found in (A), and in the references quoted therein.

complex conjugation.] The Lie algebra of $SU(1, 1)$ [or of $O(2, 1)$] is three dimensional, and the basic elements J_0, J_1, J_2 obey the commutation rules

$$\begin{aligned} [J_0, J_1] &= iJ_2, \\ [J_0, J_2] &= -iJ_1, \\ [J_1, J_2] &= -iJ_0. \end{aligned} \quad (1.2)$$

This algebra possesses a quadratic Casimir operator Q given by

$$Q = J_1^2 + J_2^2 - J_0^2. \quad (1.3)$$

The (single-valued) UIR's of $SU(1, 1)$ fall mainly into three categories: the continuous nonexceptional, the discrete, and the continuous exceptional. UIR's of the continuous exceptional type⁴ are labeled by the value of Q :

$$\begin{aligned} Q &= \frac{1}{4} - \sigma^2, \quad 0 < \sigma < \frac{1}{2}, \\ 0 < Q &< \frac{1}{4}. \end{aligned} \quad (1.4)$$

There is one such UIR for each value of the real parameter σ in the open interval given above. In these UIR's, the spectrum of the generator J_0 [which is the generator of spatial rotations in $O(2, 1)$] consists of all integers—positive, negative, and zero—each eigenvalue appearing exactly once. Thus these are single-valued UIR's of $O(2, 1)$. We may introduce a basis into the space of the UIR, corresponding to any given value of σ , made up of eigenvectors of J_0 :

$$\begin{aligned} J_0 |m\rangle &= m |m\rangle; \quad m = 0, \pm 2, \pm 1, \dots; \\ \langle m' | m \rangle &= \delta_{m'm}. \end{aligned} \quad (1.5)$$

In such a basis, the generators J_1 and J_2 have the following structure:

$$\begin{aligned} J_1 |m\rangle &= \frac{1}{2} [(m + \frac{1}{2})^2 - \sigma^2]^{\frac{1}{2}} |m + 1\rangle \\ &\quad + \frac{1}{2} [(m - \frac{1}{2})^2 - \sigma^2]^{\frac{1}{2}} |m - 1\rangle, \\ J_2 |m\rangle &= -\frac{i}{2} [(m + \frac{1}{2})^2 - \sigma^2]^{\frac{1}{2}} |m + 1\rangle \\ &\quad + \frac{i}{2} [(m - \frac{1}{2})^2 - \sigma^2]^{\frac{1}{2}} |m - 1\rangle. \end{aligned} \quad (1.6)$$

For our later use, it is necessary to mention here some remarkable facts relating to these UIR's of $O(2, 1)$. In any UIR of $O(2, 1)$, denoted by R , the unitary operators $U_{(g)}^{(R)}$, representing the elements g of $O(2, 1)$, may be specified by means of their matrix elements between eigenvectors of J_0 :

$$\mathfrak{U}_{mm'}^{(R)}(g) = \langle m | U_{(g)}^{(R)} | m' \rangle. \quad (1.7)$$

[The range of values of m and m' is appropriate to

⁴ V. Bargmann, Ref. 1, pp. 605 and 616.

the representation R .] Thus each such unitary operator is represented by an infinite-dimensional unitary matrix with discretely labeled rows and columns. Every matrix element of the form (1.7) for fixed R, m, m' constitutes a function on the group $O(2, 1)$. One can consider now the Hilbert space \mathcal{K} of all (Lebesgue) square-integrable functions on $O(2, 1)$, the integration being the usual left- and right-invariant one on the group. It has been shown by Bargmann⁵ that the set of matrix elements (1.7), with R restricted to run over the set of continuous nonexceptional UIR's and the set of discrete UIR's of $O(2, 1)$, forms a complete orthonormal basis for the Hilbert space \mathcal{K} . On the one hand, this means that, given any square integrable function $f(g)$ on the group $O(2, 1)$, one can expand it in the form

$$f(g) = \sum_{mm'} \int dR f_{mm'}(R) \mathcal{U}_{mm'}^{(R)}(g), \quad (1.8)$$

where the integration with respect to R symbolically stands for the process of summing over the discrete UIR's, as well as the process of integration with respect to the continuous parameter that labels the set of continuous nonexceptional UIR's. On the other hand, if the integration over $O(2, 1)$ is denoted by the symbol dg , one has

$$\int dg \overline{\mathcal{U}_{m_1' m_2'}^{(R')}(g)} \mathcal{U}_{m_1 m_2}^{(R)}(g) = \delta_{m_1' m_1} \delta_{m_2' m_2} \delta(R', R). \quad (1.9)$$

Once again, in (1.9), R and R' are restricted to run over the sets of UIR's that appear in (1.8); and the symbol $\delta(R', R)$ is a generalized one, denoting a Kronecker delta if R and R' are both discrete UIR's, denoting zero if one is a discrete and the other a continuous UIR, and denoting a delta-function of the Dirac type if both R and R' are continuous (nonexceptional) UIR's. Thus for this subset of UIR's of $O(2, 1)$, one has many of the properties that characterize the set of all UIR's of a compact Lie group as embodied in the Peter-Weyl theorem⁶: one has both a completeness and an orthogonality relation, (1.8) and (1.9). However, neither of these properties obtains for the matrix elements (1.7) when R stands for a UIR of the *exceptional class*. On the one hand, it is not possible to prove that these matrix elements are orthogonal to the matrix elements belonging to other inequivalent UIR's (like the continuous non-exceptional and the discrete ones) because the relevant

integrals diverge; the matrix elements (1.7) do not belong to \mathcal{K} when R is an exceptional UIR. But on the other hand, they are not needed to span the Hilbert space \mathcal{K} . These facts will be relevant in our analysis of the supplementary series of UIR's of $O(3, 1)$.

Next let us briefly summarize the UIR's of $O(3, 1)$ in which we are interested. As is well known, the group $O(3, 1)$ is the group of all real linear unimodular transformations on four real variables $(\chi_j, j = 1, 2, 3, 4)$ that leave invariant the quadratic form

$$\chi_1^2 + \chi_2^2 + \chi_3^2 - \chi_4^2,$$

and do not change the sign of χ_4 . The covering group of $O(3, 1)$ is $SL(2, C)$, the group of all complex unimodular matrices in two dimensions. Elements g in $SL(2, C)$ are in one-to-one correspondence with matrices in the following way:

$$g \rightarrow \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}, \quad \alpha\delta - \beta\gamma = 1; \quad (1.10)$$

the homomorphism from $SL(2, C)$ to $O(3, 1)$ is two-to-one. The Lie algebra of $SL(2, C)$ is six dimensional; and the basic elements $L_j, N_j, j = 1, 2, 3$ obey the following commutation rules:

$$\begin{aligned} [L_j, L_k] &= i\epsilon_{jkl} L_l; \\ [L_j, N_k] &= i\epsilon_{jkl} N_l; \\ [N_j, N_k] &= -i\epsilon_{jkl} L_l. \end{aligned} \quad (1.11)$$

The elements L_j generate the $SU(2)$ subgroup of $SL(2, C)$, while the elements L_3, N_1, N_2 generate an $SU(1, 1)$ subgroup of $SL(2, C)$. The two quadratic Casimir invariants of $SL(2, C)$ are

$$C_1 = N_j N_j - L_j L_j \quad \text{and} \quad C_2 = N_j L_j, \quad (1.12)$$

while the Casimir invariants of the $SU(2)$ and $SU(1, 1)$ subgroups are given respectively by

$$L^2 = L_j L_j \quad \text{and} \quad Q = N_1^2 + N_2^2 - L_3^2. \quad (1.13)$$

UIR's of $SL(2, C)$ fall into two categories, the principal series and the supplementary series.⁷ UIR's of the latter series are labeled by a continuous real parameter ρ which takes on values in the *open interval* $0 < \rho < 1$:

$$C_1 = 1 - \rho^2, \quad C_2 = 0; \quad 0 < \rho < 1. \quad (1.14)$$

There is one such UIR for each value of ρ in the given range. In all these UIR's, the spectrum of finite-dimensional UIR's of the subgroup $SU(2)$ is specified

⁵ V. Bargmann, Ref. 1, pp. 632-639. The statement made in the text is strictly true for single-valued representations of $O(2, 1)$. For single-valued representations of $SU(1, 1)$, however, the situation is that the matrix elements belonging to the two discrete representations $D_{\frac{1}{2}}^{\pm}$ are outside of \mathcal{K} and are not needed to span \mathcal{K} .

⁶ H. Weyl, Ann. Math. 35, 486 (1934).

⁷ M. A. Naimark, Ref. 1, p. 170. (This series is called the complementary series in this book.) I. M. Gel'fand, R. A. Minlos, and Z. Ya. Shapiro, Ref. 1, p. 247. There is an error in the equations on p. 248 of this book. The parameter ρ should be replaced by $-\rho$.

by the spectrum of the operator L^2 :

$$L^2 = l(l + 1), \quad l = 0, 1, 2, \dots, \infty. \quad (1.15)$$

Every UIR of $SU(2)$ of dimensionality $(2l + 1)$ for each nonnegative integral value of l appears exactly once. These are, then, single-valued representations of $O(3, 1)$. We may introduce a basis made up of eigenvectors of L^2 and L_3 :

$$L^2 |l, m\rangle = l(l + 1) |l, m\rangle; \quad L_3 |l, m\rangle = m |l, m\rangle;$$

$$\langle l' m' | l m \rangle = \delta_{l'l} \delta_{m'm}; \quad -l \leq m \leq l. \quad (1.16)$$

In this basis, the matrix elements of L_j and N_j are the following:

$$\left[L_0 = L_3, L_{+1} = -\frac{1}{\sqrt{2}}(L_1 + iL_2), \right.$$

$$\left. L_{-1} = \frac{1}{\sqrt{2}}(L_1 - iL_2), \text{ etc.} \right];$$

$$\langle l' m' | L_M | l m \rangle = [l(l + 1)]^{\frac{1}{2}} \delta_{l'l} C_m^l \frac{1}{M} \frac{1}{m'},$$

$$\langle l + 1, m' | N_M | l m \rangle = -[(l + 1)(2l + 1)]^{\frac{1}{2}} C_m^l \frac{1}{M} \frac{1}{m'} C_{l+1},$$

$$\langle l m' | N_M | l m \rangle = 0,$$

$$\langle l - 1, m' | N_M | l m \rangle = -[l(2l + 1)]^{\frac{1}{2}} C_m^l \frac{1}{M} \frac{1}{m'} C_l,$$

$$C_l = i \left[\frac{l^2 - \rho^2}{4l^2 - 1} \right]^{\frac{1}{2}}. \quad (1.17)$$

2. ANALYSIS OF THE EXCEPTIONAL UIR'S OF $O(2, 1)$

The UIR's of the exceptional series may be constructed explicitly in Hilbert spaces \mathcal{H}_σ consisting of a certain class of functions $f(\varphi)$ on the unit circle ($0 \leq \varphi \leq 2\pi$).⁸ The scalar product of two elements f, h and the norm of f , are given by

$$(f, h)_\sigma = (2\pi)^{-2} \int_0^{2\pi} d\varphi' \int_0^{2\pi} d\varphi f(\varphi') \overline{L_\sigma(\varphi' - \varphi) h(\varphi)};$$

$$\|f\|_\sigma = (f, f)_\sigma^{\frac{1}{2}} < \infty;$$

$$L_\sigma(\varphi' - \varphi) = (2\pi)^{\frac{1}{2}} \frac{\Gamma(\sigma + \frac{1}{2})}{2^\sigma \Gamma(\sigma)} [1 - \cos(\varphi' - \varphi)]^{\sigma - \frac{1}{2}}; \quad 0 < \sigma < \frac{1}{2}. \quad (2.1)$$

\mathcal{H}_σ consists of all functions $f(\varphi)$ with finite norm. The class of functions constituting \mathcal{H}_σ , as well as the definition of the scalar product, depends on σ . The unitary operator $U(g)$ representing the element g [given by (1.1)] acts as follows:

$$[U(g)f](\varphi) = |\bar{\alpha} - \beta e^{i\varphi}|^{-1-2\sigma} f[\psi_\sigma(\varphi)],$$

$$e^{i\psi_\sigma(\varphi)} = (\alpha e^{i\varphi} - \bar{\beta}) / (\bar{\alpha} - \beta e^{i\varphi}), \quad (2.2)$$

$$0 \leq \varphi, \quad \psi_\sigma(\varphi) \leq 2\pi.$$

The generators J_0, J_1, J_2 are represented by the following differential operators:

$$J_0 = -i \frac{d}{d\varphi},$$

$$J_1 = -i \cos \varphi \frac{d}{d\varphi} + i(\sigma + \frac{1}{2}) \sin \varphi, \quad (2.3)$$

$$J_2 = -i \sin \varphi \frac{d}{d\varphi} - i(\sigma + \frac{1}{2}) \cos \varphi.$$

The functions corresponding to the unit vectors $|m\rangle$ in (1.5) are

$$|m\rangle \rightarrow [\lambda_m(\sigma)]^{-\frac{1}{2}} e^{im\varphi}, \quad m = 0, \pm 1, \pm 2, \dots$$

and

$$\lambda_m(\sigma) = \frac{\Gamma(\frac{1}{2} + \sigma) \Gamma(|m| + \frac{1}{2} - \sigma)}{\Gamma(\frac{1}{2} - \sigma) \Gamma(|m| + \frac{1}{2} + \sigma)} > 0. \quad (2.4)$$

As noted by Bargmann,⁸ the Hilbert space \mathcal{H}_σ contains, as a *dense subset*, the set of all *square-integrable* functions of φ , namely all functions $f(\varphi)$ for which

$$\|f\|^2 \equiv (2\pi)^{-1} \int_0^{2\pi} |f(\varphi)|^2 d\varphi < \infty. \quad (2.5)$$

Such functions form a Hilbert space \mathcal{H} with respect to the norm $\|f\|$. Every function in \mathcal{H} can be expressed as a Fourier Series:

$$f(\varphi) = \sum_{m=-\infty}^{+\infty} f_m e^{im\varphi}, \quad (2.6)$$

and we have

$$\|f\|^2 = \sum_m |f_m|^2 < \infty. \quad (2.7)$$

The norm of f in the space \mathcal{H}_σ is then given by

$$\|f\|_\sigma^2 = \sum_m \lambda_m(\sigma) |f_m|^2. \quad (2.8)$$

The basic reason why \mathcal{H}_σ contains the set \mathcal{H} as a dense subset is that, for large $|m|$, $\lambda_m(\sigma)$ goes to zero like $|m|^{-2\sigma}$. The set \mathcal{H} is not closed with respect to the norm $\|f\|_\sigma$, but yields \mathcal{H}_σ on completion with respect to it. \mathcal{H}_σ consists of *all* sequences $\{f_m\}$ for which the norm $\|f\|_\sigma$, as defined in (2.8), is finite.

We now wish to diagonalize the generator J_2 . As in (A), we first change variables from φ to q , mapping the region $0 \leq \varphi \leq 2\pi$ into *two* real lines in q :

$$e^q = \tan \varphi/2: \quad 0 \leq \varphi \leq \pi \rightarrow -\infty < q < \infty$$

and

$$e^{-q} = \tan(\varphi - \pi)/2:$$

$$\pi \leq \varphi \leq 2\pi \rightarrow \infty > q > -\infty. \quad (2.9)$$

Thus the upper and the lower halves of the circumference of the unit circle get mapped onto one real

⁸ V. Bargmann, Ref. 1, pp. 616 ff.

line each. We then replace the function $f(\varphi)$ by a pair of functions of q :

$$\begin{aligned} f_1(q) &= [\cosh q]^{-\frac{1}{2}-\sigma} f(\varphi), \quad 0 \leq \varphi \leq \pi; \\ f_2(q) &= [\cosh q]^{-\frac{1}{2}-\sigma} f(\varphi), \quad \pi \leq \varphi \leq 2\pi, \end{aligned} \quad (2.10)$$

so that a given vector f in \mathcal{H}_σ is now represented by the two functions $f_r(q)$. So far these steps are the same as in (A). If we now express the scalar product (2.1) in terms of $f_r(q)$ and $h_r(q)$, we discover the appearance of cross terms between $r = 1$ and $r = 2$. These may be eliminated by working with simple linear combinations of $f_1(q)$ and $f_2(q)$:

$$f_\pm(q) = \frac{1}{2}[f_1(q) \pm f_2(q)]. \quad (2.11)$$

In terms of these combinations, we find that

$$\begin{aligned} (f, h)_\sigma &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} dq' dq [\overline{f_+(q')} K_+(q' - q) h_+(q) \\ &\quad + \overline{f_-(q')} K_-(q' - q) h_-(q)], \\ K_\pm(q) &= \mu(\sigma) [\{\cosh q - 1\}^{\sigma-\frac{1}{2}} \pm \{\cosh q + 1\}^{\sigma-\frac{1}{2}}], \\ \mu(\sigma) &= (2\pi)^{-\frac{1}{2}} 2^{-\sigma} \frac{\Gamma(\sigma + \frac{1}{2})}{\Gamma(\sigma)}. \end{aligned} \quad (2.12)$$

Functions $f_+(q)$ and $f_-(q)$ correspond, respectively, to functions $f(\varphi)$ which are even and odd under the substitution $\varphi \rightarrow 2\pi - \varphi$. (2.12) shows that the Hilbert space \mathcal{H}_σ is the direct sum of two orthogonal subspaces $\mathcal{H}_\sigma^{(\pm)}$; $\mathcal{H}_\sigma^{(+)}$ contains vectors for which $f_-(q)$ vanishes, while $\mathcal{H}_\sigma^{(-)}$ contains vectors with $f_+(q) = 0$. We will see that the kernels K_\pm are both positive definite.

The generators may now be written as differential operators acting on the column vectors:

$$\begin{pmatrix} f_+(q) \\ f_-(q) \end{pmatrix}.$$

Quite easily we find that

$$\begin{aligned} J_0 &= \left[-i \cosh q \frac{d}{dq} - i(\sigma + \frac{1}{2}) \sinh q \right] \otimes \sigma, \\ J_1 &= \left[i \sinh q \frac{d}{dq} + i(\sigma + \frac{1}{2}) \cosh q \right] \otimes \sigma_1, \quad (2.13) \\ J_2 &= -i \frac{d}{dq} \otimes \mathbf{1}; \quad \sigma_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \end{aligned}$$

The operators J_0 and J_1 do not leave $\mathcal{H}_\sigma^{(+)}$ and $\mathcal{H}_\sigma^{(-)}$ invariant. However, these subspaces are invariant under J_2 . This result follows directly if we express the operators $U(\exp i\tau J_2)$, starting from (2.2) and carrying out the changes of variable. One has

$$[U(\exp i\tau J_2) f]_\pm(q) = f_\pm(q + \tau). \quad (2.14)$$

Therefore, to diagonalize J_2 , we must pass to the

Fourier transforms of $f_\pm(q)$. Consider first the set of functions $f(q)$ which are square-integrable with respect to q .

$$(\|f\|')^2 \equiv \int_{-\infty}^{\infty} |f(q)|^2 dq < \infty. \quad (2.15)$$

This set of functions forms a Hilbert space \mathcal{H}^1 with respect to the norm $\|f\|'$. Let us pick a pair of functions $f_\pm(q)$, each of which is an element of \mathcal{H}^1 . We will show that they determine a vector f in \mathcal{H}_σ , with a finite norm $\|f\|_\sigma$. Both $f_+(q)$ and $f_-(q)$ can be written as Fourier transforms of other square-integrable functions:

$$f_\pm(q) = (2\pi)^{-\frac{1}{2}} \int_{-\infty}^{\infty} e^{iqp} f'_\pm(p) dp. \quad (2.16)$$

Substituting these in (2.12), we find that

$$\begin{aligned} \|f\|_\sigma^2 &= \int_{-\infty}^{\infty} dp \left[|f'_+(p)|^2 \int_{-\infty}^{\infty} e^{-ipq} K_+(q) dq \right. \\ &\quad \left. + |f'_-(p)|^2 \int_{-\infty}^{\infty} e^{-ipq} K_-(q) dq \right]. \end{aligned} \quad (2.17)$$

This is valid if the indicated Fourier transforms of $K_\pm(q)$ exist and are well-behaved functions of p . We find in fact that these transforms do exist and have just the desired properties⁹:

$$\begin{aligned} \int_{-\infty}^{\infty} e^{-ipq} K_\pm(q) dq &= \mu'(\sigma) \lambda_\pm(p), \\ \lambda_\pm(p) &= \Gamma(\frac{1}{2} - \sigma + ip) \Gamma(\frac{1}{2} - \sigma - ip) \\ &\quad \times [\cosh \pi p \pm \sin \pi \sigma] > 0, \\ \mu'(\sigma) &= (2\pi)^{-\frac{1}{2}} 2^{3(1-\sigma)} [\Gamma(2\sigma) \Gamma(\sigma + \frac{1}{2}) / \Gamma(\sigma)] \cos \pi \sigma. \end{aligned} \quad (2.18)$$

The asymptotic behavior of $\lambda_\pm(p)$ for large $|p|$ is given by

$$\lambda_\pm(p) \xrightarrow{|p| \rightarrow \infty} \text{const } |p|^{-2\sigma}, \quad (2.19)$$

so that the finiteness of the norm $\|f\|_\sigma$ is guaranteed because (2.15) implies that

$$\int_{-\infty}^{\infty} |f'_\pm(p)|^2 dp < \infty. \quad (2.20)$$

The positive definiteness of the functions $\lambda_\pm(p)$ implies the positive definiteness of the kernels $K_\pm(q)$. The asymptotic behavior of $\lambda_\pm(p)$ as exhibited by (2.19) is analogous to the behavior of the numbers $\lambda_m(\sigma)$ in (2.4) for large $|m|$. Just as we were able to conclude that the space \mathcal{H} of square-integrable functions of φ is contained as a dense subset in \mathcal{H}_σ , we may now conclude that the set \mathcal{H}^1 of all square-integrable functions of q is a dense subset of both

⁹ Bateman Manuscript Project, A. Erdelyi, Ed. (McGraw-Hill Book Company, Inc., New York, 1953)

$\mathcal{H}_\sigma^{(+)}$ and $\mathcal{H}_\sigma^{(-)}$. With the aid of (2.12), or (2.17), we are led in a natural way to norms $\|f\|_\sigma^{(\pm)}$ defined in $\mathcal{H}_\sigma^{(+)}$ and $\mathcal{H}_\sigma^{(-)}$, respectively. The space \mathcal{H}^1 is not closed with respect to either of these norms. Completion of \mathcal{H}^1 with respect to $\|f\|_\sigma^{(+)}$ yields the Hilbert space $\mathcal{H}_\sigma^{(+)}$; and with respect to $\|f\|_\sigma^{(-)}$ it yields $\mathcal{H}_\sigma^{(-)}$. $\mathcal{H}_\sigma^{(+)}(\mathcal{H}_\sigma^{(-)})$ consists of all functions $f'_\pm(p)$ ($f'_-(p)$) for which the norm $\|f\|_\sigma^{(+)}(\|f\|_\sigma^{(-)})$, as defined by (2.17) and (2.18), is finite. The original space \mathcal{H}_σ is obtained as the direct sum of $\mathcal{H}_\sigma^{(+)}$ and $\mathcal{H}_\sigma^{(-)}$.

The generator J_2 is diagonal in the “ p basis”:

$$[J_2 f]_\pm(p) = p f'_\pm(p). \tag{2.21}$$

We see that the spectrum of J_2 consists of all real numbers, with each eigenvalue appearing twice, as noted by Bargmann. Concerning the forms of J_0 and J_1 in this basis, we argue in the same manner as in (A) and reach the result that, if a vector f is in the domain of J_0 and J_1 , then its wavefunctions $f'_\pm(p)$ are boundary values of analytic functions of p , possessing unique analytic continuations into the complex p plane, at least up to $\text{Im } p = \pm 1$. Then the wavefunctions of the vectors $J_0 f$ and $J_1 f$ are given by

$$\begin{aligned} [J_0 f]'_\pm(p) &= \frac{1}{2}[p + i(\sigma - \frac{1}{2})]f'_\mp(p - i) \\ &\quad + \frac{1}{2}[p - i(\sigma - \frac{1}{2})]f'_\mp(p + i), \\ [J_1 f]'_\pm(p) &= \frac{1}{2}[p + i(\sigma - \frac{1}{2})]f'_\mp(p - i) \\ &\quad - \frac{1}{2}[p - i(\sigma - \frac{1}{2})]f'_\mp(p + i). \end{aligned} \tag{2.22}$$

Conversely, J_0 and J_1 can only act upon vectors f whose wavefunctions $f'_\pm(p)$ possess all these properties and are such that the right-hand sides of (2.2.2) represent normalizable wavefunctions.

3. SUPPLEMENTARY SERIES OF UIR'S OF $O(3, 1)$

The supplementary representations of $O(3, 1)$ are similar to the exceptional representations of $O(2, 1)$ in that they, too, are most naturally expressed in a space of functions with a nonlocal metric.⁷ Let \mathcal{H}_ρ be the Hilbert space of functions of a complex variable $f(\xi)$ with a finite norm $\|f\|_\rho$, this norm being defined by

$$\|f\|_\rho^2 = (f, f)_\rho = \int d^2\xi' \int d^2\xi \overline{f(\xi')} |\xi' - \xi|^{-2+2\rho} f(\xi) \tag{3.1}$$

$0 < \rho < 1.$

[These functions $f(\xi)$ are not analytic functions of ξ , but are complex-valued functions of the real and imaginary parts of ξ . The integrations in (3.1) extend over the entire ξ plane.] Once again, the possible functions $f(\xi)$ appearing in \mathcal{H}_ρ will depend on ρ . If g is an element of $SL(2, C)$ given by (1.10), it is represented by a unitary operator $U(g)$ acting on \mathcal{H}_ρ

as follows:

$$[U(g)f](\xi) = |\delta + \beta\xi|^{-2+2\rho} f\left(\frac{\alpha\xi + \gamma}{\delta + \beta\xi}\right). \tag{3.2}$$

It is possible to write down the differential operators that represent the generators L_j and N_j of $O(3, 1)$, in terms of $\xi = x + iy$. However, we shall not do so here; we merely note that they may be obtained from the formulas given in (B) by the formal replacement $\rho \rightarrow i\rho$.

It is useful to analyze first some properties of these UIR's with respect to the $SU(2)$ subgroup. For this purpose, exactly as in (B), we express the UIR in a space of functions defined on a unit sphere in three dimensions. The variable ξ is related to the angles θ, φ on the sphere by

$$\xi = x + iy = \cot \frac{\theta}{2} e^{i\varphi}; \quad 0 \leq \theta \leq \pi, \quad 0 \leq \varphi \leq 2\pi. \tag{3.3}$$

The function $f(\xi)$ is related to a function $\tilde{f}(\theta, \varphi)$ via

$$f(\xi) = 2 \left[\sin \frac{\theta}{2} \right]^{2+2\rho} \tilde{f}(\theta, \varphi). \tag{3.4}$$

Then the scalar product of two elements f and h in \mathcal{H}_ρ becomes

$$\begin{aligned} (f, h)_\rho &= 2^{-1-\rho} \iint d \cos \theta' d \cos \theta d\varphi' d\varphi \\ &\quad \times \overline{\tilde{f}(\theta', \varphi')} L_\rho(\theta' \varphi', \theta \varphi) \tilde{h}(\theta, \varphi), \\ L_\rho(\theta' \varphi', \theta \varphi) &= [1 - \cos \theta' \cos \theta - \sin \theta' \sin \theta \cos(\varphi' - \varphi)]^{-1+\rho}. \end{aligned} \tag{3.5}$$

The generators acting on $\tilde{f}(\theta, \varphi)$ are as follows¹⁰:

$$\begin{aligned} L_1 &= i \sin \varphi \frac{\partial}{\partial \theta} + i \cos \varphi \cot \theta \frac{\partial}{\partial \varphi}, \\ L_2 &= -i \cos \varphi \frac{\partial}{\partial \theta} + i \sin \varphi \cot \theta \frac{\partial}{\partial \varphi}, \\ L_3 &= -i \frac{\partial}{\partial \varphi}; \\ N_1 &= -i(\rho + 1) \sin \theta \cos \varphi \\ &\quad + i \cos \theta \cos \varphi \frac{\partial}{\partial \theta} - i \frac{\sin \varphi}{\sin \theta} \frac{\partial}{\partial \varphi}, \\ N_2 &= -i(\rho + 1) \sin \theta \sin \varphi \\ &\quad + i \cos \theta \sin \varphi \frac{\partial}{\partial \theta} + i \frac{\cos \varphi}{\sin \theta} \frac{\partial}{\partial \varphi}, \\ N_3 &= -i(\rho + 1) \cos \theta - i \sin \theta \frac{\partial}{\partial \theta}. \end{aligned} \tag{3.6}$$

¹⁰ The expressions obtained for L_j are exactly the same as for the angular momentum of a single particle in quantum mechanics. If R denotes a spatial rotation, then the unitary operator $U(R)$ has a very simple effect on the functions $\tilde{f}(\theta, \varphi)$: $[U(R)\tilde{f}]^-(\theta, \varphi) = \tilde{f}[R^{-1}(\theta, \varphi)]$.

Equations (3.5) and (3.6) are the analogs to (2.1) and (2.3).

The spectrum of UIR's of $SU(2)$ contained in this UIR of $O(3, 1)$ is made explicit by an orthogonal-function expansion of the kernel $L_\rho(\theta'\varphi', \theta\varphi)$ in terms of eigenfunctions of the differential operators L^2 and L_3 . The details of this expansion are given in Appendix A. Here we quote the result, which is as follows for $0 < \rho < 1$:

$$L_\rho(\theta'\varphi', \theta\varphi) = 2^{\rho+1}\pi \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \times \sum_{l=0}^{\infty} \sum_{m=-l}^{+l} \frac{\Gamma(l-\rho+1)}{\Gamma(l+\rho+1)} \overline{Y_l^m(\theta, \varphi)} Y_l^m(\theta', \varphi'). \quad (3.7)$$

Notice these two features: (i) the l -dependent coefficients of the spherical harmonics are always positive for $0 < \rho < 1$; (ii) for large values of l , these coefficients go to zero because

$$\frac{\Gamma(l-\rho+1)}{\Gamma(l+\rho+1)} \xrightarrow{l \rightarrow \infty} \text{const } l^{-2\rho}. \quad (3.8)$$

These properties are exactly similar to the case of $O(2, 1)$, where the coefficients $\lambda_m(\sigma)$, relating to a Fourier-series expansion of the kernel $L_\sigma(\varphi' - \varphi)$, tended to zero for $|m| \rightarrow \infty$. Let now \mathcal{H} be the space of all functions $\tilde{f}(\theta, \varphi)$ that are square-integrable over the unit sphere in the following sense:

$$\|\tilde{f}\|^2 \equiv \int_0^\pi \sin \theta \, d\theta \int_0^{2\pi} d\varphi |\tilde{f}(\theta, \varphi)|^2 < \infty. \quad (3.9)$$

The space \mathcal{H} forms a Hilbert space with respect to the norm $\|\tilde{f}\|$, and any element \tilde{f} in \mathcal{H} possesses an expansion in terms of the spherical harmonics:

$$\tilde{f}(\theta, \varphi) = \sum_{l=0}^{\infty} \sum_{m=-l}^{+l} f_{lm} Y_l^m(\theta, \varphi),$$

$$f_{lm} = \int_0^\pi \sin \theta \, d\theta \int_0^{2\pi} d\varphi \overline{Y_l^m(\theta, \varphi)} \tilde{f}(\theta, \varphi), \quad (3.10)$$

$$\|\tilde{f}\|^2 = \sum_{l=0}^{\infty} \sum_{m=-l}^{+l} |f_{lm}|^2 < \infty.$$

Such a function $\tilde{f}(\theta, \varphi)$ automatically belongs also to \mathcal{H}_ρ because, by virtue of (3.7), (3.8), and (3.10), it has a finite norm $\|f\|_\rho$ in \mathcal{H}_ρ :

$$\|f\|_\rho^2 = \pi \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \sum_l \sum_m \frac{\Gamma(l-\rho+1)}{\Gamma(l+\rho+1)} |f_{lm}|^2 < \infty. \quad (3.11)$$

The space \mathcal{H} is thus contained as a dense subset in \mathcal{H}_ρ ; \mathcal{H} is not closed with respect to the norm $\|f\|_\rho$, but yields \mathcal{H}_ρ on completion with respect to it. Vectors in \mathcal{H}_ρ correspond to all sequences f_{lm} for which (3.11)

is finite, while vectors in \mathcal{H} correspond to those sequences for which the expression appearing in (3.10) is also finite. We can identify, from (3.7), the functions of θ, φ which represent the states $|lm\rangle$ of (1.16). They are, apart from phase factors and an over-all normalization, given by

$$|lm\rangle \rightarrow \left[\frac{\Gamma(l+\rho+1)}{\Gamma(l-\rho+1)} \right]^{\frac{1}{2}} Y_l^m(\theta, \varphi). \quad (3.12)$$

All the properties elucidated above, relating to the decomposition of the UIR's under $SU(2)$, are true for all ρ in the range $0 < \rho < 1$ and are included here to display the similarity to the properties of the exceptional representations of $O(2, 1)$. Another reason for exhibiting these properties is that we would now like to achieve the decomposition of the UIR under the noncompact subgroup $O(2, 1)$ by a very similar method. As a first step, we adopt the approach in (B), and change variables from the unit sphere to two planes. We define radial and polar coordinates r, φ as follows:

$$0 \leq \theta \leq \frac{\pi}{2} : \tan \frac{\theta}{2} = \tanh \frac{r}{2}; \quad 0 \leq r < \infty \quad (3.13)$$

$$\frac{\pi}{2} \leq \theta \leq \pi : \cot \frac{\theta}{2} = \tanh \frac{r}{2}; \quad \infty > r \geq 0.$$

The angle coordinate φ remains unchanged. According to (3.13), the upper and the lower hemispheres get mapped onto one plane each. Still following the lines of (B), we define one function on each plane to replace $\tilde{f}(\theta, \varphi)$:

$$f_1(r, \varphi) = [\cosh r]^{-1-\rho} \tilde{f}(\theta, \varphi) : \quad 0 \leq \theta \leq \pi/2$$

$$f_2(r, \varphi) = [\cosh r]^{-1-\rho} \tilde{f}(\theta, \varphi) : \quad \pi/2 \leq \theta \leq \pi. \quad (3.14)$$

As we should expect, the scalar product $(f, h)_\rho$ contains cross terms between $f_1(r, \varphi)$, $h_2(r, \varphi)$, etc. These are removed by working with the combinations

$$f_\pm(r, \varphi) : \quad f_\pm(r, \varphi) = \frac{1}{2} [f_1(r, \varphi) \pm f_2(r, \varphi)]. \quad (3.15)$$

These functions, $f_+(r, \varphi)$ and $f_-(r, \varphi)$, correspond to functions $\tilde{f}(\theta, \varphi)$, which are even and odd respectively under the substitution $\theta \rightarrow \pi - \theta$. Then the scalar product is

$$(f, h)_\rho = 2^{-1-\rho} \int_0^\infty \int_0^\infty \sinh r' \sinh r \, dr' \, dr \int_0^{2\pi} \int_0^{2\pi} d\varphi' \, d\varphi$$

$$[f_+(r', \varphi') K_+(r', r, \varphi' - \varphi) h_+(r, \varphi)$$

$$+ f_-(r', \varphi') K_-(r', r, \varphi' - \varphi) h_-(r, \varphi)],$$

$$K_\pm(r', r, \varphi) = [\eta(r', r, \varphi) - 1]^{\rho-1}$$

$$\pm [\eta(r', r, \varphi) + 1]^{\rho-1},$$

$$\eta(r', r, \varphi) = \cosh r' \cosh r - \sinh r' \sinh r \cos \varphi. \quad (3.16)$$

Subject to verification of the fact that K_{\pm} are positive definite kernels, we see that \mathcal{H}_{ρ} has been expressed as the direct sum of two orthogonal subspaces $\mathcal{H}_{\rho}^{(+)}$ and $\mathcal{H}_{\rho}^{(-)}$, corresponding, respectively, to vectors in \mathcal{H}_{ρ} with $f_{-}(r, \varphi) = 0$ and $f_{+}(r, \varphi) = 0$.

The differential operators for L_j, N_j acting on $f_{\pm}(r, \varphi)$, written as a column vector

$$\begin{pmatrix} f_{+}(r, \varphi) \\ f_{-}(r, \varphi) \end{pmatrix},$$

are

$$\begin{aligned} L_3 &= -i \frac{\partial}{\partial \varphi} \otimes \mathbf{1}, \\ N_1 &= \left[i \cos \varphi \frac{\partial}{\partial r} - i \sin \varphi \coth r \frac{\partial}{\partial \varphi} \right] \otimes \mathbf{1}, \\ N_2 &= \left[i \sin \varphi \frac{\partial}{\partial r} + i \cos \varphi \coth r \frac{\partial}{\partial \varphi} \right] \otimes \mathbf{1}, \\ N_3 &= \left[-i(\rho + 1) \cosh r - i \sinh r \frac{\partial}{\partial r} \right] \otimes \sigma_1, \\ L_1 &= \left[i(\rho + 1) \sin \varphi \sinh r \right. \\ &\quad \left. + i \sin \varphi \cosh r \frac{\partial}{\partial r} + i \frac{\cos \varphi}{\sinh r} \frac{\partial}{\partial \varphi} \right] \otimes \sigma_1, \\ L_2 &= \left[-i(\rho + 1) \cos \varphi \sinh r \right. \\ &\quad \left. - i \cos \varphi \cosh r \frac{\partial}{\partial r} + i \frac{\sin \varphi}{\sinh r} \frac{\partial}{\partial \varphi} \right] \otimes \sigma_1. \end{aligned} \tag{3.17}$$

One can check in either of two ways that the subgroup $O(2, 1)$ generated by L_3, N_1 , and N_2 leaves $\mathcal{H}_{\rho}^{(+)}$ and $\mathcal{H}_{\rho}^{(-)}$ invariant. One way is to realize that these three generators, as given in (3.6), are invariant under the replacement $\theta \rightarrow \pi - \theta$. Another way is to work out explicitly the form of the unitary operator $U(g)$ for g , an element of $O(2, 1)$. Then one sees that, for any such g ,

$$[U(g)f]_{\pm}(r, \varphi) = f_{\pm}(r_g, \varphi_g), \tag{3.18}$$

the arguments r_g, φ_g of f_{\pm} on the right-hand side being some definite functions of r, φ , and g .¹¹ In particular, this "change of argument" is the same for both functions f_{+} and f_{-} , and there is no factor multiplying the functions f_{\pm} on the right-hand side of (3.18).

The decomposition of the subspaces $\mathcal{H}_{\rho}^{(+)}$ and $\mathcal{H}_{\rho}^{(-)}$ into subspaces irreducible under $O(2, 1)$ can now be achieved by an orthogonal function expansion of the kernels $K_{\pm}(r', r, \varphi' - \varphi)$ in terms of eigenfunctions of Q and L_3 . This would be analogous to (3.7). It should be pointed out that although the subgroup acts "in the

same way" on functions $f_{+}(r, \varphi)$ and $f_{-}(r, \varphi)$, as evidenced by (3.18), this does not mean that the representations of $O(2, 1)$ appearing in $\mathcal{H}_{\rho}^{(+)}$ and $\mathcal{H}_{\rho}^{(-)}$ are the same. This is just because the scalar products defined in $\mathcal{H}_{\rho}^{(+)}$ and $\mathcal{H}_{\rho}^{(-)}$ are different. The differential operator Q is

$$Q = - \left[\frac{\partial^2}{\partial r^2} + \coth r \frac{\partial}{\partial r} + \frac{1}{\sinh^2 r} \frac{\partial^2}{\partial \varphi^2} \right] \otimes \mathbf{1}, \tag{3.19}$$

and the relevant eigenfunctions of Q and L_3 are generalizations of the spherical harmonics $Y_l^m(\theta, \varphi)$ for complex l and with $\cos \theta$ replaced by $\cosh r$. The decomposition of the kernels K_{\pm} in terms of these eigenfunctions is carried out in Appendix A. Since the results of this decomposition are very different for $0 < \rho \leq \frac{1}{2}$ and for $\frac{1}{2} < \rho < 1$, we now discuss these two cases separately.

Case I: $0 < \rho \leq \frac{1}{2}$.

For such values of ρ (notice that $\rho = \frac{1}{2}$ is included), we find

$$\begin{aligned} K_{\pm}(r', r, \varphi' - \varphi) &= \frac{2^{\rho}}{\pi} \frac{\Gamma(\rho)}{\Gamma(1 - \rho)} \\ &\quad \times \int_0^{\infty} \sum_{m=-\infty}^{+\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \lambda_{\pm}(\mathcal{S}) \\ &\quad \times Y_{-\frac{1}{2}+i\mathcal{S}}^{m}(r', \varphi') \overline{Y_{-\frac{1}{2}+i\mathcal{S}}^{m}(r, \varphi)}, \\ \lambda_{\pm}(\mathcal{S}) &= \Gamma(\tfrac{1}{2} - \rho + i\mathcal{S}) \Gamma(\tfrac{1}{2} - \rho - i\mathcal{S}) \\ &\quad \times [\cosh \pi \mathcal{S} \pm \sin \pi \rho] > 0, \\ Y_{-\frac{1}{2}+i\mathcal{S}}^{m}(r, \varphi) &= \left[\frac{\cosh \pi \mathcal{S}}{\pi} \Gamma(\tfrac{1}{2} + i\mathcal{S} - m) \Gamma(\tfrac{1}{2} - i\mathcal{S} - m) \right]^{\frac{1}{2}} \\ &\quad \times e^{im\varphi} P_{-\frac{1}{2}+i\mathcal{S}}^{(m)}(\cosh r). \end{aligned} \tag{3.20}$$

The positive definiteness of the kernels K_{\pm} is thereby established. We see also that the eigenfunctions of Q appearing in the decompositions of K_{\pm} correspond to the following eigenvalue spectrum of Q :

$$Q = \tfrac{1}{4} + \mathcal{S}^2, \quad 0 \leq \mathcal{S} < \infty. \tag{3.21}$$

This spectrum is the same for both decompositions, K_{+} and K_{-} . Therefore we have the following:

Theorem I: UIR's of $O(3, 1)$ belonging to the supplementary series, for values of the parameter ρ in the range $0 < \rho \leq \frac{1}{2}$, decompose into direct integrals of UIR's of $O(2, 1)$ belonging to the continuous *nonexceptional* (integral) class, with the spectrum of the Casimir invariant Q of $O(2, 1)$ being $\frac{1}{4} \leq Q < \infty$. In the decomposition, every UIR of $O(2, 1)$ of the above-mentioned class appears *twice*.

¹¹ This is analogous to the situation described in Ref. 10.

Every UIR of $O(2, 1)$ of the continuous non-exceptional integral class appears once in the subspace $\mathcal{H}_\rho^{(+)}$ and once in $\mathcal{H}_\rho^{(-)}$. Notice that the *exceptional* representations of $O(2, 1)$ are not found in the supplementary representations of $O(3, 1)$ for $0 < \rho \leq \frac{1}{2}$. We discuss now some properties of the kernels K_\pm and the subspaces $\mathcal{H}_\rho^{(\pm)}$:

According to the statements made in Sec. 1, the matrix elements of UIR's of $O(2, 1)$, which belong to the continuous nonexceptional and the discrete classes, span the Hilbert space of square integrable functions on $O(2, 1)$. The functions $Y_{-\frac{1}{2}+iS}^m(r, \varphi)$, appearing in (3.20), are actually special cases of these matrix elements and belong to the continuous non-exceptional representations. For these matrix elements, the completeness and orthonormality properties may be stated in the following way. Let \mathcal{K} be the Hilbert space of all functions $f(r, \varphi)$ which have finite norms $\|f\|_{\mathcal{K}}$, this being defined by

$$\|f\|_{\mathcal{K}}^2 \equiv \int_0^\infty \sinh r \, dr \int_0^{2\pi} d\varphi |f(r, \varphi)|^2. \quad (3.22)$$

Then every such function $f(r, \varphi)$ may be expanded in terms of $Y_{-\frac{1}{2}+iS}^m(r, \varphi)$:

$$f(r, \varphi) = \sum_{m=-\infty}^{+\infty} \int_0^\infty dS \, S \tanh \pi S f_m(S) Y_{-\frac{1}{2}+iS}^m(r, \varphi),$$

$$f_m(S) = (2\pi)^{-1} \int_0^\infty \sinh r \, dr$$

$$\times \int_0^{2\pi} d\varphi \overline{Y_{-\frac{1}{2}+iS}^m(r, \varphi)} f(r, \varphi). \quad (3.23)$$

The norm of f is given by

$$\|f\|_{\mathcal{K}}^2 = \sum_{m=-\infty}^{+\infty} \int_0^\infty dS \, S \tanh \pi S |f_m(S)|^2. \quad (3.24)$$

If we examine the behavior of the weight functions $\lambda_\pm(S)$ appearing in the kernels K_\pm , we find that they go to zero for large S [cf. (2.18) and (2.19)]:

$$\lambda_\pm(S) \xrightarrow{S \rightarrow \infty} \text{const } S^{-2\rho}. \quad (3.25)$$

Furthermore, $\lambda_+(S)$ is finite for every finite S in the range $0 \leq S < \infty$ only if $0 < \rho < \frac{1}{2}$, while $\lambda_-(S)$ is finite for finite S if $0 < \rho \leq \frac{1}{2}$. If $\rho = \frac{1}{2}$, $\lambda_+(S)$ goes like S^{-2} near $S = 0$. We conclude that if $0 < \rho < \frac{1}{2}$, every element f in the space \mathcal{K} is also an element of $\mathcal{H}_\rho^{(+)}$ as well as of $\mathcal{H}_\rho^{(-)}$, with finite norms $\|f\|_\rho^{(\pm)}$:

$$\{\|f\|_\rho^{(\pm)}\}^2 = \frac{2^\rho}{\pi} \frac{\Gamma(\rho)}{\Gamma(1-\rho)}$$

$$\times \sum_{m=-\infty}^{\infty} \int_0^\infty dS \, S \tanh \pi S \lambda_\pm(S) |f_m(S)|^2 < \infty. \quad (3.26)$$

Consequently, for $0 < \rho < \frac{1}{2}$, the space of functions

\mathcal{K} is dense in both $\mathcal{H}_\rho^{(+)}$ and $\mathcal{H}_\rho^{(-)}$. Completion of \mathcal{K} with respect to $\|f\|_\rho^{(+)}$ yields $\mathcal{H}_\rho^{(+)}$, with respect to $\|f\|_\rho^{(-)}$ yields $\mathcal{H}_\rho^{(-)}$. If $\rho = \frac{1}{2}$, \mathcal{K} continues to be a dense subset of $\mathcal{H}_\rho^{(-)}$, but not of $\mathcal{H}_\rho^{(+)}$. [There exist elements of finite norm in \mathcal{K} which are not of finite norm in $\mathcal{H}_\rho^{(+)}$.] In any case, for all ρ in the range $0 < \rho \leq \frac{1}{2}$, $\mathcal{H}_\rho^{(+)}(\mathcal{H}_\rho^{(-)})$ consists of *all* sequences of functions $f_{+,m}(S)$ [$f_{-,m}(S)$] which have finite norms $\|f_{+,m}\|_\rho^{(+)}$ ($\|f_{-,m}\|_\rho^{(-)}$) as defined by (3.26).

The operators Q and L_3 are diagonal in the “ S, m basis,” and, denoting the transforms of $f_\pm(r, \varphi)$ according to (3.23) by $f_{\pm,m}(S)$, we have

$$[Qf]_{\pm,m}(S) = (\frac{1}{2} + S^2) f_{\pm,m}(S),$$

$$[L_3 f]_{\pm,m}(S) = m f_{\pm,m}(S). \quad (3.27)$$

In this basis, the action of the generators N_1 and N_2 is known from the structure of the UIR's of $O(2, 1)$. As for the remaining generators $N_3, L_1,$ and L_2 , we refer the reader to the discussion given in (B) for the principal series of UIR's of $O(3, 1)$. We just remark here that, for example, N_3 has for its domain the vectors f whose wavefunctions $f_{\pm,m}(S)$ are analytic functions of S ; and the vector $N_3 f$ has the following wavefunctions:

$$[N_3 f]_{\pm,m}(S)$$

$$= [m^2 - (\frac{1}{2} + iS)^2]^{\frac{1}{2}} \frac{S + i(\rho - \frac{1}{2})}{2(S - i)} f_{\mp,m}(S - i)$$

$$- [m^2 - (\frac{1}{2} - iS)^2]^{\frac{1}{2}} \frac{S + i(\rho + \frac{1}{2})}{2(S + i)} f_{\mp,m}(S + i). \quad (3.28)$$

Case II: $\frac{1}{2} < \rho < 1$.

We turn next to the remaining UIR's of $O(3, 1)$ of the supplementary series. In this case, the decomposition of the kernels K_\pm yields some important extra terms. We quote first the result derived in Appendix A, and then explain the need for these new terms. We find

$$K_+(r', r, \varphi' - \varphi)$$

$$= \frac{2^\rho}{\pi} \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \sum_{m=-\infty}^{\infty} \int_0^\infty dS \, S \tanh \pi S \lambda_+(S)$$

$$\times Y_{-\frac{1}{2}+iS}^m(r', \varphi') \overline{Y_{-\frac{1}{2}+iS}^m(r, \varphi)} + 2^{-\rho} \sqrt{\pi} \frac{\Gamma(\rho)}{\Gamma(\rho - \frac{1}{2})}$$

$$\times \sum_{m=-\infty}^{\infty} Y_{\rho-1}^m(r', \varphi') \overline{Y_{\rho-1}^m(r, \varphi)}, \quad (3.29a)$$

$$K_-(r', r, \varphi' - \varphi)$$

$$= \frac{2^\rho}{\pi} \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \sum_{m=-\infty}^{\infty} \int_0^\infty dS \, S \tanh \pi S \lambda_-(S)$$

$$\times Y_{-\frac{1}{2}+iS}^m(r', \varphi') \overline{Y_{-\frac{1}{2}+iS}^m(r, \varphi)}. \quad (3.29b)$$

The weight functions $\lambda_{\pm}(\mathcal{S})$ are the same functions of \mathcal{S} and ρ as before [cf. (3.20)]. Once again the positive definiteness of the kernels is evident.

Let us explain briefly the appearance of extra terms in (3.29a). These decompositions have been made by starting with an expansion of the function

$$(\zeta - \mu)^{-1+\rho}; \quad -1 \leq \mu \leq 1, \quad |\zeta| > 1 \quad (3.30)$$

in a series of Legendre polynomials $P_l(\mu)$ and then applying the Watson–Sommerfeld–Regge transform to the expansion. Thus this function is expressed in terms of an integral in the complex l plane, with a “background term” where the integration is along the line $\text{Re } l = -\frac{1}{2}$, together with Regge-pole type terms. Once the function has been expressed in this way, one may allow μ to vary in the range $1 \leq \mu < \infty$. Taking the limits $\zeta \rightarrow \pm 1$ and setting $\mu = \eta(r', r, \varphi' - \varphi)$ [cf. (3.16)], the decomposition of K_{\pm} is achieved. Now for small enough ρ , namely $0 < \rho \leq \frac{1}{2}$, the behavior of the function (3.30) as $\mu \rightarrow \infty$ is such that this behavior can be reproduced correctly by the background term alone. This is the result of Case I as treated above, and the background term represents precisely the continuous nonexceptional UIR’s of $O(2, 1)$. However, if ρ increases beyond $\frac{1}{2}$ and lies in the range $\frac{1}{2} < \rho < 1$, then the background term cannot account for the asymptotic behavior of (3.30) for large μ . Indeed, one finds a “Regge pole” at the point $l = \rho - 1$ in the complex l plane, and this pole contributes the extra terms in (3.29a).

The interpretation of this Regge pole is quite clear: it represents the UIR of $O(2, 1)$ of the exceptional class corresponding to the value $\sigma = \rho - \frac{1}{2}$. [Notice that the open interval $\frac{1}{2} < \rho < 1$ corresponds precisely to the range $0 < \sigma < \frac{1}{2}$, which occurs when exceptional UIR’s of $O(2, 1)$ exist.] It is also natural that this pole term contributes only to K_+ and not to K_- ; the behaviors of K_+ and K_- for large $\eta(r', r, \varphi' - \varphi)$ are different.

We express these results by means of the following:

Theorem II: UIR’s of $O(3, 1)$ belonging to the supplementary series, for values of the parameter ρ in the open interval $\frac{1}{2} < \rho < 1$, decompose into direct integrals of UIR’s of $O(2, 1)$ belonging to the continuous nonexceptional (integral) class, together with a single UIR of $O(2, 1)$ of the exceptional class, corresponding to the parameter $\sigma = \rho - \frac{1}{2}$. Every UIR of the continuous nonexceptional class appears *twice*, while the exceptional UIR appears *once*; and the spectrum of Q consists of the discrete point $Q = \frac{1}{4} - (\rho - \frac{1}{2})^2$ and the region of the real line

$\frac{1}{4} \leq Q < \infty$, the latter eigenvalues appearing twice each.

The subspace $\mathcal{K}_\rho^{(-)}$ contains exactly the same UIR’s of $O(2, 1)$ as in Case I—namely, every nonexceptional continuous class UIR appears once. As before, it follows that the space \mathcal{K} is dense in $\mathcal{K}_\rho^{(-)}$ and that \mathcal{K} yields $\mathcal{K}_\rho^{(-)}$ on completion with respect to the norm $\|f\|_\rho$ restricted to $\mathcal{K}_\rho^{(-)}$. As in (3.26), let us write this norm defined in $\mathcal{K}_\rho^{(-)}$ as $\|f\|_\rho^{(-)}$. Then $\mathcal{K}_\rho^{(-)}$ consists of *all* sequences of functions $f_{-,m}(\mathcal{S})$ for which the norm $\|f_{-}\|_\rho^{(-)}$ defined in (3.26) is finite.

The subspace $\mathcal{K}_\rho^{(+)}$ contains every nonexceptional continuous class UIR of $O(2, 1)$ once, and, in addition, contains the single exceptional UIR with $\sigma = \rho - \frac{1}{2}$. Since the exceptional UIR corresponds to a discrete point in the spectrum of Q , the corresponding vectors in $\mathcal{K}_\rho^{(+)}$ must have finite norm. The question whether or not \mathcal{K} appears as a dense subset of $\mathcal{K}_\rho^{(+)}$ is a little harder to answer. The problem is that the decomposition of K_+ as given in (3.29a) is not always valid; i.e., in order to use it to compute the norm of a vector $f_+(r, \varphi)$ in $\mathcal{K}_\rho^{(+)}$, one has to interchange the order of the integrations involved, and this may not always be justified. One can analyse the structure of $\mathcal{K}_\rho^{(+)}$ in the following way. We have proved that the space \mathcal{K} of all square-integrable functions $\tilde{f}(\theta, \varphi)$ forms a dense set in \mathcal{K}_ρ . Restricting oneself to the subset of \mathcal{K} made up of functions with the property

$$\tilde{f}(\theta, \varphi) = \tilde{f}(\pi - \theta, \varphi),$$

one is led, via (3.14) and (3.15), to a set of functions $f_+(r, \varphi)$ which form a dense subset \mathcal{D} of $\mathcal{K}_\rho^{(+)}$. That $\tilde{f}(\theta, \varphi)$ belongs to \mathcal{K} may be expressed in the following manner:

$$\int_0^\infty \sinh r \, dr \int_0^{2\pi} d\varphi [\cosh r]^{2\rho} |f_+(r, \varphi)|^2 < \infty. \quad (3.31)$$

This implies that both $f_+(r, \varphi)$ and $[\cosh r]^\rho f_+(r, \varphi)$ are elements in \mathcal{K} , so that the vectors $f_+(r, \varphi)$ in \mathcal{D} form a proper subset of \mathcal{K} .

For the vectors of $\mathcal{K}_\rho^{(+)}$ that lie in \mathcal{D} , we can prove that the decomposition (3.29a) of K_+ may be used. Since any $f_+(r, \varphi)$ in \mathcal{D} is also in \mathcal{K} and since $\lambda_+(\mathcal{S}) \rightarrow 0$ for large \mathcal{S} , it follows that

$$\sum_{m=-\infty}^\infty \int_0^\infty d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \lambda_+(\mathcal{S}) |f_{+,m}(\mathcal{S})|^2 < \infty, \quad (3.32)$$

$$f_{+,m}(\mathcal{S}) = (2\pi)^{-1} \int_0^\infty \sinh r \, dr \int_0^{2\pi} d\varphi$$

$$\times \overline{Y_{-\frac{1}{2}+i\mathcal{S}}^m(r, \varphi)} f_+(r, \varphi).$$

On the other hand, since $[\cosh r]^\rho f_+(r, \varphi)$ also belongs to \mathcal{K} , we can go further. Even though the functions

$$Y_{\rho-1}^m(r, \varphi), \quad \frac{1}{2} < \rho < 1$$

do not belong to \mathcal{K} (see Sec. 1), the functions

$$[\cosh r]^{-\rho} Y_{\rho-1}^m(r, \varphi)$$

do. Consequently, for elements $f_+(r, \varphi)$ in \mathcal{D} , using the fact that $[\cosh r]^\rho f_+(r, \varphi)$ belongs to \mathcal{K} , we have

$$\sum_{m=-\infty}^{\infty} |f_{+,m}^{(\rho-1)}|^2 < \infty,$$

$$f_{+,m}^{(\rho-1)} = (2\pi)^{-1} \int_0^\infty \sinh r \, dr \int_0^{2\pi} d\varphi \overline{Y_{\rho-1}^m(r, \varphi)} f_+(r, \varphi). \quad (3.33)$$

(3.32) and (3.33) prove that the decomposition of the kernel K_+ given in (3.29a) definitely can be used for the dense subset \mathcal{D} in $\mathcal{H}_\rho^{(+)}$. Combining (3.32), (3.33), and (3.29a), we may say the following: Every vector f_+ in the dense subset \mathcal{D} of $\mathcal{H}_\rho^{(+)}$ is specified by the following quantities:

$$f_+ \rightarrow \{f_{+,m}(\mathcal{S}), f_{+,m}^{(\rho-1)}\}. \quad (3.34)$$

$f_{+,m}(\mathcal{S})$ and $f_{+,m}^{(\rho-1)}$ are obtained via (3.32) and (3.33) from the function $f_+(r, \varphi)$ corresponding to the vector f_+ . [This function $f_+(r, \varphi)$ obeys (3.31).] The norm $\|f_+\|_\rho^{(+)}$ of f_+ is given by

$$\begin{aligned} \{\|f_+\|_\rho^{(+)}\}^2 &= \frac{2^\rho}{\pi} \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \sum_m \int_0^\infty d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \\ &\times \lambda_+(\mathcal{S}) |f_{+,m}(\mathcal{S})|^2 + 2^{-\rho} \pi^{\frac{1}{2}} \frac{\Gamma(\rho)}{\Gamma(\rho - \frac{1}{2})} \sum_m |f_{+,m}^{(\rho-1)}|^2. \end{aligned} \quad (3.35)$$

The completion of \mathcal{D} with respect to $\|f_+\|_\rho^{(+)}$ yields $\mathcal{H}_\rho^{(+)}$.

The nature of $\mathcal{H}_\rho^{(+)}$ will be revealed by the process of completion of \mathcal{D} . We must first examine the properties of the vectors f_+ in \mathcal{D} . The fact that, for an element f_+ of \mathcal{D} , the function

$$[\cosh r]^\rho f_+(r, \varphi)$$

also belongs to \mathcal{K} may be used to reach the following conclusions regarding the quantities $f_{+,m}(\mathcal{S})$ and $f_{+,m}^{(\rho-1)}$ of (3.34):

(i) $f_{+,m}(\mathcal{S})$ is the boundary value as $t \rightarrow 0$ of an analytic function $f_{+,m}(\zeta)$ of $\zeta = \mathcal{S} + it$ (boundary value in the sense of the limit in the mean); this analytic function has so singularities in the open strip $-\rho < t = \text{Im } \zeta < \rho$.

(ii) For any fixed value of $\text{Im } \zeta = t$ in this range,

$$\sum_m \int_0^\infty d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} |f_{+,m}(\mathcal{S} + it)|^2 < \infty. \quad (3.36)$$

(iii) The quantities $f_{+,m}^{(\rho-1)}$ are determined in terms of these analytic functions by

$$f_{+,m}^{(\rho-1)} = f_{+,m}(-i\sigma), \quad \sigma = \rho - \frac{1}{2}. \quad (3.37)$$

These statements follow in a straightforward way from an examination of the behavior of the functions $Y_{\rho-\frac{1}{2}+i\mathcal{S}}^m(r, \varphi)$ for large values of r , and from the fact that these are entire functions of \mathcal{S} .

Thus if one computes the norm $\|f_+\|_\rho^{(+)}$ of a vector f_+ in \mathcal{D} according to (3.35), one finds that the two parts $f_{+,m}(\mathcal{S})$ and $f_{+,m}^{(\rho-1)}$ of f_+ are not independent; in particular, because of (3.37), there can be no vector in \mathcal{D} for which $f_{+,m}(\mathcal{S})$ vanishes while $f_{+,m}^{(\rho-1)}$ does not. Nevertheless, on completion of \mathcal{D} with respect to $\|f_+\|_\rho^{(+)}$, one finds the following structure for $\mathcal{H}_\rho^{(+)}$ ¹²: Every vector f_+ in $\mathcal{H}_\rho^{(+)}$ consists of a sequence of functions $f_{+,m}(\mathcal{S})$ and a sequence of complex numbers $f_{+,m}^{(\rho-1)}$, subject only to the condition that the norm $\|f_+\|_\rho^{(+)}$ be finite— $\|f_+\|_\rho^{(+)}$ defined in (3.35). For a general vector f_+ in $\mathcal{H}_\rho^{(+)}$, there is no relation between these functions of \mathcal{S} and the complex numbers $f_{+,m}^{(\rho-1)}$; to “recover” all of $\mathcal{H}_\rho^{(+)}$, one must choose these independently of one another. This means that $\mathcal{H}_\rho^{(+)}$ can itself be split up into a direct sum of two Hilbert spaces $\mathcal{H}_{\rho,c}^{(+)}$ and $\mathcal{H}_{\rho,d}^{(+)}$:

$$\mathcal{H}_\rho^{(+)} = \mathcal{H}_{\rho,c}^{(+)} \oplus \mathcal{H}_{\rho,d}^{(+)}. \quad (3.38)$$

The subscripts c and d refer respectively to the continuous and the discrete eigenvalues of Q . Vectors f_+ in $\mathcal{H}_{\rho,c}^{(+)}$ correspond to all sequences of functions $f_{+,m}(\mathcal{S})$ with

$$\{\|f_+\|_{\rho,c}^{(+)}\}^2 = \sum_m \int_0^\infty d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \lambda_+(\mathcal{S}) |f_{+,m}(\mathcal{S})|^2 < \infty, \quad (3.39)$$

while vectors f_+ in $\mathcal{H}_{\rho,d}^{(+)}$ correspond to all sequences of complex numbers $f_{+,m}^{(\rho-1)}$ with

$$\{\|f_+\|_{\rho,d}^{(+)}\}^2 = \sum_m |f_{+,m}^{(\rho-1)}|^2 < \infty. \quad (3.40)$$

This structure for $\mathcal{H}_\rho^{(+)}$ is proved by showing that, on the one hand, if one assumes $\mathcal{H}_\rho^{(+)}$ to be characterized in this way, then \mathcal{D} is dense in $\mathcal{H}_\rho^{(+)}$; and, on the other hand, we know that such a $\mathcal{H}_\rho^{(+)}$ is in fact a Hilbert space. [The argument is outlined in Appendix B.]

The question whether and, if so, in what way \mathcal{K} is contained in $\mathcal{H}_\rho^{(+)}$ is now answered as follows: \mathcal{K} is contained as a dense subset in $\mathcal{H}_{\rho,c}^{(+)}$ only. This follows from the expression (3.39) for the norm in $\mathcal{H}_{\rho,c}^{(+)}$.

The subspace $\mathcal{H}_\rho^{(-)}$ of \mathcal{H}_ρ has a more or less uniform behavior for all ρ in the range $0 < \rho < 1$. In all

¹² The author is indebted to Dr. N. J. Papastamatiou and Professor L. O’Raifeartaigh for help in clarifying this situation.

cases $\mathcal{H}_\rho^{(-)}$ contains \mathcal{K} as a dense subset and carries only the continuous nonexceptional UIR's of $O(2, 1)$. On the other hand, for $0 < \rho < \frac{1}{2}$, $\mathcal{H}_\rho^{(+)}$ carries only the continuous nonexceptional UIR's of $O(2, 1)$ and contains \mathcal{K} as a dense subset. For $\frac{1}{2} < \rho < 1$, $\mathcal{H}_\rho^{(+)}$ breaks up into two orthogonal subspaces $\mathcal{H}_{\rho,c}^{(+)}$ and $\mathcal{H}_{\rho,d}^{(+)}$. The former again carries all the continuous nonexceptional UIR's of $O(2, 1)$, and contains \mathcal{K} in a natural way as a dense subset. $\mathcal{H}_{\rho,d}^{(+)}$ carries the single exceptional UIR of $O(2, 1)$ corresponding to the value $\sigma = \rho - \frac{1}{2}$. On $\mathcal{H}_{\rho,d}^{(+)}$, Q reduces to the number $\frac{1}{4} - (\rho - \frac{1}{2})^2$. For $\rho = \frac{1}{2}$, $\mathcal{H}_\rho^{(+)}$ carries the continuous nonexceptional UIR's of $O(2, 1)$ only, but does not contain \mathcal{K} as a dense subset in any natural way.

It would be desirable to check by direct calculation whether the wavefunctions $Y_{\rho-1}^m(r, \varphi)$ and $Y_{-\frac{1}{2}+\rho-1}^m(r, \varphi)$ are orthonormal with respect to the kernel K_+ of (3.16). If this could be done, then the exceptional UIR's of $O(2, 1)$ would appear somewhat less strange. Even though orthonormality of these functions does not obtain in the ordinary sense, as explained in Sec. 1, we might have orthonormality with respect to an invariant positive definite "two-point measure" on $O(2, 1)$. If indeed this is the case, we would then be able to discuss in these UIR's also the action of the generators L_1, L_2, N_3 in the basis where Q and L_3 are diagonal.

CONCLUSION

We have discussed the reduction of the exceptional class of unitary representations of $O(2, 1)$, and of the supplementary series of unitary representations of $O(3, 1)$, under the noncompact subgroups $O(1, 1)$ and $O(2, 1)$, respectively. In the former case, we recovered the result that the spectrum of the generator of $O(1, 1)$ covers the real line twice. In the latter case, we found that the reduction gives different results depending on whether the parameter ρ labeling these representations obeys $0 < \rho \leq \frac{1}{2}$ or $\frac{1}{2} < \rho < 1$. In the first case, the only representations of $O(2, 1)$ that are present are those of the continuous nonexceptional class, and each of these occurs twice. In the second case, in addition to this double occurrence of the continuous nonexceptional representations, we have a single exceptional unitary representation of $O(2, 1)$. The value of the parameter σ describing this exceptional representation is related to ρ by $\sigma = \rho - \frac{1}{2}$.

In (B) we have shown that the principal series of UIR's of $O(3, 1)$ contains only continuous non-exceptional and discrete class UIR's of $O(2, 1)$. It is gratifying that the exceptional UIR's of $O(2, 1)$ have been found in part of the supplementary series of UIR's of $O(3, 1)$, since, on general grounds, one

knows that every UIR of $O(2, 1)$ must appear in some UIR of $O(3, 1)$.

It is also interesting that, in the problem considered here, one is able to associate a unitary representation of $O(2, 1)$ with a Regge pole in the complex l plane in an unambiguous and concrete fashion. This pole moves with changing ρ , and corresponds to a unitary representation only when $\frac{1}{2} < \rho < 1$.

Lastly, we should draw attention to the problems mentioned at the end of Sec. 3. It seems likely that the complicated expressions for the scalar products in the relevant Hilbert spaces are necessary if we want a standard algebraic form for the generators in all classes of UIR's. [Compare, for instance, (2.13) with the similar equation in (A) for the generators of $O(2, 1)$ in the continuous and discrete representations.] It should be possible to write the scalar product in a simple way, but with different-looking differential operators for the generators in each class of UIR's. We hope to treat these questions elsewhere.

ACKNOWLEDGEMENTS

The author thanks the members of the Theoretical High Energy Group at Syracuse University, as well as the members of the Irving Institute, for many useful discussions.

APPENDIX A

We outline here the derivation of Eqs. (3.7), (3.20), and (3.29). Consider the function $L(\mu)$ defined by

$$L(\mu) = (\zeta - \mu)^{-1+\rho}; \quad -1 \leq \mu \leq 1, \quad |\zeta| > 1. \quad (A1)$$

We can expand $L(\mu)$ in a series of Legendre polynomials for μ in the indicated range:

$$L(\mu) = \frac{1}{2} \sum_{l=0}^{\infty} (2l+1) a_l(\zeta, \rho) P_l(\mu), \quad (A2)$$

$$a_l(\zeta, \rho) = \int_{-1}^1 d\mu P_l(\mu) (\zeta - \mu)^{-1+\rho}.$$

Using the standard property of Legendre functions of the second kind,¹³ that

$$Q_l(\mu + i\epsilon) - Q_l(\mu - i\epsilon) = -i\pi P_l(\mu), \quad -1 \leq \mu \leq 1, \quad (A3)$$

we can express the coefficients a_l in terms of a contour integral in the complex l plane:

$$a_l(\zeta, \rho) = \frac{i}{\pi} \oint_{C'} d\mu Q_l(\mu) (\zeta - \mu)^{-1+\rho}. \quad (A4)$$

¹³ This is derived easily from Neumann's formula, Ref. 9, p. 154.

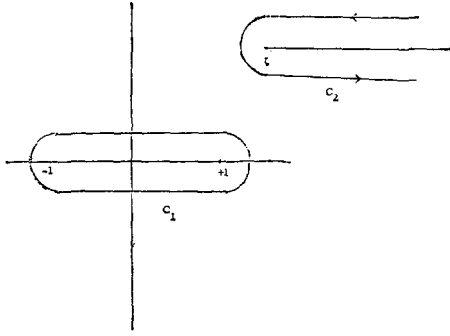


FIG. 1

The contour C_1 is given in Fig. 1; one can now deform the contour C_1 into the contour C_2 because the behavior of $Q_i(\hbar)$ for large $|\hbar|$ permits the neglect of the contribution from infinity, and because there are no other singularities of the integrand except a branch point at $\hbar = \zeta$:

$$a_i(\zeta, \rho) = \frac{i}{\pi} \oint_{C_1} d\hbar Q_i(\hbar) (\zeta - \hbar)^{-1+\rho}. \quad (\text{A5})$$

On the contour C_2 , $\hbar - \zeta$ is a real positive number. By keeping track of the phase of the integrand while deforming the contour C_1 into C_2 , we find

$$a_i(\zeta, \rho) = \frac{2}{\pi} \sin \pi \rho \int_{\zeta}^{\infty} d\hbar Q_i(\hbar) (\hbar - \zeta)^{\rho-1}. \quad (\text{A6})$$

We are interested in the range $0 < \rho < 1$ for ρ . For this set of values of ρ , the integral appearing in (A6) can be evaluated explicitly,¹⁴ and we finally get

$$a_i(\zeta, \rho) = \frac{2}{\pi} \sin \pi \rho (\zeta^2 - 1)^{\rho/2} e^{i\pi\rho} \Gamma(\rho) Q_i^{-\rho}(\zeta). \quad (\text{A7})$$

We can now write (A2) in the compact form

$$\begin{aligned} & (\zeta - \mu)^{-1+\rho} \\ &= \frac{e^{i\pi\rho}}{\Gamma(1-\rho)} (\zeta^2 - 1)^{\rho/2} \sum_{l=0}^{\infty} (2l+1) P_l(\mu) Q_l^{-\rho}(\zeta), \\ & \quad 0 < \rho < 1; \quad -1 < \mu < 1, \quad |\zeta| > 1. \end{aligned} \quad (\text{A8})$$

For $\rho \rightarrow 0$, we recover the well-known formula due to Heine.

(A8) leads directly to Eq. (3.7) of the text. We first take the limit $\zeta \rightarrow 1$ through values greater than unity, and then identify μ to be

$$\mu = \cos \theta' \cos \theta + \sin \theta' \sin \theta \cos(\varphi' - \varphi). \quad (\text{A9})$$

This gives, using the notation of (3.5),

$$\begin{aligned} & L_{\rho}(\theta' \varphi', \theta \varphi) \\ &= 2^{\rho-1} \frac{\Gamma(\rho)}{\Gamma(1-\rho)} \sum_{l=0}^{\infty} (2l+1) \frac{\Gamma(l-\rho+1)}{\Gamma(l+\rho+1)} \\ & \times P_l(\cos \theta' \cos \theta + \sin \theta' \sin \theta \cos(\varphi' - \varphi)). \end{aligned} \quad (\text{A10})$$

¹⁴ Reference 9, p. 160, Eq. (31).

An application of the spherical-harmonic addition theorem¹⁵ then gives (3.7).

The partial wave expansion (A8) is not appropriate for large values of μ , and therefore we apply the usual Watson-Sommerfeld-Regge transform to it. We first write it in the form of a contour integral in the complex l plane:

$$\begin{aligned} (\zeta - \mu)^{-1+\rho} &= \frac{e^{i\pi\rho}}{\Gamma(1-\rho)} (\zeta^2 - 1)^{\rho/2} \frac{1}{2i} \oint_C \frac{(2l+1) dl}{\sin \pi l} \\ & \times P_l(-\mu) Q_l^{-\rho}(\zeta). \end{aligned} \quad (\text{A11})$$

The contour C encircles the positive real l axis, running from $+\infty$ to 0 above it, and from 0 to $+\infty$ below it. The poles of the function $Q_l^{-\rho}(\zeta)$ lie at the points¹⁶

$$l - \rho = -1, -2, -3, \dots, \quad (\text{A12})$$

and so, as long as ρ is in the open interval $0 < \rho < 1$, we can always make the contour C go around the origin $l = 0$ without going through these poles. If we now try to deform the contour C into the straight line $l = -\frac{1}{2} + i\mathcal{S}$, parallel to the imaginary axis, we see from (A12) that we will encounter no poles of the integrand in doing so, for $0 < \rho \leq \frac{1}{2}$. Even for $\rho = \frac{1}{2}$, the (simple) pole of $Q_l^{-\frac{1}{2}}$ at $l = -\frac{1}{2}$ is cancelled by the factor $(2l+1)$ in (A11). Therefore we get

$$\begin{aligned} 0 < \rho \leq \frac{1}{2}: \quad & (\zeta - \mu)^{-1+\rho} \\ &= \frac{e^{i\pi\rho}}{\Gamma(1-\rho)} (\zeta^2 - 1)^{\rho/2} \frac{1}{2i} \int_{-\frac{1}{2}+i\infty}^{-\frac{1}{2}-i\infty} \frac{(2l+1) dl}{\sin \pi l} \\ & \times P_l(-\mu) Q_l^{-\rho}(\zeta). \end{aligned} \quad (\text{A13})$$

It is now permissible to replace μ by $-\mu$ and to allow μ to vary in the range $1 \leq \mu < \infty$, since the right-hand side has the correct asymptotic behavior in μ . We replace l by $-\frac{1}{2} + i\mathcal{S}$ and use standard properties of the $Q_l^{-\rho}$ functions in order to write (A13) as¹⁷

$$\begin{aligned} (\zeta + \mu)^{-1+\rho} &= \frac{(\zeta^2 - 1)^{\rho/2}}{\Gamma(1-\rho)} \int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \\ & \times \Gamma(\frac{1}{2} - \rho + i\mathcal{S}) \Gamma(\frac{1}{2} - \rho - i\mathcal{S}) \\ & \times P_{-\frac{1}{2}+i\mathcal{S}}(\mu) P_{-\frac{1}{2}+i\mathcal{S}}(\zeta). \end{aligned} \quad (\text{A14})$$

All that remains to be done is to take the limits $\zeta \rightarrow \pm 1$. We find

$$\begin{aligned} & (\zeta^2 - 1)^{\rho/2} P_{-\frac{1}{2}+i\mathcal{S}}(\zeta) \xrightarrow{\zeta \rightarrow 1} 2^{\rho} \frac{\Gamma(\rho)}{\pi} \sin \pi \rho \\ & \xrightarrow{\zeta \rightarrow -1} 2^{\rho} \frac{\Gamma(\rho)}{\pi} \cos \pi \mathcal{S}. \end{aligned} \quad (\text{A15})$$

¹⁵ Reference 9, p. 168.

¹⁶ L. Robin, *Fonctions sphériques de Legendre et fonctions sphéroidales* (Gauthier-Villars, Paris, 1958), Tome II, p. 88.

¹⁷ The formula (A14), valid for $0 < \rho \leq \frac{1}{2}$, may be found also in L. Robin, Ref. 16, Tome III, p. 167.

Identifying μ with $\eta(r', r, \varphi' - \varphi)$ of (3.16),

$$\begin{aligned} \mu &= \cosh r' \cosh r - \sinh r' \sinh r \cos(\varphi' - \varphi) \\ &= \eta(r', r, \varphi' - \varphi). \end{aligned} \tag{A16}$$

Combining (A14) with (A15), we can expand the kernels $K_{\pm}(r', r, \varphi' - \varphi)$ of (3.16) in the following fashion:

$$\begin{aligned} K_{\pm}(r', r, \varphi' - \varphi) &= \frac{2^{\rho}}{\pi} \frac{\Gamma(\rho)}{\Gamma(1 - \rho)} \int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \lambda_{\pm}(\mathcal{S}) P_{-\frac{1}{2} + i\mathcal{S}} \\ &\quad \times [\eta(r', r, \varphi' - \varphi)]. \end{aligned} \tag{A17}$$

The functions $\lambda_{\pm}(\mathcal{S})$ have been given in Eq. (3.20). Lastly, we use the addition theorem for the Legendre function appearing in (A17),¹⁵ and this leads to the decompositions of the kernels K_{\pm} given in (3.20).

In the case $\frac{1}{2} < \rho < 1$, all the above manipulations go through except for the fact that the integrand in (A11) has a ‘‘Regge pole’’ at $l = \rho - 1$; and this contributes an extra piece to the right-hand side of (A13) since now $-\frac{1}{2} < \rho - 1 < 0$. Therefore we have, in place of (A13),

$$\begin{aligned} \frac{1}{2} < \rho < 1: \\ (\zeta - \mu)^{-1 + \rho} &= \frac{e^{i\pi\rho}}{\Gamma(1 - \rho)} \frac{(\zeta^2 - 1)^{\rho/2}}{2i} \int_{-\frac{1}{2} + i\infty}^{-\frac{1}{2} - i\infty} \\ &\quad \times \frac{(2l + 1) dl}{\sin \pi l} P_l(-\mu) Q_l^{-\rho}(\zeta) \\ &\quad + 2^{1-\rho} \sqrt{\pi} \frac{\Gamma(\rho)}{\Gamma(\rho - \frac{1}{2})} P_{\rho-1}(-\mu). \end{aligned} \tag{A18}$$

Notice the interesting fact that the pole-term has no dependence on ζ . Also it has just the right behavior as $\mu \rightarrow \infty$ to agree with the left-hand side. As before, we replace $\mu \rightarrow -\mu$, identify μ with $\eta(r', r, \varphi' - \varphi)$, and take the limits $\zeta \rightarrow \pm 1$. Since the pole term is independent of ζ , it contributes only to the kernel K_+ . Use of the addition theorem for the Legendre functions¹⁵ then yields (3.29).

APPENDIX B

We outline the arguments that lead to the form for the Hilbert space $\mathcal{K}_{\rho}^{(+)}$ derived from its dense subset \mathcal{D} in Sec 3, Case II. The presence of the variable φ is irrelevant, and so we will omit it completely. The superscripts and subscripts (+) can be omitted also.

Let \mathcal{K} be the Hilbert space of function $\tilde{f}(x)$ which have finite norm in the following sense:

$$\|f\|_{\mathcal{K}}^2 = \int_1^{\infty} |\tilde{f}(x)|^2 dx < \infty. \tag{B1}$$

Such functions $\tilde{f}(x)$ may always be expressed in the form

$$\begin{aligned} \tilde{f}(x) &= \int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} P_{-\frac{1}{2} + i\mathcal{S}}(x) f(\mathcal{S}), \\ f(\mathcal{S}) &= \int_1^{\infty} dx P_{-\frac{1}{2} + i\mathcal{S}}(x) \tilde{f}(x). \end{aligned} \tag{B2}$$

Then,

$$\|f\|_{\mathcal{K}}^2 = \int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} |f(\mathcal{S})|^2 \equiv \|f(\mathcal{S})\|_{\mathcal{K}}^2 < \infty, \tag{B3}$$

so that we can equally well consider \mathcal{K} as the space of all functions $f(\mathcal{S})$ with finite norms as computed according to (B3).

Among all elements $\tilde{f}(x)$ in \mathcal{K} , those obeying the relation

$$\int_1^{\infty} x^{2\rho} |\tilde{f}(x)|^2 dx < \infty, \tag{B4}$$

where ρ is some positive real number, certainly constitute a dense set $\mathcal{D}^{(0)}$ in \mathcal{K} . (We are interested only in $\frac{1}{2} < \rho < 1$.) How can the functions $f(\mathcal{S})$ corresponding to $\tilde{f}(x) \in \mathcal{D}^{(0)}$ be characterized? The behavior of $P_{-\frac{1}{2} + i\mathcal{S}}(x)$ for large x shows quite easily that if $f(\mathcal{S})$ belongs to $\mathcal{D}^{(0)}$, then:

(i) $f(\mathcal{S})$ is the boundary value (in the l. i. m. sense) of an analytic function $f(\zeta)$ of $\zeta = \mathcal{S} + it$, as $t \rightarrow 0$; $f(\zeta)$ is free of singularities in the open strip $-\rho < t < \rho$.

(ii) For any fixed value of $t = \text{Im } \zeta$ in this range

$$\int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} |f(\mathcal{S} + it)|^2 < \infty. \tag{B5}$$

In general, if $f(\mathcal{S})$ belongs to $\mathcal{D}^{(0)}$, then $f[-i(\rho - \frac{1}{2})]$ may not vanish. [Note that, for $\rho > \frac{1}{2}$, the point $\mathcal{S} = -i(\rho - \frac{1}{2})$ is contained within the strip of analyticity of $f(\mathcal{S})$.] Let us define a subset $\mathcal{D}^{(1)}$ in $\mathcal{D}^{(0)}$ as follows: A vector $f(\mathcal{S})$ belongs to $\mathcal{D}^{(1)}$ if (i) $f(\mathcal{S})$ belongs to $\mathcal{D}^{(0)}$, and (ii) $f[-i(\rho - \frac{1}{2})] = 0$. We first show that $\mathcal{D}^{(1)}$ is also dense in \mathcal{K} . If $f(\mathcal{S})$ belongs to $\mathcal{D}^{(0)}$, the function

$$f'(\mathcal{S}) = \frac{\mathcal{S} + i(\rho - \frac{1}{2})}{\mathcal{S}^2 + a^2} f(\mathcal{S}), \quad a > \rho \tag{B6}$$

clearly obeys all conditions to belong to $\mathcal{D}^{(1)}$. If a vector $h(\mathcal{S})$ is orthogonal to all elements of $\mathcal{D}^{(1)}$, it must, in particular, obey the condition that

$$\int_0^{\infty} d\mathcal{S} \mathcal{S} \tanh \pi \mathcal{S} \left[h(\mathcal{S}) \frac{\mathcal{S} - i(\rho - \frac{1}{2})}{\mathcal{S}^2 + a^2} \right]^* f(\mathcal{S}) = 0 \tag{B7}$$

for all $f(\mathcal{S})$ in $\mathcal{D}^{(0)}$. Together with $h(\mathcal{S})$, the bracketed

expression in (B7) is certainly an element of \mathcal{K} and, since $\mathcal{D}^{(0)}$ is dense in \mathcal{K} , it must vanish. It follows that $h(\mathcal{S})$ must itself vanish and that $\mathcal{D}^{(1)}$ is dense in \mathcal{K} .

Let us now define another Hilbert space \mathcal{K} as the direct sum of a one-dimensional space \mathcal{K}_a and an infinite-dimensional space \mathcal{K}_c :

$$\mathcal{K} = \mathcal{K}_c \oplus \mathcal{K}_a. \tag{B8}$$

\mathcal{K}_c consists of all functions $f(\mathcal{S})$ for which

$$\|f\|_c^2 = \int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} \lambda(\mathcal{S}) |f(\mathcal{S})|^2 < \infty. \tag{B9}$$

Then every vector f in \mathcal{K} corresponds to a pair

$$f \rightarrow \{f(\mathcal{S}), f_1\} \tag{B10}$$

made up of a function $f(\mathcal{S})$ in \mathcal{K}_c , and a complex number f_1 , with the norm in \mathcal{K} given by

$$\|f\|^2 = \int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} \lambda(\mathcal{S}) |f(\mathcal{S})|^2 + |f_1|^2. \tag{B11}$$

The weight function $\lambda(\mathcal{S})$ in (B9) and (B11) is the real positive-definite function $\lambda_+(\mathcal{S})$ of (3.20). We will use the fact that, as $\mathcal{S} \rightarrow \infty$, $\lambda(\mathcal{S})$ goes to zero like $\mathcal{S}^{-2\rho}$.

We define a subset \mathcal{D} in \mathcal{K} as follows; a vector f lies in \mathcal{D} if

- (i) $f(\mathcal{S}) \in \mathcal{D}^{(0)}$,
- (ii) $f_1 = f[-i(\rho - \frac{1}{2})]$. (B12)

Since $f(\mathcal{S})$ belongs to \mathcal{K} if it belongs to $\mathcal{D}^{(0)}$, and since $\lambda(\mathcal{S})$ vanishes as $\mathcal{S} \rightarrow \infty$, the first term in (B11) will

be finite for such an $f(\mathcal{S})$. We will now show that \mathcal{D} is dense in \mathcal{K} . We assume the existence of a vector $h \in \mathcal{K}$ orthogonal to all elements f of \mathcal{D} . It obeys

$$\int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} \lambda(\mathcal{S}) \overline{h(\mathcal{S})} f(\mathcal{S}) + h_1 f[-i(\rho - \frac{1}{2})] = 0;$$

$$h = \{h(\mathcal{S}), h_1\}, f = \{f(\mathcal{S}), f[-i(\rho - \frac{1}{2})]\} \in \mathcal{D}. \tag{B13}$$

Restricting $f(\mathcal{S})$ to the dense set $\mathcal{D}^{(1)}$ in \mathcal{K} , we have

$$\int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} \lambda(\mathcal{S}) \overline{h(\mathcal{S})} f(\mathcal{S}) = 0, \quad f(\mathcal{S}) \in \mathcal{D}^{(1)}. \tag{B14}$$

Since $h(\mathcal{S})$ belongs to \mathcal{K} , the integral

$$\int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} \lambda(\mathcal{S}) |h(\mathcal{S})|^2$$

is finite; so then also is

$$\int_0^\infty d\mathcal{S} \tanh \pi\mathcal{S} |\lambda(\mathcal{S}) h(\mathcal{S})|^2.$$

Hence $\lambda(\mathcal{S}) h(\mathcal{S})$ belongs to \mathcal{K} . Since $\mathcal{D}^{(1)}$ is dense in \mathcal{K} and $\lambda(\mathcal{S})$ is nonvanishing, then $h(\mathcal{S})$ must vanish identically. We now go back to (B13) and pick a vector f in \mathcal{D} for which f_1 is nonvanishing. Then we conclude $h_1 = 0$, or that the vector h vanishes identically. Thus the subset \mathcal{D} in \mathcal{K} is dense in \mathcal{K} , and completion of \mathcal{D} with respect to $\|f\|$ of (B11) yields \mathcal{K} .

Following essentially these same arguments—but with the variable φ included—we reach the conclusions stated in Sec. 3 of the text.

Asymptotic Properties of Perturbation Theory

JOSEPH B. KRIEGER

Polytechnic Institute of Brooklyn, Brooklyn, New York

(Received 27 October 1966)

The perturbation expansions are derived by a technique which does not assume that convergent expansions exist. The theory is shown to be asymptotic, and criteria are developed to determine if a finite number of terms underestimates or overestimates the exact result for sufficiently small values of the coupling constant.

I. INTRODUCTION

Conventional derivations in time-independent perturbation theory of the eigenfunctions and energy eigenvalues¹ generally proceed by assuming that both have a power series expansion in the coupling constant and then finding the coefficients in the power series by setting the coefficients of arbitrary powers of the coupling constant equal to zero in Schrödinger's equation. This treatment is rigorous provided that the perturbation series converges. However, if the series has no radius of convergence, from a purely mathematical point of view the derivation is meaningless. Nevertheless, perturbation theory is constantly used with no attempt made to show that the series actually converges. Furthermore, in those instances where it has been shown that the series diverges for all nonzero values of the coupling constant (i.e., quantum

electrodynamics), the series is said to be asymptotic, which thus allows us to feel assured that taking only the first few terms gives a reasonably close approximation to the exact result. However, this conjecture does not follow directly from the mathematical derivation mentioned above, but rather it arises intuitively from the notion that for a small enough coupling constant, the solution should approach the unperturbed solution. The divergence of the expansion should not be thought of as due to some subtlety involved in the interaction, as one might believe for quantum electrodynamics, but actually arises for a large class of simple single-particle potentials. For example, if we perturb a one-dimensional harmonic oscillator with an anharmonic term αx^4 , the perturbation series for the ground state energy diverges.² This follows from the fact that a typical term in the n th order

$$\alpha^n \frac{\langle 0 | x^4 | 4 \rangle \langle 4 | x^4 | 8 \rangle \cdots \langle 4n | x^4 | 4n \rangle \langle 4n | x^4 | 4n - 4 \rangle \cdots \langle 4 | x^4 | 0 \rangle}{(\epsilon_0 - \epsilon_4)(\epsilon_0 - \epsilon_8) \cdots (\epsilon_0 - \epsilon_{4n})(\epsilon_0 - \epsilon_{4n-4}) \cdots (\epsilon_0 - \epsilon_4)} \sim \alpha^n \frac{[(2n)!]^4}{(n!)^2} \sim \alpha^n (n!)^2 \xrightarrow{n \rightarrow \infty} \infty$$

for any nonzero α .

Without even doing this calculation, it is clear that the series must diverge by the following line of reasoning. Suppose the series converged for $0 \leq \alpha \leq \bar{\alpha}$. Then it follows from the theory of functions that the series will also converge for $-\bar{\alpha} \leq \alpha \leq 0$. But this is impossible since there are no solutions for the problem for negative α . Hence the series must diverge for positive α . Obviously the same reasoning can be applied to any potential, and we reach the conclusion that if by changing the sign of the coupling constant there are no bound-state solutions to the Schrödinger equation, then the perturbation series has no radius of convergence for either sign of the coupling constant. Since so many potentials have this property, it is worthwhile to investigate perturbation

theory by an alternate derivation which sheds some light on its asymptotic nature.

We show that perturbation theory does yield asymptotic expansions. Furthermore, we show that for a sufficiently small coupling constant the expansion for the energy always improves in accuracy, if successive terms have the same sign; and if they have different signs, the exact result lies between the two partial sums. Finally, it is observed that the partial sums up through an odd order of perturbation theory do not necessarily overestimate the energy, as has been stated elsewhere.

II. BASIC THEOREM

Consider

$$H(\lambda)\psi_i(\lambda) = \epsilon_i(\lambda)\psi_i(\lambda), \tag{1}$$

¹ E. C. Kemble, *The Fundamental Principles of Quantum Mechanics* (McGraw-Hill Book Company, Inc., New York, 1937), pp. 380-388.

² H. A. Kramers, *Quantum Mechanics* (Dover Publications, Inc., New York, 1964), p. 194.

where

$$H(\lambda) = H_0 + \lambda H'.$$

Then

$$\left. \frac{\partial \epsilon_i^{p+1}(\lambda)}{\partial \lambda^{p+1}} \right|_{\lambda=0} \quad \text{and} \quad \left. \frac{\partial^p \psi_i(\lambda)}{\partial \lambda^p} \right|_{\lambda=0}$$

are given by the results of perturbation theory, provided that all terms up to and including the $(p + 1)$ th terms in the perturbation expansion converge, even though the perturbation series expansion may ultimately diverge for $\epsilon_i(\lambda)$ and $\psi_i(\lambda)$.

Proof: Let ϕ_n be the complete orthonormal set of functions satisfying

$$H_0 \phi_n = \epsilon_n^0 \phi_n. \tag{2}$$

Then, from the completeness of the ϕ_n , we can write

$$\psi_i(\lambda) = \sum_{n'} b_{n'}^{(i)}(\lambda) \phi_{n'}, \tag{3}$$

and since as $\lambda \rightarrow 0$, $\psi_i(\lambda) \rightarrow \phi_i$, then

$$b_n^{(i)}(\lambda) \Big|_{\lambda \rightarrow 0} = \delta_{in}. \tag{4}$$

From here on the superscript i will be understood.

Substituting Eq. (3) into Eq. (1), multiplying both sides of the resulting equation by ϕ_n^* , and integrating over all particle coordinates, we obtain

$$\lambda \sum_{n'} H'_{nn'} b_{n'}(\lambda) = [\epsilon_i(\lambda) - \epsilon_n^0] b_n(\lambda), \tag{5}$$

where

$$H'_{nn'} \equiv \langle \phi_n | H' | \phi_{n'} \rangle.$$

Equation (5) is just the usual set of coupled equations for the $b_n(\lambda)$ that are conventionally solved by iteration to obtain the perturbation-theory series expansions.

Equation (5) is sufficient to determine $\epsilon_i(\lambda)$ and the $\{b_n(\lambda)\}$, but it is more convenient to make use of Feynman's theorem,³ which in terms of the $\{b_n\}$ can be written

$$\sum_{nn'} b_n^*(\lambda) b_{n'}(\lambda) H'_{nn'} = \frac{\partial \epsilon_i(\lambda)}{\partial \lambda}. \tag{6}$$

Then from Eq. (5), using (4), we have

$$\epsilon_i^0 = \epsilon_i(\lambda = 0), \tag{7}$$

provided only that

$$H'_{ii}$$

exists, and from Eq. (6),

$$H'_{ii} = \left. \frac{\partial \epsilon_i(\lambda)}{\partial \lambda} \right|_{\lambda=0}, \tag{8}$$

which is just the usual perturbation-theory result.

To obtain results to the next order we differentiate Eqs. (5) and (6) with respect to λ and obtain

$$\sum_{n'} H'_{nn'} b_{n'}(\lambda) + \lambda \sum_{n'} H'_{nn'} \frac{\partial b_{n'}(\lambda)}{\partial \lambda} = \frac{\partial \epsilon_i(\lambda)}{\partial \lambda} b_n(\lambda) + [\epsilon_i(\lambda) - \epsilon_n^0] \frac{\partial b_n(\lambda)}{\partial \lambda} \tag{5'}$$

and

$$\sum_{n,n'} \left[\frac{\partial b_n^*(\lambda)}{\partial \lambda} b_{n'}(\lambda) + b_n^*(\lambda) \frac{\partial b_{n'}(\lambda)}{\partial \lambda} \right] H'_{nn'} = \frac{\partial^2 \epsilon_i(\lambda)}{\partial \lambda^2}. \tag{6'}$$

Then, letting $\lambda \rightarrow 0$ in Eq. (5') and using Eqs. (4) and (7), we have

$$H'_{ni} = (\epsilon_i^0 - \epsilon_n^0) \left. \frac{\partial b_n}{\partial \lambda} \right|_0 \quad n \neq i,$$

provided only that

$$\sum_{n'} H'_{nn'} \frac{\partial b_{n'}}{\partial \lambda}$$

exists (at this point it is necessary to assume that if ϵ_i was originally degenerate, then the unperturbed wave functions have been chosen so that $H'_{in} = 0$ for all $n \neq i$), i.e., provided

$$\sum_{n'} \frac{H'_{nn'} H'_{n'i}}{\epsilon_i^0 - \epsilon_n^0}$$

exists, where the prime on the sum means $i = n'$ is omitted. (*Note:* Since

$$|b_i(\lambda)|^2 = 1 - \sum_{n \neq i} |b_n(\lambda)|^2,$$

then differentiating with respect to λ and using $b_n(\lambda) \rightarrow 0$, we obtain $(\partial b_i / \partial \lambda)|_0 = 0$.) But this is just the second-order correction to b_n in perturbation theory. Hence if the second-order correction to b_n exists in perturbation theory, then

$$\left. \frac{\partial b_n}{\partial \lambda} \right|_{\lambda=0} = \frac{H'_{ni}}{\epsilon_i^0 - \epsilon_n^0} \quad n \neq i, \tag{9}$$

which is, of course, just the usual perturbation-theory result. Furthermore, from Eqs. (6') and (9),

$$\sum_{n'} \frac{H'_{in'} H'_{n'i}}{\epsilon_i - \epsilon_{n'}} = \frac{1}{2} \left. \frac{\partial^2 \epsilon_i(\lambda)}{\partial \lambda^2} \right|_{\lambda=0}, \tag{10}$$

provided only that the sum exists.

Higher-order terms may be obtained by differentiating (5') with respect to λ . Then

$$2 \sum_{n'} H'_{nn'} \frac{\partial b_{n'}}{\partial \lambda} + \lambda \sum_{n'} H'_{nn'} \frac{\partial^2 b_{n'}}{\partial \lambda^2} = \frac{\partial^2 \epsilon_i(\lambda)}{\partial \lambda^2} b_n + 2 \frac{\partial \epsilon_i(\lambda)}{\partial \lambda} \frac{\partial b_n}{\partial \lambda} + [\epsilon_i(\lambda) - \epsilon_n^0] \frac{\partial^2 b_n}{\partial \lambda^2}. \tag{5''}$$

³ W. Pauli, *Encyclopedia of Physics*, S. Flügge, Ed. (Springer-Verlag, Berlin, 1958), Vol. V, Part I, p. 83.

Letting $\lambda \rightarrow 0$ and using Eqs. (4), (7), (8), and (10), we obtain

$$\frac{1}{2} \frac{\partial^2 b_n}{\partial \lambda^2} \Big|_{\lambda=0} = \sum_{n'} \frac{H'_{nn'} H'_{n'i}}{(\epsilon_i^0 - \epsilon_{n'}^0)(\epsilon_i^0 - \epsilon_n^0)} - \frac{H'_{ii} H'_{ni}}{(\epsilon_i^0 - \epsilon_n^0)^2}, \quad (11)$$

provided only that $\sum_{n'} H'_{nn'} (\partial^2 b_{n'} / \partial \lambda^2)$ converges for $\lambda \rightarrow 0$.

In order to evaluate the sum we first obtain from $\sum_n |b_n|^2 = 1$ (choosing b_i real) that

$$\frac{\partial^2 b_i}{\partial \lambda^2} \Big|_0 = - \sum_{n'} \frac{H'_{in'} H'_{n'i}}{(\epsilon_i^0 - \epsilon_{n'}^0)^2}.$$

Then $\sum_{n'} H'_{nn'} (\partial^2 b_{n'} / \partial \lambda^2)$ can be evaluated, and it is easily shown that it exists if the third-order correction to b_n exists in perturbation theory.

In general, we can obtain all higher derivatives of b_n by differentiating Eq. (5) and requiring that, for $\partial^p b_n / \partial \lambda^p$ to exist, the following must also exist:

$$\lambda \sum_{n'} H'_{nn'} \frac{\partial^p b_{n'}}{\partial \lambda^p} \rightarrow 0 \quad \text{as } \lambda \rightarrow 0$$

or

$$\sum_{n'} H'_{nn'} \frac{\partial^p b_{n'}}{\partial \lambda^p}.$$

The resulting equation for $(1/p!) (\partial^p b_n / \partial \lambda^p)$ is exactly the same equation that one gets for b_{np} if we had assumed a solution of the form

$$b_n(\lambda) = \sum_p b_{np} \lambda^p, \\ \epsilon_i(\lambda) = \sum_p \epsilon_{ip} \lambda^p,$$

and equated powers of λ to zero, as one does in a conventional derivation. Hence one finds that

$$\frac{1}{p!} \frac{\partial^p b_n}{\partial \lambda^p} \Big|_0 = b_{np}, \\ \frac{1}{(p+1)!} \frac{\partial^{p+1} \epsilon_i(\lambda)}{\partial \lambda^{p+1}} \Big|_0 = \epsilon_{i,p+1}, \quad (12)$$

where b_{np} and $\epsilon_{i,p+1}$ are the results from perturbation theory, provided only that these quantities exist and that

$$\sum_{n'} H'_{nn'} \frac{\partial^p b_{n'}}{\partial \lambda^p} \Big|_0 \quad \text{converges,}$$

i.e.,

$$F_p \equiv \sum_{n'} H'_{nn'} b_{n',p} \quad \text{converges.} \quad (13)$$

Let us assume that the quantities in Eq. (12) exist and determine the conditions under which F_p exists. Equating powers of λ^{p+1} in Eq. (5) to zero, we note that F_p can be rewritten as a sum of products of $\epsilon_{im} b_{n,p+1-m}$, $m=0,1,\dots,p+1$, all of which exist by Eq. (12)—except possibly $b_{n,p+1}$. Hence, if the quantities

in Eq. (12) exist, F_p exists if and only if the $(p+1)$ th correction to b_n exists in perturbation theory. Hence the theorem is proven.

Now if $g(\lambda)|_0, g'(\lambda)|_0, \dots, g^{(p)}(\lambda)|_0$ all exist, then Taylor's formula with remainder is

$$g(\lambda) = g|_0 + g'|_0 \lambda + \frac{g''}{2} \Big|_0 \lambda^2 + \dots + \frac{g^{(p)}}{p!} \Big|_0 \lambda^p \\ + \frac{g^{(p+1)}(\theta \lambda)}{(p+1)!} \lambda^{p+1}, \quad 0 < \theta < 1. \quad (14)$$

Hence $\epsilon_i(\lambda)$ and $\psi_i(\lambda)$ can be written as such an expansion with the coefficients being given by perturbation theory. From this we can deduce several corollaries.

Corollary 1: If $\epsilon_{i,p}$ and b_{np} exist for all finite p in perturbation theory, then the perturbation series is an asymptotic expansion for both the energy and the wavefunction.

Proof: We want to prove that if

$$g = \sum_{n=0}^N \frac{g^{(n)}}{n!} \Big|_0 \lambda^n + \frac{g^{(N+1)}(\theta \lambda)}{(N+1)!} \lambda^{N+1} \quad (15)$$

for arbitrary N , then

$$\lim_{\lambda \rightarrow 0} \lambda^{-p} \left[g - \sum_{n=0}^p \frac{g^{(n)}}{n!} \Big|_0 \lambda^n \right] \rightarrow 0 \quad (16)$$

for arbitrary p .

Take $N = p$. Then

$$\lim_{\lambda \rightarrow 0} \lambda^{-p} \left[g - \sum_{n=0}^p \frac{g^{(n)}}{n!} \Big|_0 \lambda^n \right] \\ = \lim_{\lambda \rightarrow 0} \lambda^{-p} \left[\lambda^{p+1} \frac{g^{(p+1)}}{(p+1)!}(\theta \lambda) \right] = \lim_{\lambda \rightarrow 0} 0(\lambda) = 0, \quad (17)$$

which is Eq. (16).

Corollary 2: If $\epsilon_{i,p}$ and b_{np} exist for all finite p in perturbation theory, then the perturbation series is unique, even though it may not converge.

The proof of this statement follows from the first corollary, since the asymptotic expansion of a function is unique. It should be noted, however, that if by some means the perturbation expansion can be "summed"; i.e., if we can find an analytic expression whose expansion is exactly the same as the perturbation series term-by-term, it does not follow that the expression is the exact solution of the problem. This statement follows from the fact that two different functions may have the same asymptotic expansions, or, said another way, there are functions whose

asymptotic expansion is identically zero, i.e., $\exp - (1/\lambda^2)$, which can be added to the result of the "summation" and which gives the same asymptotic expansion. However, if the perturbation series actually converges, then the series defines a unique function and there is no ambiguity concerning the appropriate exact solution.

Corollary 3: If $\epsilon_{i,p+1}$ is greater (less) than zero, then, for sufficiently small λ , the sum of the series up to and including $\epsilon_{i,p}$ underestimates (overestimates) $\epsilon_i(\lambda)$. The proof follows immediately from the sign of the remainder term in Eq. (14), since, for sufficiently small λ , $g^{(p+1)}(\theta\lambda)$ always has the same sign as $g^{(p+1)}(0)$, and hence the remainder has the same sign as $g^{(p+1)}(0)$ (for $\lambda > 0$). Then if the remainder is positive (negative), taking only terms up to λ^p in the series must underestimate (overestimate) the true result.

From the corollary it immediately follows that if two consecutive terms in the perturbation expansion have the same sign, then, for sufficiently small λ , the partial sum including the first term is always closer to the exact result than the partial sum up to, but not including, the first term. Furthermore, it also follows from the corollary that if two consecutive terms in the perturbation series have opposite sign, then for sufficiently small λ , the exact result lies between the partial sum including the first term and the partial sum up to, but not including, the first term.

It is sometimes stated that the odd orders of perturbation theory overestimate the perturbed ground state energy.^{4,5} We see from our work that this can be true only if all the even terms are negative. This is, of course, the case for $p = 2$, but we see no reason why it is true in general. The "proof" generally given that the odd orders overestimate the energy is merely the statement that the terms through $O(\lambda^{2n+1})$ in perturbation theory are the same as those given by the expectation value of H if we take as our wavefunction the perturbation-theory result to $O(\lambda^n)$. However, with such a wavefunction we obtain an upper bound on the perturbed ground-state energy containing terms to $O(\lambda^{2n+1})$, which is not the same expression as the $(2n + 1)$ th-order result in perturbation theory. This follows from the fact that in perturbation theory there are contributions to the energy to $O(\lambda^{2n+1})$ from terms in the wavefunction of orders as high as $O(\lambda^{2n})$. On the other hand, if we take the expectation value of H with the perturbation-theory wavefunction through $O(\lambda^{2n})$, we obtain an expression for the energy having terms up to $O(\lambda^{4n+1})$, and it is this expression [and not the sum of the terms to $O(\lambda^{2n+1})$] which, by the variational theorem, overestimates the energy. Once we drop the terms of higher order than λ^{2n+1} , we can no longer be certain that what is left still overestimates the energy.

⁴ B. F. Gray, *J. Chem. Phys.* **29**, 276 (1958).

⁵ P. M. Morse and H. Feshbach, *Methods of Theoretical Physics* (McGraw-Hill Book Company, Inc., New York 1953), Vol. II, pp. 1119-20.

Numerical Solution of a Singular Integral Equation Encountered in Polymer Physics*

D. W. SCHLITT
University of Nebraska, Lincoln, Nebraska

(Received 14 July 1967)

A numerical method for the solution of an integral equation of the type encountered in the Kirkwood-Riseman theory of intrinsic viscosities of flexible macromolecules is investigated. The absolute accuracy and rate of convergence of the method are evaluated for a special case and the results of this method are compared with another recently proposed method of solution.

I. INTRODUCTION

A theory of the intrinsic viscosities and translational diffusion constants of flexible macromolecules has been developed by Kirkwood and Riseman.¹ The theory requires the solution of a linear integral equation of the Fredholm type with an unbounded kernel. The solutions of this equation have been discussed by Auer and Gardner,² and numerical methods have been applied to the problem by Ullman.^{3,4} The numerical methods, which have been applied, involve the usual reduction of the integral equation to a set of linear algebraic equations. The fact that the kernel is unbounded introduces complications which were resolved by Ullman by averaging over the singular region. In this paper, an alternative method is considered and the results of the two methods are compared.

II. EQUATION AND METHODS OF SOLUTION

The integral equation of interest is of the form

$$\varphi(x) = f(x) - \lambda \int_{-1}^1 \varphi(t) |x - t|^{-\alpha} dt \quad (1)$$

with $0 < \alpha < 1$. The approximation procedure which we use is based on the Gaussian quadrature formula. The integral is converted into a sum of weights times the value of the integrand at points $t = x_i$ and we evaluate the equation for values of $x = x_i$. This means that, for one term in the sum, $|x_i - t|$ will vanish. Ullman and Ullman⁴ replace that term by an appropriate average value for the integral in the vicinity of $t = x_i$. An alternative method⁵ is to rewrite the integral equation so that the singular contribution to the integrand cancels out. If we add

and subtract

$$\lambda \varphi(x) \int_{-1}^1 |x - t|^{-\alpha} dt$$

to Eq. (1), we get

$$\varphi(x) = f(x) - \lambda \int_{-1}^1 [\varphi(t) - \varphi(x)] |x - t|^{-\alpha} dt - \lambda \varphi(x) \int_{-1}^1 |x - t|^{-\alpha} dt$$

and if we let

$$g_\alpha(x) = \int_{-1}^1 |x - t|^{-\alpha} dt, \quad (2)$$

we have

$$\varphi(x) = f(x)[1 + \lambda g_\alpha(x)]^{-1} - \lambda [1 + \lambda g_\alpha(x)]^{-1} \times \int_{-1}^1 [\varphi(t) - \varphi(x)] |x - t|^{-\alpha} dt. \quad (3)$$

The kernel in Eq. (2) is considerably more complicated, but if $\varphi(x)$ is continuous and has a continuous first derivative for $-1 < x < 1$, then the integrand is bounded. The integral defining g_α can be evaluated for the values of α which are of interest, giving

$$g_\alpha(x) = (1 - \alpha)^{-1} [(x - 1)^{-\alpha+1} + (1 - x)^{-\alpha+1}]. \quad (2')$$

The integral in Eq. (3) can be converted to a sum using a Gaussian quadrature formula which reduces the integral equation to a set of linear algebraic equations for the value of φ at a set of fixed points defined by the quadrature formula. If x_i denotes the i th abscissa and w_i its weight, then we must solve the equations

$$\sum_{j=1}^N A_{ij} y_j = z_i \quad (4)$$

for the y_i where

$$y_i = \varphi(x_i),$$

$$z_i = f(x_i)[1 + \lambda g_\alpha(x_i)]^{-1}, \quad (5)$$

and

$$A_{ij} = \lambda w_j [1 + \lambda g_\alpha(x_i)]^{-1} |x_i - x_j|^{-\alpha}; \quad i \neq j,$$

$$A_{ii} = 1 - \sum_{j=1}^N (1 - \delta_{ij}) A_{ij}. \quad (6)$$

* Supported in part by National Science Foundation Grant GP-5373.

¹ J. G. Kirkwood and J. Riseman, *J. Chem. Phys.* **16**, 565 (1948).

² P. L. Auer and C. S. Gardner, *J. Chem. Phys.* **23**, 1545 (1955); **23**, 1546 (1955).

³ R. Ullman, *J. Chem. Phys.* **40**, 2193 (1964).

⁴ N. Ullman and R. Ullman, *J. Math. Phys.* **7**, 1743 (1966).

⁵ The method was suggested to the author by Dr. Grahame Frye (private communication) several years ago in connection with an entirely different problem and is probably very widely known.

III. RESULTS

In order to evaluate the usefulness of this method and to allow comparison with the published results of Ullman and Ullman,⁴ Eq. (4) was solved for the case where

$$f(x) = x^2 \quad \text{and} \quad \alpha = \frac{1}{2}.$$

A code was written in FORTRAN for the IBM 360 model 50 using a standard matrix inversion subroutine with full pivotal condensation. Execution time was not excessive: roughly 0.04 h was required to compute the result for one value of λ with $N = 80$.

The use of a Gaussian integration method has one drawback; the abscissas depend on the value of N . This complicates the comparison of the results for different N . In order to avoid this problem an attempt was made to use Eq. (1) for the evaluation of $\varphi(x)$ for $x \neq x_i$. The integral on the right in Eq. (1) depends only on the y_i and thus can be replaced by a sum with all terms finite. This method proved to be extremely inaccurate and was of no value.

The reliability of the numerical work was checked in two ways. First, the ordering of the linear equations was changed in some check calculations and the variations in the values of $\varphi(x)$ were observed, and second, the computed values of $\varphi(x)$ and $\varphi(-x)$ were compared (they should be equal). Both tests lead to the same conclusion that errors from this source are not detectable in the numbers reported here, although they

are quite easy to see in the single-precision output from the computer.

A number of comparisons were made in order to estimate the absolute accuracy of the method and its accuracy relative to the method of Ref. 4. These comparisons are indicated in the graphical and tabular materials. Figures 1 and 2 show a direct comparison between the results in this paper and Ref. 4. There is a definite systematic difference between the two methods. Examination of Tables I-IV shows that the differences are generally between 1 and 0.1 per cent. Since the applications of the solutions of Eq. (1) to polymer physics involve an integral of $\varphi(x)$ over x , errors of this magnitude could be serious. The effect is probably minimized because the differences change sign as can be seen in Figs. 1 and 2.

The absolute error can be investigated by comparing the solutions obtained here with an exact solution for $\lambda \rightarrow \infty$. This solution is^{2,4}

$$\varphi_a(x) = \sqrt{2} (4x^2 - 1) / 3\pi\lambda(1 - x^2)^{\frac{1}{2}}.$$

The comparison is found in Table V. Notice that both the results reported here and those of Ref. 4 increase in magnitude as λ increases. In addition, the results of Ref. 4 are all larger in magnitude than the asymptotic value while the results reported here are smaller in magnitude and appear to approach the asymptotic value.

The convergence of the solution with increasing N

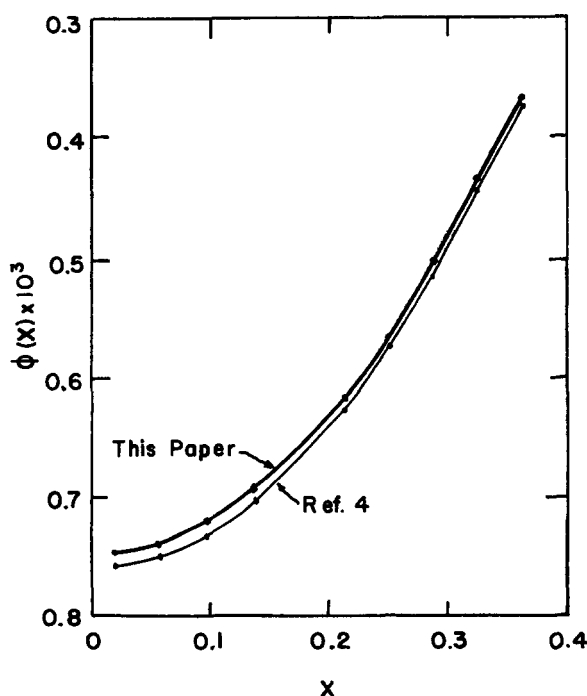


FIG. 1. A plot of $\varphi(x)$ vs x for $0 < x < 0.4$ and $\lambda = 200$. The upper line is the present work, the lower is Ref. 4.

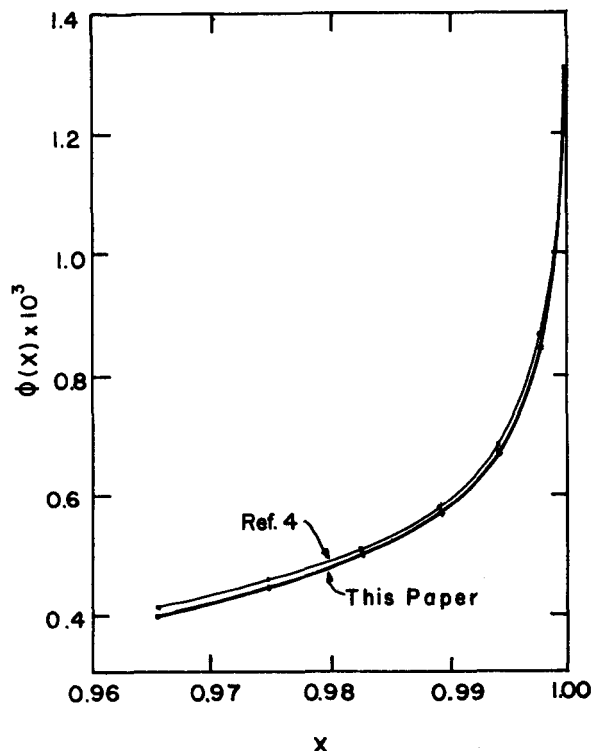


FIG. 2. A plot of $\varphi(x)$ vs x for $0.96 < x < 1.00$ and $\lambda = 200$. The lower line is the present work, the upper is Ref. 4.

TABLE I. A comparison of two methods of solution for $\lambda = 0.5$ and $N = 20, 40,$ and $80.$

x		$N = 20$	$N = 40$	$N = 80$
0.019511	This paper Ref. 4			-0.081002 -0.081611
0.116084	This paper Ref. 4	-0.07370	-0.073751 -0.074555	-0.07376
0.502804	This paper Ref. 4.	0.06252	0.06254	0.062541 0.062666
0.778305	This paper Ref. 4.	0.28371	0.28398 0.28551	0.28399
0.999554	This paper Ref. 4.			0.67454 0.67470

TABLE II. A comparison of two methods of solutions for $\lambda = 5$ and $N = 20, 40,$ and $80.$

x		$N = 20$	$N = 40$	$N = 80$
0.019511	This paper Ref. 4.			-0.024896 -0.025238
0.116084	This paper Ref. 4.	-0.02351	-0.023537 -0.024002	-0.23543
0.502804	This paper Ref. 4.	0.00273	0.00271	0.002707 0.002638
0.778305	This paper Ref. 4.	0.05011	0.050288 0.050836	0.05028
0.999554	This paper Ref. 4.			0.25986 0.25978

TABLE III. A comparison of two methods of solution for $\lambda = 20$ and $N = 20, 40,$ and $80.$

x		$N = 20$	$N = 40$	$N = 80$
0.019511	This paper Ref. 4.			-0.0071355 -0.0075413
0.116084	This paper Ref. 4.	-0.006762	-0.0067718 -0.0069158	-0.006773
0.502804	This paper Ref. 4.	0.000282	0.000272	0.00027054 0.00024289
0.778305	This paper Ref. 4.	0.01322	0.013269 0.013422	0.01327
0.999554	This paper Ref. 4.			0.101496 0.101630

TABLE IV. A comparison of two methods of solution for $\lambda = 200$ and $N = 20, 40,$ and $80.$

x		$N = 20$	$N = 40$	$N = 80$
0.019511	This paper Ref. 4.			-0.00074540 -0.00075671
0.116084	This paper Ref. 4.	-0.0007070	-0.00070822 -0.00072363	0.0007085
0.502804	This paper Ref. 4.4	0.0000127	0.0000114	0.000011018 0.0000079162
0.778305	This paper Ref. 4.	0.001341	0.0013461 0.0013618	0.001346
0.999554	This paper Ref. 4.			0.012984 0.013045

TABLE V. A comparison of two methods of solution with $N = 80$ and $\lambda = 100$ and 200 with the exact solution for λ infinite.

x		$\lambda\varphi, \lambda = 100$	$\lambda\varphi, \lambda = 200$	$\lambda\varphi_\infty$
0.019511	This paper Ref. 4.	-0.14835 -0.15059	-0.14908 -0.15134	-0.14984
0.116084	This paper Ref. 4.	-0.14099	-0.14170	-0.14247
0.250952	This paper Ref. 4.	-0.11276 -0.11462	-0.11342 -0.11529	-0.11409
0.502804	This paper Ref. 4.	0.002597	0.002220 0.001583	0.001815
0.778305	This paper Ref. 4.	0.26871	0.26911	0.26948
0.982849	This paper Ref. 4.	0.99057 1.00001	0.99614 1.00574	1.00072
0.999554	This paper Ref. 4.	2.51135 2.52160	2.59671 2.60904	2.60159

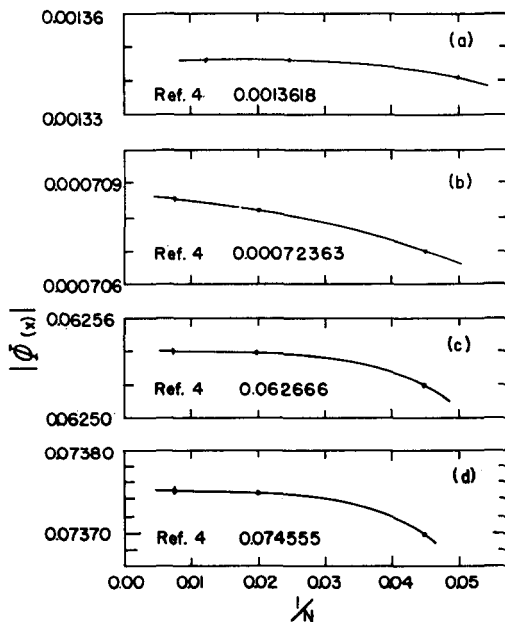


FIG. 3. Typical plots of $\varphi(x)$ vs $1/N$ for selected values of x and λ . Curve (a) $x = 0.77, \lambda = 200,$ (b) $x = 0.11, \lambda = 200,$ (c) $x = 0.50, \lambda = 0.5,$ and (d) $x = 0.11, \lambda = 0.5.$

was investigated and the results are shown in Tables II-V and in Fig. 3. As was pointed out earlier, the Gaussian integration method does not give values of φ at the same x for different N . In order to make these comparisons, values of φ at $x = 0.116084$, 0.802804 , and 0.778305 were determined by an Aitken four-point interpolation. The accuracy of the interpolation was checked by interpolating to a value of x for which φ was determined and comparing the two results. This confirmed that the interpolations were accurate to the number of places shown in the tables. As can be seen from Fig. 3, the relative convergence is good. If the convergence is to the exact result, then the results of Ref. 4 are less accurate for $N = 80$ than those reported here are for $N = 20$.

From these observations we conclude that the

method proposed here is probably a better approach to the numerical solution of Eq. (1) than that proposed in Ref. 4 on the basis of speed of convergence and on overall accuracy.

Insofar as could be determined by the present author, the method of solution suggested in Ref. 2 has not been investigated as a method of numerical solution of the equations. It would be interesting to compare that method with this one.

ACKNOWLEDGMENTS

The author acknowledges with gratitude the help of the University of Nebraska Computer Center, in particular, for providing a matrix inversion routine. The help of Ronald McFee in performing the interpolations is also acknowledged.

JOURNAL OF MATHEMATICAL PHYSICS VOLUME 9, NUMBER 3 MARCH 1968

Scattering of Scalar Waves by a Convex Transparent Object with Statistical Surface Irregularities

YUNG MING CHEN

*Department of Applied Analysis, State University of New York at Stony Brook,
Stony Brook, New York*

(Received 27 March 1967)

The scattering of time-harmonic spherical scalar waves by a large, convex, transparent, dense, and three-dimensional object with statistically corrugated surface is considered. The maximum deviation of the corrugated surface from the smooth one is assumed to be small, and hence the boundary-perturbation technique is utilized in this study. First, the scattering of scalar waves by a large, transparent, and dense sphere with statistical surface irregularities is treated as a canonical problem in the general discussion. After the perturbation solution is expanded asymptotically for large ka , it is found that the higher-order solutions can be obtained from the zeroth-order solution in a simple and straightforward manner. Then this relationship is generalized to scattering by a large, convex, transparent, and dense object with statistical surface irregularities; a general recipe is given. Finally, the asymptotic expressions of mean values of the scattered wavefunction and the scattered intensity are given for the general problem.

I. INTRODUCTION

In recent years, an increasing amount of attention has been devoted to the study of the effect of statistical surface irregularities on propagation and scattering of various types of waves. Although the general problem of scattering waves from corrugated surfaces appears to be difficult, a number of investigators have been able to make progress by applying either probability theory or perturbation theory to the problem of scattering of waves by statistically corrugated plane surfaces under a few suitable assumptions. For a systematic classification of existing theories developed for random irregular surfaces, and a rather complete

bibliography, readers are referred to an excellent text by Beckmann and Spizzichino.¹

Here we shall be concerned with the scattering of time-harmonic spherical scalar waves by a large, convex, transparent, and dense object with a statistically corrugated surface. The case of a small transparent sphere with a statistically corrugated surface has been treated by Chen and Kim.² The ratio of the maximum deviation of the corrugated surface from the unperturbed one to the local radius of the scatterer

¹ P. Beckmann and A. Spizzichino, *The Scattering of Electromagnetic Waves from Rough Surfaces* (The Macmillan Company, New York, 1963).

² Y. M. Chen and S. J. Kim, *J. Acoust. Soc. Am.* **42**, 1 (1967).

was investigated and the results are shown in Tables II-V and in Fig. 3. As was pointed out earlier, the Gaussian integration method does not give values of φ at the same x for different N . In order to make these comparisons, values of φ at $x = 0.116084$, 0.802804 , and 0.778305 were determined by an Aitken four-point interpolation. The accuracy of the interpolation was checked by interpolating to a value of x for which φ was determined and comparing the two results. This confirmed that the interpolations were accurate to the number of places shown in the tables. As can be seen from Fig. 3, the relative convergence is good. If the convergence is to the exact result, then the results of Ref. 4 are less accurate for $N = 80$ than those reported here are for $N = 20$.

From these observations we conclude that the

method proposed here is probably a better approach to the numerical solution of Eq. (1) than that proposed in Ref. 4 on the basis of speed of convergence and on overall accuracy.

Insofar as could be determined by the present author, the method of solution suggested in Ref. 2 has not been investigated as a method of numerical solution of the equations. It would be interesting to compare that method with this one.

ACKNOWLEDGMENTS

The author acknowledges with gratitude the help of the University of Nebraska Computer Center, in particular, for providing a matrix inversion routine. The help of Ronald McFee in performing the interpolations is also acknowledged.

JOURNAL OF MATHEMATICAL PHYSICS VOLUME 9, NUMBER 3 MARCH 1968

Scattering of Scalar Waves by a Convex Transparent Object with Statistical Surface Irregularities

YUNG MING CHEN

*Department of Applied Analysis, State University of New York at Stony Brook,
Stony Brook, New York*

(Received 27 March 1967)

The scattering of time-harmonic spherical scalar waves by a large, convex, transparent, dense, and three-dimensional object with statistically corrugated surface is considered. The maximum deviation of the corrugated surface from the smooth one is assumed to be small, and hence the boundary-perturbation technique is utilized in this study. First, the scattering of scalar waves by a large, transparent, and dense sphere with statistical surface irregularities is treated as a canonical problem in the general discussion. After the perturbation solution is expanded asymptotically for large ka , it is found that the higher-order solutions can be obtained from the zeroth-order solution in a simple and straightforward manner. Then this relationship is generalized to scattering by a large, convex, transparent, and dense object with statistical surface irregularities; a general recipe is given. Finally, the asymptotic expressions of mean values of the scattered wavefunction and the scattered intensity are given for the general problem.

I. INTRODUCTION

In recent years, an increasing amount of attention has been devoted to the study of the effect of statistical surface irregularities on propagation and scattering of various types of waves. Although the general problem of scattering waves from corrugated surfaces appears to be difficult, a number of investigators have been able to make progress by applying either probability theory or perturbation theory to the problem of scattering of waves by statistically corrugated plane surfaces under a few suitable assumptions. For a systematic classification of existing theories developed for random irregular surfaces, and a rather complete

bibliography, readers are referred to an excellent text by Beckmann and Spizzichino.¹

Here we shall be concerned with the scattering of time-harmonic spherical scalar waves by a large, convex, transparent, and dense object with a statistically corrugated surface. The case of a small transparent sphere with a statistically corrugated surface has been treated by Chen and Kim.² The ratio of the maximum deviation of the corrugated surface from the unperturbed one to the local radius of the scatterer

¹ P. Beckmann and A. Spizzichino, *The Scattering of Electromagnetic Waves from Rough Surfaces* (The Macmillan Company, New York, 1963).

² Y. M. Chen and S. J. Kim, *J. Acoust. Soc. Am.* **42**, 1 (1967).

is assumed to be small in this study. Hence the scattered wave can be determined by the boundary-perturbation technique,³ which is based on the Taylor expansion of boundary conditions at the perturbed boundary and the representation of the field as a power series in the aforementioned ratio.

First, we shall treat scattering of scalar waves by a transparent sphere with a statistically corrugated surface as a canonical problem in our general discussion. This has the advantage of illustrating the method without introducing extraneous geometrical details; it is also a case for which we can obtain the exact perturbation solution, assuming that the perturbation series converges. The exact perturbation solution is then asymptotically expanded for large ka . It is found that the higher-order solutions can be systematically obtained from the zeroth-order solution in a rather straightforward manner. Hence a general recipe based on the zeroth-order solution, which, in general, can be constructed from the geometrical theory of diffraction,⁴⁻⁷ is given for the treatment of scattering by a large, transparent, and convex object of arbitrary shape with a statistically corrugated surface. Finally, the asymptotic expressions of mean values of the scattered wavefunction and the scattered intensity are given for the general problem.

II. BOUNDARY PERTURBATION

Let any point in a three-dimensional physical space be denoted by vector $\mathbf{r} = (\xi_1, \xi_2, \xi_3)$, where ξ_i are generalized curvilinear coordinates.

Let a random surface \mathcal{S} be defined by

$$\xi_1 = A[1 + \epsilon f(\xi_2, \xi_3, q)], \quad (1)$$

where A is a constant; ϵ is a small parameter such that $|\epsilon f(\xi_2, \xi_3, q)| < 1$; $f(\xi_2, \xi_3, q)$ is a smooth, continuous function of ξ_2 and ξ_3 , such that $|\partial f / \partial \xi_2| < 1$ and $|\partial f / \partial \xi_3| < 1$ for all values of ξ_2, ξ_3 , and q ; q is a random variable ranging over a space X in which a probability density $P(q)$ is defined such that the mean (or average) value of a random function $W(\mathbf{r}, q)$ is defined as

$$\langle W \rangle = \int_X W(\mathbf{r}, q) P(q) dq. \quad (2)$$

Let the random surface \mathcal{S} separate the entire physical space into region 1 and region 2. Let the solution of

a linear differential equation satisfy the following system of equations:

$$LU(\mathbf{r}, q) = F(\mathbf{r}) \quad \text{for } \mathbf{r} \text{ in region 1,} \quad (3)$$

$$LV(\mathbf{r}, q) = G(\mathbf{r}) \quad \text{for } \mathbf{r} \text{ in region 2,} \quad (4)$$

$$B(U, V, \mathbf{n} \cdot \nabla U, \mathbf{n} \cdot \nabla V, \dots) = 0 \quad \text{for } \mathbf{r} \text{ on } \mathcal{S}, \quad (5)$$

and some uniqueness conditions of U and V in region 1 and region 2, respectively, where \mathbf{n} is the unit normal vector of \mathcal{S} .

Now assuming that both U and V are analytic functions of parameter ϵ , we can expand U and V in a power series of ϵ :

$$U(\mathbf{r}, q) = V_0(\mathbf{r}) + \sum_{j=1}^{\infty} \epsilon^j U_j(\mathbf{r}, q) \quad (6)$$

and

$$V(\mathbf{r}, q) = V_0(\mathbf{r}) + \sum_{j=1}^{\infty} \epsilon^j V_j(\mathbf{r}, q). \quad (7)$$

Since

$$\begin{aligned} \frac{\partial}{\partial n} &= \mathbf{n} \cdot \nabla = \left| \sum_{i=1}^3 \mathbf{e}_i h_i^{-1} \frac{\partial \mathcal{S}}{\partial \xi_i} \right|^{-1} \sum_{i=1}^3 \mathbf{e}_i h_i^{-1} \frac{\partial \mathcal{S}}{\partial \xi_i} \cdot \sum_{i=1}^3 h_i^{-1} \frac{\partial}{\partial \xi_i} \\ &= h_1^{-1} \left\{ 1 - \epsilon^2 \frac{h_1^2 A^2}{2} \left[\frac{1}{h_2^2} \left(\frac{\partial f}{\partial \xi_2} \right)^2 + \frac{1}{h_3^2} \left(\frac{\partial f}{\partial \xi_3} \right)^2 \right] \right\} \frac{\partial}{\partial \xi_1} \\ &\quad - \epsilon \frac{h_1}{h_2^2} A \frac{\partial f}{\partial \xi_2} \frac{\partial}{\partial \xi_2} - \epsilon \frac{h_1}{h_3^2} A \frac{\partial f}{\partial \xi_3} \frac{\partial}{\partial \xi_3} + O(\epsilon^3), \quad (8) \end{aligned}$$

expressions for $\mathbf{n} \cdot \nabla U$, $\mathbf{n} \cdot \nabla V$, and others can be easily derived. Next we expand U , V , $\mathbf{n} \cdot \nabla U$, $\mathbf{n} \cdot \nabla V$, \dots , in a Taylor series about the unperturbed boundary $\xi_1 = A$. By setting the coefficients of similar powers of ϵ in the boundary condition (5) equal to each other, we obtain boundary conditions for the zeroth-order and j th-order solutions, respectively:

$$B(U_0, V_0, \mathbf{n}_0 \cdot \nabla U_0, \mathbf{n}_0 \cdot \nabla V_0, \dots) \Big|_{\xi_1=A} = 0 \quad (9)$$

$$B_j \left[S_{2j-1} \left(U_0, \mathbf{n}_0 \cdot \nabla U_0, V_0, \mathbf{n}_0 \cdot \nabla V_0, \dots, U_j, \mathbf{n}_0 \cdot \nabla U_j, V_j, \mathbf{n}_0 \cdot \nabla V_j, \dots, f, \frac{\partial f}{\partial \xi_2}, \frac{\partial f}{\partial \xi_3}, \dots \right), \right.$$

$$\left. S_{2j} \left(U_0, \mathbf{n}_0 \cdot \nabla U_0, V_0, \mathbf{n}_0 \cdot \nabla V_0, \dots, U_j, \mathbf{n}_0 \cdot \nabla U_j, V_j, \mathbf{n}_0 \cdot \nabla V_j, \dots, f, \frac{\partial f}{\partial \xi_2}, \frac{\partial f}{\partial \xi_3}, \dots \right) \right] \Big|_{\xi_1=A} = 0$$

$$j = 1, 2, 3, \dots, \quad (10)$$

where \mathbf{n}_0 is the unit normal vector of surface $\xi_1 = A$. Upon substituting (6) and (7) into (3) and (4), respectively, we get

$$\begin{cases} LU_0 = F & \text{for } \mathbf{r} \text{ in region 1,} \\ LU_j = 0 & j = 1, 2, 3, \dots, \end{cases} \quad (11)$$

³ P. M. Morse and H. Feshbach, *Methods of Theoretical Physics* (McGraw-Hill Book Company, Inc., New York, 1953), Vol. II.

⁴ J. B. Keller, "A Geometrical Theory of Diffraction," Symposium on Microwave Optics, Eaton Electronics Res. Lab., McGill Univ., Montreal, Canada (1953).

⁵ J. B. Keller, *J. Opt. Soc. Am.* **52**, 116 (1962).

⁶ Y. M. Chen, *J. Math. Phys.* **5**, 820 (1964).

⁷ Y. M. Chen, *J. Math. Phys.* **6**, 1332 (1965).

and

$$\begin{cases} LV_0 = G & \text{for } \mathbf{r} \text{ in region 2,} \\ LV_j = 0 & j = 1, 2, 3, \dots \end{cases} \quad (12)$$

In view of the above derivation, the j th-order solution is generated by the equivalent sources S_{2j-1} and S_{2j} , arising from the interaction of all the lower-order solutions with surface irregularities. Since the boundary condition for the solution of each order is matched at the unperturbed surface $\xi_1 = A$, rather than at S , the essential effect of this boundary-perturbation technique is to transform the original boundary-value problem to an equivalent boundary-value problem with the unperturbed body as the scatterer on which equivalent sources are induced.

III. PERTURBATION SOLUTION OF A ROUGH SPHERE

The perturbation solution of the scattering of a time-harmonic ($e^{i\omega t}$) spherical wave by a penetrable sphere with statistically corrugated surface has been obtained by Chen and Kim,² but we will reproduce it here briefly for our convenience.

Let the surface of an almost spherical obstacle be defined by

$$\tilde{r} = a + bf(\theta, \varphi, q), \quad (13)$$

where $(\tilde{r}, \theta, \varphi)$ are the spherical coordinates of a typical point on the surface and where $f(\theta, \varphi, q)$ is a smooth continuous function of θ and φ which satisfies conditions $|(b/a)f(\theta, \varphi, q)| < 1$, $|\partial f/\partial \theta| < 1$, and $|\partial f/\partial \varphi| < 1$ for all values of θ, φ , and q .

The wavefunction $u(r, \theta, \varphi)$, caused by a point source located at $(r_0, 0, 0)$ (see Fig. 1), satisfies the following equations:

$$(\nabla^2 + k_1^2)[u_i + u_s(q)] = \frac{\delta(r - r_0)\delta(\theta)}{2\pi r^2 \sin \theta}, \quad r \geq \tilde{r}, \quad (14)$$

$$(\nabla^2 + k_2^2)u_t(q) = 0, \quad r \leq \tilde{r}, \quad (15)$$

$$u_t(q) < \infty \quad \text{at } r = 0, \quad (16)$$

$$\lim_{r \rightarrow \infty} r[\partial(u_i + u_s)/\partial r + ik_1(u_i + u_s)] = 0, \quad (17)$$

$$(u_i + u_s) = \alpha u_t \quad \text{at } r = \tilde{r}, \quad (18)$$

and

$$\partial(u_i + u_s)/\partial n = \beta \partial u_t/\partial n \quad \text{at } r = \tilde{r}. \quad (19)$$

If the parameter $\epsilon = b/a$ is small enough, from (6) and (7) we obtain

$$u_i + u_s = u_i + u_{s0} + \epsilon u_{s1} + \epsilon^2 u_{s2} + O(\epsilon^3) \quad (20)$$

and

$$u_t = u_{t0} + \epsilon u_{t1} + \epsilon^2 u_{t2} + O(\epsilon^3). \quad (21)$$

From Eqs. (9) and (10) the boundary conditions

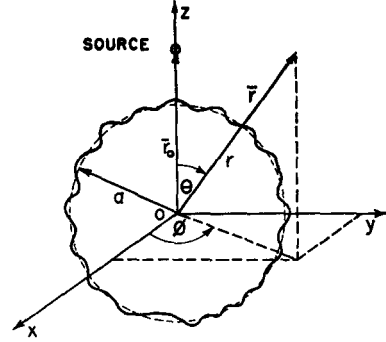


Fig. 1. The geometry of a statistically corrugated sphere is shown.

for the zeroth-order, first-order, and second-order solutions are, respectively,

$$\left. \begin{aligned} u_{0\alpha} &\equiv u_i + u_{s0} - \alpha u_{t0} = 0, \\ \frac{\partial u_{0\beta}}{\partial r} &\equiv \partial(u_i + u_{s0} - \beta u_{t0})/\partial r = 0, \end{aligned} \right\} \quad \text{at } r = a, \quad (22)$$

$$\left. \begin{aligned} [u_{s1} - \alpha u_{t1} &= -af\partial u_{0\alpha}/\partial r]_{r=a} \equiv S_1(\theta, \varphi, q) \\ [\partial u_{s1}/\partial r - \beta \partial u_{t1}/\partial r &= -af\partial^2 u_{0\beta}/\partial r^2 \\ &+ a^{-1}(\partial f/\partial \theta)\partial u_{0\beta}/\partial \theta]_{r=a} \equiv S_2(\theta, \varphi, q) \end{aligned} \right\}, \quad (23)$$

and

$$\left. \begin{aligned} [u_{s2} - \alpha u_{t2} &= \frac{1}{2}a^2 f^2 \partial^2 u_{0\alpha}/\partial r^2]_{r=a} \equiv S_3(\theta, \varphi, q) \\ [\partial u_{s2}/\partial r - \beta \partial u_{t2}/\partial r &= \frac{1}{2}a^2 f^2 \partial^3 u_{0\beta}/\partial r^3 \\ &- a^{-1}f(\partial f/\partial \theta)(\partial u_{0\beta}/\partial \theta) - a^{-1}(\partial f/\partial \theta)\partial u_{1\beta}/\partial \theta \\ &+ a^{-1} \sin^{-2} \theta (\partial f/\partial \varphi)\partial u_{1\beta}/\partial \varphi]_{r=a} \equiv S_4(\theta, \varphi, q) \end{aligned} \right\}. \quad (24)$$

The other S 's can be obtained in a similar way.

By the method of eigenfunction expansion, the zeroth-order and higher-order solutions can be easily obtained as

$$u_i = -\frac{e^{-ik_1|r-r_0|}}{4\pi|\mathbf{r}-\mathbf{r}_0|} = \sum_{n=0}^{\infty} A_n P_n(\cos \theta) j_n(k_1 r_<) h_n^{(2)}(k_1 r_>), \quad (25)$$

$$u_{s0} = \sum_{n=0}^{\infty} B_n P_n(\cos \theta) h_n^{(2)}(k_1 r_0) h_n^{(2)}(k_1 r), \quad (26)$$

$$u_{t0} = \sum_{n=0}^{\infty} C_n P_n(\cos \theta) j_n(k_2 r), \quad (27)$$

$$u_{sj} = \sum_{n=0}^{\infty} \sum_{m=0}^{\infty} (D_{nm}^j Y_{nm}^e + E_{nm}^j Y_{nm}^o) h_n^{(2)}(k_1 r), \quad j = 1, 2, 3, \dots, \quad (28)$$

and

$$u_{tj} = \sum_{n=0}^{\infty} \sum_{m=0}^{\infty} (F_{nm}^j Y_{nm}^e + G_{nm}^j Y_{nm}^o) j_n(k_2 r), \quad j = 1, 2, 3, \dots, \quad (29)$$

where

$$A_n = ik_1(4\pi)^{-1}(2n+1), \quad (30)$$

$$B_n = -A_n \Delta_n^{-1} [\alpha j_n(k_2 a) j'_n(k_1 a) - \beta N j'_n(k_2 a) j_n(k_1 a)], \quad (31)$$

$$C_n = -ik_1^{-2} a^{-2} A_n \Delta_n^{-1} h_n^{(2)}(k_1 r_0), \quad (32)$$

$$D_{nm}^j = \Delta_n^{-1} [\alpha k_1^{-1} j_n(k_2 a) L_{nm}^{2j} - \beta N j'_n(k_2 a) L_{nm}^{2j-1}], \quad (33)$$

$$E_{nm}^j = \Delta_n^{-1} [\alpha k_1^{-1} j_n(k_2 a) M_{nm}^{2j} - \beta N j'_n(k_2 a) M_{nm}^{2j-1}], \quad (34)$$

$$F_{nm}^j = \Delta_n^{-1} [k_1^{-1} h_n^{(2)}(k_1 a) L_{nm}^{2j} - h_n^{(2)'}(k_1 a) L_{nm}^{2j-1}], \quad (35)$$

$$G_{nm}^j = \Delta_n^{-1} [k_1^{-1} h_n^{(2)}(k_1 a) M_{nm}^{2j} - h_n^{(2)'}(k_1 a) M_{nm}^{2j-1}], \quad (36)$$

$$L_{nm}^j = \Gamma_{nm} \int_0^{2\pi} \int_0^\pi S_j(\theta', \varphi', q) Y_{nm}^e \sin \theta' d\theta' d\varphi', \quad (37)$$

$$M_{nm}^j = \Gamma_{nm} \int_0^{2\pi} \int_0^\pi S_j(\theta', \varphi', q) Y_{nm}^0 \sin \theta' d\theta' d\varphi', \quad (38)$$

$$r_< = \min(r, r_0), \quad r_> = \max(r, r_0), \quad (39)$$

$$Y_{nm}^e = \cos m\varphi P_n^m(\cos \theta), \quad (40)$$

$$Y_{nm}^0 = \sin m\varphi P_n^m(\cos \theta), \quad (41)$$

$$N = k_2/k_1 > 1, \quad (42)$$

$$\Delta_n = \alpha h_n^{(2)'}(k_1 a) j_n(k_2 a) - \beta N h_n^{(2)}(k_1 a) j'_n(k_2 a), \quad (43)$$

$$\Gamma_{nm} = \epsilon_m \frac{(2n+1)(n-m)!}{4\pi(n+m)!}, \quad (44)$$

and

$$\epsilon_0 = 1, \quad \epsilon_1 = \epsilon_2 = \epsilon_3 = \dots = 2. \quad (45)$$

In view of the above calculation, the first-order solution is essentially generated by equivalent sources S_1 and S_2 arising from the interaction of the zeroth-order wave with surface irregularities, and the second-order solution is essentially generated by equivalent sources S_3 and S_4 arising from the interaction of the zeroth-order and the first-order waves with surface irregularities. Obviously this branching process will go on indefinitely. The essential effect of the boundary-perturbation technique is to transform the original boundary-value problem to an equivalent boundary-value problem with the unperturbed body as the scatterer which has induced sources on its surface.

IV. ZERO-ORDER SOLUTION IN THE EXTERIOR OF A LARGE ROUGH SPHERE

The zeroth-order solution is the solution for the case of smooth penetrable sphere. The asymptotic

expansion of the exact solution for large $k_1 a$ has been partially obtained by Rubinow,⁸ but we shall derive it in a more complete manner here. By means of the Poisson summation formula, (26) can be rewritten as

$$u_{s0} = -i4^{-1}(rr_0)^{-\frac{1}{2}} \sum_{l=-\infty}^{\infty} e^{i\pi l} \int_0^\pi \psi_\nu P_{\nu-\frac{1}{2}}(\cos \theta) \mathfrak{G} e^{-i2\pi l \nu} \nu d\nu, \quad (46)$$

where

$$\psi_\nu = \frac{\alpha J_\nu(k_2 a) J'_\nu(k_1 a) - \beta N J'_\nu(k_2 a) J_\nu(k_1 a)}{\alpha J_\nu(k_2 a) H_\nu^{(2)'}(k_1 a) - \beta N J'_\nu(k_2 a) H_\nu^{(2)}(k_1 a)}, \quad (47)$$

and

$$\mathfrak{G} = H_\nu^{(2)}(k_1 r_0) H_\nu^{(2)}(k_1 r). \quad (48)$$

To exhibit the physical interpretation of (46) we expand ψ_ν into a geometric series:

$$\psi_\nu = \frac{1}{2} \left\{ 1 - \mathfrak{H}_1 \left[R_{11} + T_{12} T_{21} \sum_{p=1}^{\infty} R_{22}^{p-1} \mathfrak{H}_2^p \right] \right\}, \quad (49)$$

where

$$R_{11} = -[\alpha \log' H_\nu^{(1)}(k_1 a) - \beta N \log' H_\nu^{(1)}(k_2 a)] \times [\alpha \log' H_\nu^{(2)}(k_1 a) - \beta N \log' H_\nu^{(1)}(k_2 a)]^{-1}, \quad (50)$$

$$R_{22} = -[\alpha \log' H_\nu^{(2)}(k_1 a) - \beta N \log' H_\nu^{(2)}(k_2 a)] \times [\alpha \log' H_\nu^{(2)}(k_1 a) - \beta N \log' H_\nu^{(1)}(k_2 a)]^{-1}, \quad (51)$$

$$\alpha T_{12} = 1 + R_{11}, \quad (52)$$

$$T_{21} = \alpha(1 + R_{22}), \quad (53)$$

$$\log' H_\nu^{(\delta)}(z) = H_\nu^{(\delta)'}(z)/H_\nu^{(\delta)}(z), \quad \delta = 1, 2, \quad (54)$$

$$\mathfrak{H}_1 = H_\nu^{(1)}(k_1 a)/H_\nu^{(2)}(k_1 a), \quad (55)$$

and

$$\mathfrak{H}_2 = H_\nu^{(2)}(k_2 a)/H_\nu^{(1)}(k_2 a). \quad (56)$$

The form (49) is introduced because the first term $(1 - \mathfrak{H}_1 R_{11})$ in the integrand of (46) represents the wave u_{s0} , externally reflected from the sphere in the lit region and negative of incident wave in the shadow region. The p th term in the sum, denoted by u_{s0p} , represents a wave transmitted into the sphere, reflected $p-1$ times internally from the interface, having passed p times through the sphere, and finally transmitted out into the surrounding medium. This interpretation is borne out when the various terms are expanded asymptotically for a short wavelength and evaluated properly.

A. Geometric Optics Wave

The criterion for the proper identification of the geometric optics wave from the asymptotic evaluation of (46) is the existence of real saddle points $\nu_{0\rho}$ ($\rho = 1, 2, \dots$) such that $0 < \nu_{0\rho} < k_1 r_<$. Before evaluating u_{s0} asymptotically by the saddle-point

⁸ S. I. Rubinow, Ann. Phys. 14, 305 (1961).

method, we rewrite (46) as

$$u_{s0} = (i4)^{-1}(rr_0)^{-\frac{1}{2}} \sum_{l=-\infty}^{\infty} e^{i\pi l} \int_0^{\infty} \psi_\nu [Q_{\nu-\frac{1}{2}}^{(1)}(\cos \theta) + Q_{\nu-\frac{1}{2}}^{(2)}(\cos \theta)] \mathcal{G} e^{-i2\pi l \nu} \nu d\nu, \quad (57)$$

where $Q_\nu^{(\delta)}(\cos \theta)$, $\delta = 1, 2$, are given in the appendix.

We now consider

$$u_{s0r} = (i8)^{-1}(rr_0)^{-\frac{1}{2}} \sum_{l=-\infty}^{\infty} e^{i\pi l} (J_{s0rl}^{(1)} + J_{s0rl}^{(2)}), \quad (58)$$

where

$$J_{s0rl}^{(\delta)} = \int_0^{\infty} (1 - \mathcal{R}_{11}) Q_{\nu-\frac{1}{2}}^{(\delta)}(\cos \theta) \mathcal{G} e^{-i\pi l \nu} \nu d\nu, \quad \delta = 1, 2. \quad (59)$$

After using the proper asymptotic forms (see Appendix) in (57), we investigate each term of (57) and find that only $J_{s0r0}^{(1)}$ has any real saddle points. It has two real saddle points, which fall in the ranges $0 < \nu < k_1 a$ and $k_1 a < \nu < k_1 r < \dots$. The corresponding saddle-point equations are

$$\cos^{-1} \frac{\nu_0}{k_1 r} + \cos^{-1} \frac{\nu_0}{k_1 r_0} - 2 \cos^{-1} \frac{\nu_0}{k_1 a} = \theta \quad (60)$$

and

$$\cos^{-1} \frac{\nu_0}{k_1 r} + \cos^{-1} \frac{\nu_0}{k_1 r_0} = \theta, \quad (61)$$

respectively.

Equation (60) has an unique, real solution only for the field point \mathbf{r} lying in the lit region, and the solution is

$$\nu_0 = k_1 a \sin \eta_0, \quad (62)$$

where η_0 is the angle of incidence or reflection. Equation (61) has an unique real solution only for the field point lying in the shadow, and the solution is

$$\nu_0 = k_1 r_0 \sin \varphi_0, \quad (63)$$

where φ_0 is the angle between $\mathbf{r} - \mathbf{r}_0$ and \mathbf{r}_0 .

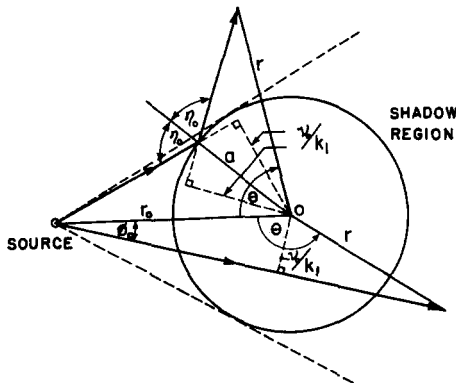


FIG. 2. The geometric rays of u_{s0r} is shown for the case of $N > 1$. The geometrical relations between θ and η_0 and between θ and φ_0 are also shown.

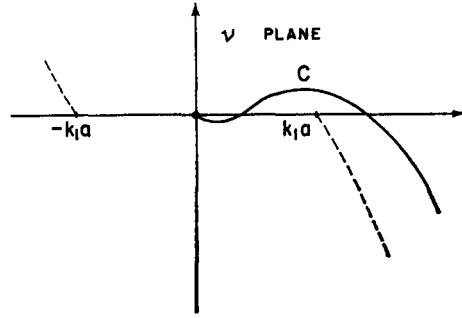


FIG. 3. The saddle-point path C is shown schematically.

By substituting (62) and (63) into (60) and (61), respectively, one obtains the correct geometrical relations between θ and η_0 and between θ and φ_0 (Fig. 2), respectively:

$$2\eta_0 - \sin^{-1} [(a/r) \sin \eta_0] - \sin^{-1} [(a/r_0) \sin \eta_0] = \theta \quad (64)$$

and

$$\frac{1}{2}\pi - \varphi_0 + \cos^{-1} [(r_0/r) \sin \varphi_0] = \theta. \quad (65)$$

The saddle-point path C is shown in Fig. 3. It is found that the end-point contribution of integrals is asymptotically small in comparison with the saddle-point contribution; hence it can be neglected. Finally, the saddle-point contribution of u_{s0r} is

$$u_{s0r}^G \sim \frac{-1}{4\pi} \mathcal{R}_{11}(a \sin \eta_0 \cos \eta_0)^{\frac{1}{2}} \times \{ [2a^{-1}(r^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}}(r_0^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}} - (r^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}} \cos \eta_0 - (r_0^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}} \cos \eta_0] r r_0 \sin \theta \}^{-\frac{1}{2}} \times \exp \{ -ik_1 [(r^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}} + (r_0^2 - a^2 \sin^2 \eta_0)^{\frac{1}{2}} - 2a \cos \eta_0] \} \quad \begin{matrix} \text{in lit region} \\ \text{in shadow region,} \end{matrix} \quad (66)$$

where ζ_0 is the angle of refraction such that

$$\sin \eta_0 = N \sin \zeta_0, \quad (67)$$

and \mathcal{R}_{11} is the Fresnel reflection coefficient in medium 1 such that

$$\mathcal{R}_{11} = \frac{\alpha \cos \eta_0 - \beta N \cos \zeta_0}{\alpha \cos \eta_0 + \beta N \cos \zeta_0}. \quad (68)$$

u_{s0r}^G gives the contribution to the external reflected wave in the lit region and the negative of the incident wave in the shadow.

For the once-transmitted wave, we have to evaluate asymptotically the following equation by the saddle-point method:

$$u_{s01} = (i8)^{-1}(rr_0)^{-\frac{1}{2}} \sum_{l=-\infty}^{\infty} e^{i\pi l} (J_{s01l}^{(1)} + J_{s01l}^{(2)}), \quad (69)$$

where

$$J_{s01}^{(\delta)} = \int_0^\infty \mathcal{H}_1 T_{12} T_{21} \mathcal{H}_2 Q_{\nu-\frac{1}{2}}^{(\delta)}(\cos \theta) \mathcal{G} e^{-i2\pi\nu} \nu d\nu, \quad \delta = 1, 2. \quad (70)$$

It is found that only $J_{s01}^{(1)}$ may have real saddle points such that $0 < \nu_{0\rho} < k_1 a$, $\rho = 1, 2, \dots$. The corresponding saddle-point equations are

$$\begin{aligned} \cos^{-1} \frac{\nu_{0\rho}}{k_1 r} + \cos^{-1} \frac{\nu_{0\rho}}{k_1 r_0} + 2 \cos^{-1} \frac{\nu_{0\rho}}{N k_1 a} \\ - 2 \cos^{-1} \frac{\nu_{0\rho}}{k_1 a} = \theta, \quad \rho = 1, 2, \dots, \end{aligned} \quad (71)$$

whose solutions are also given by (62) with an additional subscript ρ . Equation (71) gives the correct geometric relations:

$$\begin{aligned} 2\eta_{0\rho} + (\pi - 2\zeta_{0\rho}) - \sin^{-1}[(a/r) \sin \eta_{0\rho}] \\ - \sin^{-1}[(a/r_0) \sin \eta_{0\rho}] = \theta, \quad \rho = 1, 2, \dots \end{aligned} \quad (72)$$

Again the end-point contribution of integrals is asymptotically small in comparison with the saddle-point contribution, and a typical saddle-point contribution of u_{s01} is

$$\begin{aligned} u_{s01}^{(\rho)} \sim \left(\frac{-i}{4\pi}\right) \mathcal{G}_{12}^{(\rho)} \mathcal{G}_{21}^{(\rho)} (a \cos \eta_{0\rho} \sin \eta_{0\rho})^{\frac{1}{2}} \\ \times \{ [2a^{-1}(r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}}(r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \\ - (r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \cos \eta_{0\rho} \\ - (r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \cos \eta_{0\rho} \\ - (r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}}(r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \\ \times (Na \cos \zeta_{0\rho})^{-1} 2 \cos \eta_{0\rho} r r_0 \sin \theta]^{-\frac{1}{2}} \\ \times \exp \{ -ik_1 [(r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \\ + (r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} - 2a \cos \eta_{0\rho}] \\ - i2k_2 a \cos \zeta_{0\rho} \}, \end{aligned} \quad (73)$$

where $\mathcal{G}_{12}^{(\rho)}$ is the Fresnel transmission coefficient from medium 1 to medium 2 and $\mathcal{G}_{21}^{(\rho)}$ is the Fresnel transmission coefficient from medium 2 to medium 1, such that

$$\begin{aligned} \mathcal{G}_{12}^{(\rho)} &= \alpha^{-1}(1 + \mathcal{R}_{11}^{(\rho)}) \\ &= \alpha^{-1} \left(1 + \frac{\alpha \cos \eta_{0\rho} - \beta N \cos \zeta_{0\rho}}{\alpha \cos \eta_{0\rho} + \beta N \cos \zeta_{0\rho}} \right) \end{aligned} \quad (74)$$

and

$$\mathcal{G}_{21}^{(\rho)} = \alpha(1 - \mathcal{R}_{11}^{(\rho)}). \quad (75)$$

For the p th transmitted field, we evaluate asymptotically the following equation by the saddle-point method:

$$u_{s0p} = (i8)^{-1} (r r_0)^{-\frac{1}{2}} \sum_{l=-\infty}^{\infty} e^{i\pi l} (J_{s0pl}^{(1)} + J_{s0pl}^{(2)}), \quad (76)$$

where

$$J_{s0pl}^{(\delta)} = \int_0^\infty \mathcal{H}_1 T_{12} T_{21} \mathcal{R}_{22}^{p-1} \mathcal{H}_2 Q_{\nu-\frac{1}{2}}^{(\delta)}(\cos \theta) \mathcal{G} e^{-i2\pi\nu} \nu d\nu, \quad \delta = 1, 2. \quad (77)$$

Only $J_{s0pl}^{(1)}$ may have proper saddle points, and their corresponding saddle-point equations are

$$\begin{aligned} \cos^{-1} \frac{\nu_{0\rho}}{k_1 r} + \cos^{-1} \frac{\nu_{0\rho}}{k_1 r_0} + 2p \cos^{-1} \frac{\nu_{0\rho}}{N k_1 a} \\ - 2 \cos^{-1} \frac{\nu_{0\rho}}{k_1 a} = 2\pi l + \theta, \quad \rho = 1, 2, \dots \end{aligned} \quad (78)$$

Again the relations between $\nu_{0\rho}$ and $\eta_{0\rho}$ are given by (62). By substituting (62) into (78), we have

$$\begin{aligned} 2\eta_{0\rho} + p(\pi - 2\zeta_{0\rho}) - \sin^{-1}[(a/r) \sin \eta_{0\rho}] \\ - \sin^{-1}[(a/r_0) \sin \eta_{0\rho}] = 2\pi l + \theta, \quad \rho = 1, 2, \dots, \end{aligned} \quad (79)$$

where l 's are properly chosen such that (79) is identical with the geometric relations between θ and $\eta_{0\rho}$. A typical saddle-point contribution of u_{s0p} is

$$\begin{aligned} u_{s0p}^{(\rho)} \sim \left(\frac{-1}{4\pi}\right) \mathcal{G}_{12}^{(\rho)} \mathcal{G}_{21}^{(\rho)} \mathcal{R}_{22}^{(p)p-1} (a \sin \eta_{0\rho} \cos \eta_{0\rho})^{\frac{1}{2}} \\ \times \{ [2a^{-1}(r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}}(r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \\ - (r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \cos \eta_{0\rho} \\ - (r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \cos \eta_{0\rho} \\ - (r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}}(r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \\ \times (Na \cos \zeta_{0\rho})^{-1} 2p \cos \eta_{0\rho} r r_0 \sin \theta]^{-\frac{1}{2}} \\ \times \exp \left\{ -ik_1 [(r^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} \right. \\ \left. + (r_0^2 - a^2 \sin^2 \eta_{0\rho})^{\frac{1}{2}} - 2a \cos \eta_{0\rho}] \right. \\ \left. - i2k_2 a p \cos \zeta_{0\rho} + ip \frac{\pi}{2} \right\}, \end{aligned} \quad (80)$$

where \mathcal{R}_{22} is the Fresnel reflection coefficient in medium 2 such that

$$\mathcal{R}_{22}^{(\rho)} = -\mathcal{R}_{11}^{(\rho)}. \quad (81)$$

Except for the difference between the amplitude factors in the two-dimensional and three-dimensional cases, the expressions of $u_{s0r}^{(\rho)}$ and $u_{s0p}^{(\rho)}$ are exactly the same as the case of a transparent circular cylinder,⁶ including the Fresnel coefficients. Hence these expressions can also be obtained by the geometrical theory of diffraction.⁴⁻⁶

B. Diffracted Wave

Any integral in (57) which does not have proper real saddle points may still be evaluated by the method of residues. In order to make the asymptotic evaluation, we first extend the integrals along the positive-real ν axis to the entire real ν axis. To this end, we

find that

$$\begin{aligned}
 u_{s0r} &\cong (8)^{-1}(rr_0)^{-\frac{1}{2}} \left[\sum_{l=1}^{\infty} e^{i\pi l} \int_{-\infty}^{\infty} \mathcal{G}P_{v-\frac{1}{2}}(-\cos \theta) \right. \\
 &\quad \times e^{-i2\pi(l+\frac{1}{2})v} dv \\
 &\quad - \sum_{l=0}^{\infty} e^{i\pi l} \int_{-\infty}^{\infty} R_{11} \mathcal{H}_1 \mathcal{G}P_{v-\frac{1}{2}}(-\cos \theta) \\
 &\quad \times e^{-i2\pi(l+\frac{1}{2})v} dv \left. \right] - u_i^G \quad \text{in the shadow region,} \\
 &\cong u_{s0r}^G - (8)^{-1}(rr_0)^{-\frac{1}{2}} \\
 &\quad \times \left[\sum_{l=0}^{\infty} e^{i\pi l} \int_{-\infty}^{\infty} R_{11} \mathcal{H}_1 \mathcal{G}Q_{v-\frac{1}{2}}^{(2)}(-\cos \theta) \right. \\
 &\quad \times e^{-i2\pi(l+\frac{1}{2})v} dv \\
 &\quad + \sum_{l=0}^{\infty} e^{i\pi(l+1)} \int_{-\infty}^{\infty} R_{11} \mathcal{H}_1 \mathcal{G}Q_{v-\frac{1}{2}}^{(1)}(-\cos \theta) \\
 &\quad \times e^{-i2\pi(l+\frac{3}{2})v} dv \left. \right] \quad \text{in the lit region.} \quad (82)
 \end{aligned}$$

It is easy to see that u_{s0r} has simple poles in the integrand. The positions of the poles of all the integrals are determined by

$$\alpha \log' H_v^{(2)}(k_1 a) = \beta N \log' H_v^{(1)}(k_2 a). \quad (83)$$

The approximate positions of the root v_λ of (83) in the lower complex v plane are

$$v_\lambda = k_1 a + t_\lambda \left(\frac{k_1 a}{6} \right)^{\frac{1}{3}} e^{-i(\pi/3)} + \dots, \quad (84)$$

where t_λ is a number determined by

$$e^{i(\pi/3)} \left(\frac{6}{k_1 a} \right)^{\frac{1}{3}} \frac{A'(t_\lambda)}{A(t_\lambda)} \cong -i \frac{\beta}{\alpha} \left[1 - \left(\frac{v_\lambda}{k_2 a} \right)^2 \right]^{\frac{1}{2}} N, \quad (85)$$

and $A(z)$ is the Airy function (see Appendix).

After closing the contour in the lower complex v plane (Fig. 4), evaluating the residues, and neglecting

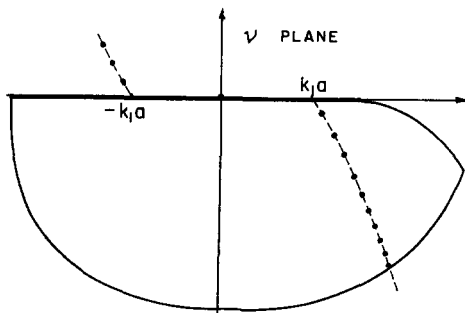


FIG. 4. This figure shows schematically the positions of the poles of the integrand of (82), as well as the path of integration in the v plane.

terms of $O[(k_1 a)^{-1}(N^2 - 1)^{-\frac{1}{2}}]$, we obtain

$$\begin{aligned}
 u_{s0r} &\cong -u_i^G + \sum_{\lambda} \left(\frac{-1}{4\pi} \right) \mathcal{D}_{11\lambda}^2 \left(\frac{a}{rr_0 \sin \theta} \right)^{\frac{1}{2}} \\
 &\quad \times [(r^2 - a^2)(r_0^2 - a^2)]^{-\frac{1}{2}} (1 + e^{-i2\pi v_\lambda})^{-1} \\
 &\quad \times [e^{-i v_\lambda \theta} + e^{-i v_\lambda (2\pi - \theta) + i(\pi/2)}] \\
 &\quad \times \exp \left\{ -i k_1 [(r^2 - a^2)^{\frac{1}{2}} + (r_0^2 - a^2)^{\frac{1}{2}}] \right. \\
 &\quad \left. + i v_\lambda \left(\cos^{-1} \frac{a}{r} + \cos^{-1} \frac{a}{r_0} \right) \right\} \\
 &\quad \text{in the shadow region,} \\
 &\cong u_{s0r}^G + \sum_{\lambda} \left(\frac{1}{4\pi i} \right) \mathcal{D}_{11\lambda}^2 \left(\frac{a}{rr_0 \sin \theta} \right)^{\frac{1}{2}} \\
 &\quad \times [(r^2 - a^2)(r_0^2 - a^2)]^{-\frac{1}{2}} (1 + e^{-i2\pi v_\lambda})^{-1} \\
 &\quad \times [e^{-i v_\lambda (2\pi - \theta)} + e^{-i v_\lambda (2\pi + \theta) + i(\pi/2)}] \\
 &\quad \times \exp \left\{ -i k_1 [(r^2 - a^2)^{\frac{1}{2}} + (r_0^2 - a^2)^{\frac{1}{2}}] \right. \\
 &\quad \left. + i v_\lambda \left(\cos^{-1} \frac{a}{r} + \cos^{-1} \frac{a}{r_0} \right) \right\} \\
 &\quad \text{in the lit region,} \quad (86)
 \end{aligned}$$

where

$$\begin{aligned}
 \mathcal{D}_{11\lambda}^2 &= \pi \alpha^2 \beta^{-2} (N^2 - 1)^{-1} (36 k_1 a)^{-\frac{1}{2}} \\
 &\quad \times (2\pi/k_1)^{\frac{1}{2}} [A(t_\lambda)]^{-2} e^{-i(5\pi/12)} \quad (87)
 \end{aligned}$$

and $\mathcal{D}_{11\lambda}$ is the diffraction coefficient in medium 1.

Similarly, u_{s0p} has poles of $(p + 1)$ order whose positions can also be determined by (83). After evaluating them asymptotically, the expression of a typical term of u_{s0p} is

$$\begin{aligned}
 u_{s0p}^D &\cong \sum_{\lambda} \left(\frac{-1}{4\pi} \right) \mathcal{D}_{11\lambda}^2 \mathcal{D}_{12} \mathcal{D}_{21} \left(\frac{a}{rr_0 \sin \theta} \right)^{\frac{1}{2}} \\
 &\quad \times [(r^2 - a^2)(r_0^2 - a^2)]^{-\frac{1}{2}} \\
 &\quad \times \left[\omega_p + (p - 1)(\mathcal{D}_{12} \mathcal{D}_{21}) \frac{\omega_p^2}{2!} \right. \\
 &\quad \left. + \frac{(p - 1)(p - 2)}{2!} (\mathcal{D}_{12} \mathcal{D}_{21})^2 \frac{\omega_p^3}{3!} \right. \\
 &\quad \left. + \dots + \frac{(p - 1)!}{(p - 1 - \mu)! \mu!} (\mathcal{D}_{12} \mathcal{D}_{21})^\mu \frac{\omega_p^\mu}{(\mu + 1)!} \right. \\
 &\quad \left. + \dots + (\mathcal{D}_{12} \mathcal{D}_{21})^{p-1} \frac{\omega_p^p}{p!} \right] \\
 &\quad \times \exp \left\{ -i k_1 [(r^2 - a^2)^{\frac{1}{2}} + (r_0^2 - a^2)^{\frac{1}{2}}] \right. \\
 &\quad \left. - i v_\lambda \omega_p - i 2 p k_2 a \cos \zeta_{0c} + i p \pi / 2 \right\}, \quad (88)
 \end{aligned}$$

where the angular distance is

$$\omega_p = 2\pi l + \theta - \cos^{-1} \frac{a}{r} - \cos^{-1} \frac{a}{r_0} - 2p \cos^{-1} \frac{1}{N}, \quad (89)$$

the angle of critical incidence is

$$\zeta_{oc} = \sin^{-1}(1/N), \quad (90)$$

and the product of the diffraction coefficients from medium 1 to medium 2 and from medium 2 to medium 1 is

$$\mathcal{D}_{12}\mathcal{D}_{21} = 2(\alpha/\beta)(N^2 - 1)^{-\frac{1}{2}}. \quad (91)$$

Except for the difference between the amplitude factors in the two-dimensional and three-dimensional cases, the expressions of

$$u_{s0r}^D = (u_{s0r} - u_{s0r}^G) \quad \text{and} \quad u_{s0pl}^D$$

are exactly the same as the case of transparent circular cylinder,⁶ including the diffraction coefficients. Hence their detailed physical interpretations are given in Ref. 6. These expressions can also be obtained by the geometrical theory of diffraction.^{4,5}

V. HIGHER-ORDER SOLUTIONS IN THE EXTERIOR OF A LARGE ROUGH SPHERE

Upon applying the addition theorem for Legendre polynomials³ to (28) and assuming that the integration and the summation are interchangeable, we obtain

$$u_{sj} = \frac{1}{4\pi k_1} \int_0^{2\pi} \int_0^\pi \left\{ \sum_{n=0}^{\infty} (2n+1) \Delta_n^{-1} [\alpha j_n(k_2 a) S_{2j}(\theta', \varphi') - \beta k_{2j}'(k_2 a) S_{2j-1}(\theta', \varphi')] \times P_n(\cos \Omega) h_n^{(2)}(k_1 r) \right\} \sin \theta' d\theta' d\varphi', \quad j = 1, 2, 3, \dots, \quad (92)$$

where

$$\cos \Omega = \cos \theta \cos \theta' + \sin \theta \sin \theta' \cos(\varphi - \varphi'), \quad (93)$$

with Ω being the angle between \mathbf{r} and the vector passing through the point (a, θ', φ') on a spherical surface (Fig. 5).

By means of the Poisson summation formula, (92) can be rewritten as

$$u_{sj} = \frac{1}{2\pi k_1} \left(\frac{a}{r}\right)^{\frac{1}{2}} \int_0^{2\pi} \int_0^\pi \left[\sum_{l=-\infty}^{\infty} e^{i\pi l} \int_0^\infty \psi_{vj} H_v^{(2)}(k_1 r) \times P_{v-\frac{1}{2}}(\cos \Omega) e^{-i\pi l v} dv \right] \sin \theta' d\theta' d\varphi', \quad j = 1, 2, 3, \dots, \quad (94)$$

where

$$\psi_{vj} = \frac{\alpha J_v(k_2 a) S_{2j} - \beta k_{2j}'(k_2 a) S_{2j-1}}{\alpha H_v^{(2)}(k_1 a) J_v(k_2 a) - \beta N H_v^{(2)}(k_1 a) J_v'(k_2 a)}, \quad j = 1, 2, 3, \dots, \quad (95)$$

To exhibit the physical interpretation of (94) we again expand ψ_{vj} into a geometric series;

$$\psi_{vj} = \frac{-k_1 S_{2j-1}}{H_v^{(2)}(k_1 a)} \left[R_{11j} + T_{12j} T_{21} \sum_{p=1}^{\infty} R_{22}^{p-1} \mathcal{J} \mathcal{C}_2^p \right], \quad j = 1, 2, 3, \dots, \quad (96)$$

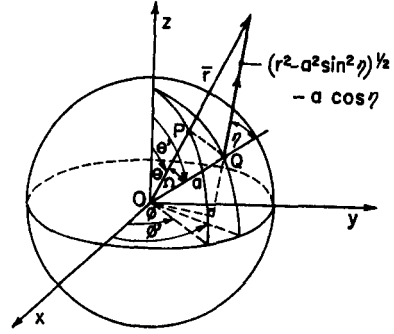


FIG. 5. The geometry of the angle Ω is shown.

where

$$R_{11j} = -[\alpha k_1^{-1} S_{2j-1}^{-1} S_{2j} - \beta N \log' H_v^{(1)}(k_2 a)] \times [\alpha \log' H_v^{(2)}(k_1 a) - \beta N \log' H_v^{(1)}(k_2 a)]^{-1} \quad (97)$$

and

$$T_{12j} = \alpha^{-1}(1 + R_{11j}), \quad j = 1, 2, 3, \dots, \quad (98)$$

The form (96) is introduced because the first term in the integrand of (94) represents mainly the wave radiated directly from the equivalent source distributions S_{2j} and S_{2j-1} [at a point (a, θ', φ') on the spherical surface] to the field point \mathbf{r} . The p th term in the sum represents mainly a wave radiated from the source distributions S_{2j} and S_{2j-1} [at a point (a, θ', φ') on the spherical surface] into the sphere, reflected $p-1$ times internally from the interface, having passed p times through the sphere, and finally transmitted out to the field point \mathbf{r} in the surrounding medium. This interpretation is borne out when the various terms are expanded asymptotically for short wavelengths and evaluated properly.

A. Geometric Optics Wave

The criterion for the proper identification of the geometric optics wave from the asymptotic evaluation of (94) is the existence of real saddle points $v_{j\rho}$ ($\rho = 1, 2, \dots$) such that $0 < v_{j\rho} < k_1 a$. For convenience we write (94) as

$$u_{sj} = u_{sjr} + \sum_{p=1}^{\infty} u_{sjp}, \quad j = 1, 2, 3, \dots, \quad (99)$$

where

$$u_{sjr} = (2\pi k_1)^{-1} (a/r)^{\frac{1}{2}} \int_0^{2\pi} \int_0^\pi \left(\sum_{l=-\infty}^{\infty} e^{i\pi l} \sum_{\delta=1}^2 \mathcal{J}_{sjrl}^{(\delta)} \right) \times \sin \theta' d\theta' d\varphi', \quad (100)$$

$$u_{sjp} = (2\pi k_1)^{-1} (a/r)^{\frac{1}{2}} \int_0^{2\pi} \int_0^\pi \left(\sum_{l=-\infty}^{\infty} e^{i\pi l} \sum_{\delta=1}^2 \mathcal{J}_{sjpl}^{(\delta)} \right) \times \sin \theta' d\theta' d\varphi', \quad (101)$$

$$\mathcal{J}_{sjrl}^{(\delta)} = -k_1 S_{2j-1} \int_0^\infty R_{11j} [H_v^{(2)}(k_1 a)]^{-1} \times H_v^{(2)}(k_1 r) \mathcal{Q}_{v-\frac{1}{2}}^{(\delta)}(\cos \Omega) e^{-i2\pi l v} dv, \quad (102)$$

and

$$J_{s_j p l}^{(\delta)} = -k_1 S_{2j-1} \int_0^\infty T_{12j} T_{21} R_{22}^{2j-1} \mathcal{H}_2^{(2)} [H_v^{(2)}(k_1 a)]^{-1} \times H_v^{(2)}(k_1 r) Q_{v-\frac{1}{2}}^{(\delta)}(\cos \Omega) e^{-i2\pi l v} v dv. \quad (103)$$

To find the contribution of the waves radiated directly from the equivalent sources S_{2j} and S_{2j-1} [as a point (a, θ', φ') on the spherical surface] to the field point \mathbf{r} , we have to evaluate $J_{s_j r l}^{(\delta)}$ asymptotically by the saddle-point method. After using the proper asymptotic forms, it is found that only $J_{s_j r 0}^{(1)}$ has a proper real saddle point when (a, θ', φ') is in the lit region of \mathbf{r} . The saddle-point equation is

$$\cos^{-1} \frac{v_j}{k_1 r} - \cos^{-1} \frac{v_j}{k_1 a} = \Omega. \quad (104)$$

Its solution can be expressed in terms of physical entities as

$$v_j = k_1 a \sin \eta, \quad (105)$$

and (104) becomes

$$\eta - \sin^{-1} \left(\frac{a}{r} \sin \eta \right) = \Omega, \quad (106)$$

where η is the equivalent angle of incidence or reflection (Fig. 6). Since the end-point contribution of $J_{s_j r l}^{(\delta)}$ is asymptotically small in comparison with the saddle-point contribution, the major contribution of $u_{s_j r}$ is

$$u_{s_j r}^\alpha \cong -k_1 a^2 (2\pi)^{-1} \iint_{\mathcal{L}} S_{2j-1} \mathcal{R}_{11j} \cos \eta (\sin \eta)^{\frac{1}{2}} \{[(r^2 - a^2 \sin^2 \eta)^{\frac{1}{2}} - a \cos \eta] r \sin \Omega\}^{-\frac{1}{2}} \times \exp \left\{ -ik_1 [(r^2 - a^2 \sin^2 \eta)^{\frac{1}{2}} - a \cos \eta] + i \frac{\pi}{2} \right\} \times \sin \theta' d\theta' d\varphi', \quad (107)$$

where the equivalent Fresnel reflection coefficient in

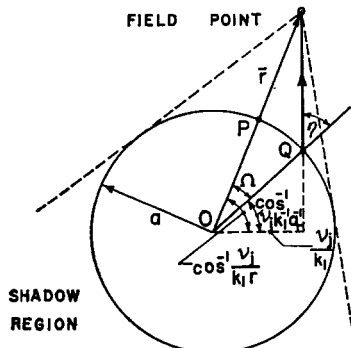


FIG. 6. The geometric ray from the equivalent source at a point Q on the surface to the field point is shown, and so is the geometrical relation between Ω and η .

medium 1 is

$$\mathcal{R}_{11j} = \frac{-i\alpha k_1^{-1} S_{2j-1}^{-1} S_{2j} - \beta N \cos \zeta}{\alpha \cos \eta + \beta N \cos \zeta}, \quad (108)$$

ζ is the equivalent angle of refraction such that

$$\sin \eta = N \sin \zeta, \quad (109)$$

and \mathcal{L} means that the integration is carried out only on the portion of spherical surface in the lit region of \mathbf{r} .

$u_{s_j r}^\alpha$ represents mainly the total contribution from the waves radiated directly from the equivalent sources S_{2j} and S_{2j-1} everywhere on the portion of the spherical surface which is in the lit region of \mathbf{r} to the field point \mathbf{r} (see Fig. 6).

To find the contribution of the waves radiated from the equivalent sources S_{2j} and S_{2j-1} at a point (a, θ', φ') on the spherical surface into the sphere, reflected $p - 1$ times internally from the interface, having passed p times through the sphere, and finally transmitted out to the field point \mathbf{r} in the surrounding medium, we have to evaluate $J_{s_j p l}^{(\delta)}$ asymptotically by the saddle-point method. After using the proper asymptotic forms, it is found that only $J_{s_j p l}^{(1)}$ may have proper saddle points, and their corresponding saddle-point equations are

$$\cos^{-1} \frac{v_{j\rho}}{kr_1} + 2p \cos^{-1} \frac{v_{j\rho}}{Nk_1 a} - \cos^{-1} \frac{v_{j\rho}}{k_1 a} = 2\pi l + \Omega, \quad \rho = 1, 2, \dots \quad (110)$$

Again the relations between $v_{j\rho}$ and η_ρ are given by

$$v_{j\rho} = k_1 a \sin \eta_\rho, \quad \rho = 1, 2, \dots \quad (111)$$

Equation (110) then gives the correct geometric relations:

$$\eta_\rho + p(\pi - 2\zeta_\rho) - \sin^{-1} \left(\frac{a}{r} \sin \eta_\rho \right) = 2\pi l + \Omega, \quad \rho = 1, 2, \dots, \quad (112)$$

where the l 's have to be properly chosen such that (112) is identical with the geometric relations between Ω and η_ρ (see Fig. 7). A typical saddle-point

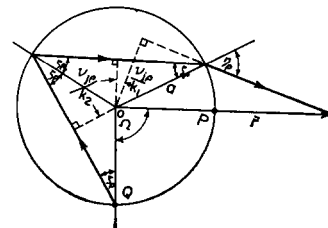


FIG. 7. The geometric ray from the equivalent source at Q , reflected once internally from the interface, having passed twice through the sphere, and transmitted out to the field point in the surrounding medium, is shown, and so is the geometrical relation between Ω and η_ρ .

contribution of $u_{s_j p}$ is

$$\begin{aligned}
 u_{s_j p}^{(\rho)} &\cong -k_1 a^2 (2\pi)^{-1} \int_{\mathcal{L}_{p\rho}} S_{2j-1} \mathcal{T}_{12j}^{(\rho)} \mathcal{T}_{21}^{(\rho)} \mathcal{R}_{22}^{(\rho)p-1} \\
 &\times \cos \eta_\rho (\sin \eta_\rho)^{\frac{1}{2}} \{ [r^2 - a^2 \sin^2 \eta_\rho]^{\frac{1}{2}} - a \cos \eta_\rho \\
 &- 2p(r^2 - a^2 \sin^2 \eta_\rho)^{\frac{1}{2}} (N \cos \zeta_\rho)^{-1} \cos \eta_\rho \} r \sin \Omega \}^{\frac{1}{2}} \\
 &\times \exp \left\{ -ik_1 [(r^2 - a^2 \sin^2 \eta_\rho)^{\frac{1}{2}} - a \cos \eta_\rho] \right. \\
 &\left. - i2pk_2 a \cos \zeta_\rho + i(p+1) \frac{\pi}{2} \right\} \sin \theta' d\theta' d\varphi', \quad (113)
 \end{aligned}$$

where $\mathcal{L}_{p\rho}$ means that the integration is carried only in the region in which S_{2j-1} and S_{2j} give contributions to $u_{s_j p}^{(\rho)}$ and

$$\begin{aligned}
 \mathcal{T}_{12j}^{(\rho)} &= \alpha^{-1} (1 + \mathcal{R}_{11j}^{(\rho)}) \\
 &= \alpha^{-1} \left(1 + \frac{-i\alpha k_1^{-1} S_{2j-1} S_{2j} - \beta N \cos \zeta_\rho}{\alpha \cos \eta_\rho + \beta N \cos \zeta_\rho} \right). \quad (114)
 \end{aligned}$$

Hence $\sum_p u_{s_j p}^{(\rho)}$ represents mainly the total contribution from the waves radiated from the equivalent sources S_{2j} and S_{2j-1} everywhere on the portions of the spherical surface into the sphere, reflected $p-1$ times internally from the interface, having passed p times through the sphere, and finally transmitted out to the field point \mathbf{r} (see Fig. 7).

B. Diffracted Wave

Similar to the zeroth-order solution, any integral of (102) and (103) which does not have proper real saddle points may still be evaluated by the method of residues. In order to make the asymptotic evaluation, we must first extend the integrals along the positive real ν axis to the entire real ν axis. Hence we find that

$$\begin{aligned}
 u_{s_j r} &= u_{s_j r}^G - \frac{i}{2\pi} \left(\frac{a}{r} \right)^{\frac{1}{2}} \left\{ \int_{\mathcal{L}'} \int S_{2j-1} \sum_{l=0}^{\infty} e^{i\pi l} \int_{-\infty}^{\infty} R_{11j} \right. \\
 &\times [H_\nu^{(2)}(k_1 a)]^{-1} H_\nu^{(2)}(k_1 r) P_{\nu-\frac{1}{2}}(-\cos \Omega) \\
 &\times e^{-i2\pi(i+\frac{1}{2})\nu} d\nu \sin \theta' d\theta' d\varphi' \\
 &+ \int_{\mathcal{L}'} \int S_{2j-1} \sum_{l=0}^{\infty} e^{i\pi l} \int_{-\infty}^{\infty} R_{11j} [H_\nu^{(2)}(k_1 a)]^{-1} \\
 &\times H_\nu^{(1)}(k_1 r) [Q_{\nu-\frac{1}{2}}^{(2)}(-\cos \Omega) + Q_{\nu-\frac{1}{2}}^{(1)}(-\cos \Omega) \\
 &\times e^{-i2\pi(\nu-\frac{1}{2})}] e^{-i2\pi(l+\frac{1}{2})\nu} d\nu \sin \theta' d\theta' d\varphi' \left. \right\}, \quad (115)
 \end{aligned}$$

where \mathcal{L}' means that the integration is carried out only on the portion of spherical surface in the shadow region of \mathbf{r} .

Since $u_{s_j r}$ has exactly the same singularities as those of $u_{s_0 r}$, after closing the integration path in the lower

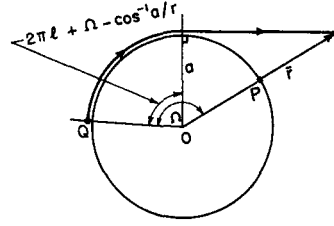


FIG. 8. The diffracted ray from the equivalent source at Q , travelling along the interface on the side of medium 1 and then leaving the surface tangentially to the field point, is shown.

complex ν plane and evaluating the residues, we obtain

$$\begin{aligned}
 u_{s_j r} &\sim u_{s_j r}^G + \sum_{\lambda} \frac{a}{2\beta^2(N^2 - 1)} [A(t_\lambda)]^{-1} (6/k_1 a)^{-\frac{1}{2}} \\
 &\times e^{-i\frac{\pi}{2}} (r^2 - a^2)^{-\frac{1}{2}} (1 + e^{-i\pi\nu\lambda})^{-1} \\
 &\times e^{-ik_1(r^2 - a^2)^{\frac{1}{2}} + i\nu\lambda \cos^{-1}(a/r)} \\
 &\times \left\{ \int_{\mathcal{L}'} \int [\alpha^2 k_1^{-1} S_{2j} - i\alpha\beta S_{2j-1} (N^2 - 1)^{\frac{1}{2}}] \right. \\
 &\times (r \sin \Omega)^{-\frac{1}{2}} \\
 &\times [e^{-i\nu\lambda\Omega} + e^{-i\nu\lambda(2\pi - \Omega) + i\frac{\pi}{2}}] \\
 &\times \sin \theta' d\theta' d\varphi' \left. \right\} + \left\{ \int_{\mathcal{L}'} \int [\alpha^2 k_1^{-1} S_{2j} - i\alpha\beta S_{2j-1} \right. \\
 &\times (N^2 - 1)^{\frac{1}{2}}] (r \sin \Omega)^{-\frac{1}{2}} \\
 &\times [e^{-i\nu\lambda(2\pi - \Omega)} + e^{-i\nu\lambda(2\pi + \Omega) + i\frac{\pi}{2}}] \\
 &\times \sin \theta' d\theta' d\varphi' \left. \right\}. \quad (116)
 \end{aligned}$$

$u_{s_j r}^D = (u_{s_j r} - u_{s_j r}^G)$ represents mainly the total contribution from the waves which are radiated from the equivalent sources S_{2j} and S_{2j-1} , which travel along the interface an angular distance [e.g., $2\pi l + \Omega - \cos^{-1}(a/r)$ on the side of medium 1], which decay exponentially, and then leave the surface tangentially toward the field point \mathbf{r} (see Fig. 8). Finally, $u_{s_j p}$ can be evaluated in a manner similar to the case of $u_{s_0 p}$; hence their analogous physical interpretation can be obtained.

VI. PERTURBATION SOLUTION IN THE EXTERIOR OF A LARGE CONVEX ROUGH SCATTERER

The zeroth-order solution is the solution for the special case of a large convex smooth and deterministic scatterer for which an asymptotic solution in the exterior of the scatterer can be constructed from the geometrical theory of diffraction.⁴⁻⁷ The higher-order solutions can then be constructed by observing the following relation between the zeroth-order solution of spherical scatterer and its higher-order solutions and then using it.

Upon comparing (66) to (107) and (80) to (113), we find that

$$u_{s_{jr}}^G = 2i \iint_{\Gamma} k_1 a^2 S_{2j-1} \cos \eta \left[\begin{array}{c} \lim_{r_0 \rightarrow a, \theta \rightarrow \Omega} (u_{s_{0r}}^G) \\ \mathcal{R}_{11} \rightarrow \mathcal{R}_{11j} \end{array} \right] \times \sin \theta' d\theta' d\varphi' \quad (117)$$

and

$$u_{s_{jp}}^{(\rho)} = 2i \iint_{\Gamma_{pp}} k_1 a^2 S_{2j-1} \cos \eta \times \left[\begin{array}{c} \lim_{r_0 \rightarrow a, \theta \rightarrow \Omega} (u_{s_{0p}}^{(\rho)}) \\ \mathcal{R}_{11}^{(\rho)} \rightarrow \mathcal{R}_{11j}^{(\rho)}, \mathcal{C}_{11}^{(\rho)} \rightarrow \mathcal{C}_{12j}^{(\rho)} \end{array} \right] \sin \theta' d\theta' d\varphi'. \quad (118)$$

The combination of (117) and (118) implies that the geometric optics wave of u_{s_j} is given by

$$u_{s_j}^G = 2i \iint k_1 a^2 S_{2j-1} \cos \eta \left[\begin{array}{c} \lim_{r_0 \rightarrow a, \theta \rightarrow \Omega} (u_{s_0}^G) \\ \mathcal{R}_{11} \rightarrow \mathcal{R}_{11j}, \mathcal{C}_{12} \rightarrow \mathcal{C}_{12j} \end{array} \right] \times \sin \theta' d\theta' d\varphi', \quad j = 1, 2, 3, \dots \quad (119)$$

To generalize (119) to the case of diffraction of waves by a large, convex, rough object, we replace each "a" in (119) by the local radius of curvature \tilde{a} . Hence the geometric optics wave of the higher-order solution is given by

$$U_{s_j}^G = 2i \iint k_1 \tilde{a}^2 \tilde{S}_{2j-1} \cos \tilde{\eta} \times \left[\begin{array}{c} \lim_{r_0 \rightarrow \tilde{a}} (U_{s_0}^G) \\ \mathcal{R}_{11} \rightarrow \tilde{\mathcal{R}}_{11j}, \mathcal{C}_{12} \rightarrow \tilde{\mathcal{C}}_{12j} \end{array} \right] d\Lambda, \quad j = 1, 2, 3, \dots, \quad (120)$$

where $d\Lambda$ is the increment of solid angle and the tilde denotes the local behavior.

The diffracted wave of the higher-order solution can be obtained from the diffracted wave of the zeroth-order solution in a similar way. Because of the complexity we only give the formula for $U_{s_{jr\lambda}}^D$:

$$U_{s_{jr\lambda}}^D = - \left(\frac{2k_1}{\pi} \right)^{\frac{1}{2}} e^{i(\pi/12)} \iint \left\{ \lim_{r_0 \rightarrow \tilde{a}} \tilde{a} A(\tilde{r}_\lambda) (6k_1^2 \tilde{a}^2)^{\frac{1}{2}} \cdot [k_1^{-1} \tilde{S}_{2j} - i\alpha^{-1} \beta \tilde{S}_{2j-1} (N^2 - 1)^{\frac{1}{2}}] \times (r_0^2 - \tilde{a}^2)^{\frac{1}{2}} U_{s_{0r\lambda}}^D \right\} d\Lambda, \quad j = 1, 2, 3, \dots \quad (121)$$

VII. MEAN WAVEFUNCTION AND MEAN INTENSITY

The mean values of the wavefunction and intensity in the exterior of the scatterer are, respectively,

$$\langle U_i + U_s \rangle = U_i + U_{s_0} + \sum_{j=1}^{\infty} \epsilon^j \langle U_{s_j} \rangle, \quad (122)$$

and

$$\langle |U_s|^2 \rangle = |U_{s_0}|^2 + \epsilon 2 \operatorname{Re} (\bar{U}_{s_0} \langle U_{s_1} \rangle) + \epsilon^2 [2 \operatorname{Re} (\bar{U}_{s_0} \langle U_{s_2} \rangle) + \langle |U_{s_1}|^2 \rangle] + O(\epsilon^3). \quad (123)$$

Upon examining our problem, it is found that the statistics of this problem are contained only in the equivalent sources \tilde{S}_{2j-1} and \tilde{S}_{2j} ($j = 1, 2, 3, \dots$). Therefore, from (120) and (121), we have

$$\langle U_{s_j}^G \rangle = 2i \iint k_1 \tilde{a}^2 \cos \tilde{\eta} \left\langle \begin{array}{c} \lim_{r_0 \rightarrow \tilde{a}} (U_{s_0}^G) \tilde{S}_{2j-1} \\ \mathcal{R}_{11} \rightarrow \tilde{\mathcal{R}}_{11j}, \mathcal{C}_{12} \rightarrow \tilde{\mathcal{C}}_{12j} \end{array} \right\rangle d\Lambda, \quad j = 1, 2, 3, \dots, \quad (124)$$

and

$$\langle U_{s_{jr\lambda}}^D \rangle = \left(\frac{2k_1}{\pi} \right)^{\frac{1}{2}} e^{-i(11\pi/12)} \iint \left\{ \lim_{r_0 \rightarrow \tilde{a}} \tilde{a} (6k_1^2 \tilde{a}^2)^{\frac{1}{2}} A(\tilde{r}_\lambda) [k_1^{-1} \langle \tilde{S}_{2j} \rangle - i\alpha^{-1} \beta \langle \tilde{S}_{2j-1} \rangle (N^2 - 1)^{\frac{1}{2}}] (r_0^2 - \tilde{a}^2)^{\frac{1}{2}} U_{s_{0r\lambda}}^D \right\} d\Lambda, \quad j = 1, 2, 3, \dots \quad (125)$$

Similarly, the asymptotic expansion of $\langle |U_s|^2 \rangle$ can be obtained from (120), (121), (123), (124), and (125)—except that it is very complex. Since \tilde{S}_{2j-1} and \tilde{S}_{2j} contain various combinations of f , $(\partial f / \partial \xi_2)$, and $(\partial f / \partial \xi_3)$ in sums and products, once all of the moments of these various combinations are known, then the mean values of the scattered wavefunction and the scattered intensity are determined.

ACKNOWLEDGMENT

This research was partially supported by grant No. GU1533 of the National Science Foundation. The author wishes to thank Professor J. B. Keller for the invitation to Courant Institute of Mathematical Sciences, New York University, where part of this research was performed. The author also wishes to acknowledge the Danish Statens teknisk-videnskabelige Fond for the fellowship at the Technical University of Denmark, where the last stage of this research was completed. Finally, the author wants to express his thanks to Professor Eric Hansen for his kind interest in this work.

APPENDIX

Debye asymptotic forms for large argument and index, with $\nu < z$, are

$$H_\nu^{(1)}(z) \sim \left(\frac{\pi}{2} z \sin \tau \right)^{-\frac{1}{2}} \times \exp \left[iz(\sin \tau - \tau \cos \tau) - i \frac{\pi}{4} \right],$$

$$H_\nu^{(2)}(z) \sim \left(\frac{\pi}{2} z \sin \tau \right)^{-\frac{1}{2}} \times \exp \left[-iz(\sin \tau - \tau \cos \tau) + i \frac{\pi}{4} \right],$$

$$H_\nu^{(1)'}(z) \sim i \sin \tau H_\nu^{(1)}(z),$$

and

$$H_v^{(2)'}(z) \sim -i \sin \tau H_v^{(2)}(z),$$

where

$$\tau = \cos^{-1}(\nu/z).$$

For argument and index both large, and $\nu \simeq z$, we have the Airy function representations for the Hankel functions as

$$H_v^{(2)}(z) \sim (2/\pi)ZA(t),$$

$$H_v^{(2)'}(z) \sim -(2/\pi)Z^2A'(t),$$

$$(\partial/\partial\nu)H_v^{(2)}(z) \sim (2/\pi)Z^2A'(t),$$

and

$$(\partial/\partial\nu)H_v^{(2)'}(z) \sim (2/\pi)(Z^3/3)tA(t),$$

where

$$Z = (6/z)^{1/3} e^{i(\pi/3)}, \quad t = Z(\nu - z).$$

In connection with Legendre functions, $Q_v^{(1)}(\cos \theta)$ and $Q_v^{(2)}(\cos \theta)$ are defined as

$$Q_v^{(1)}(\cos \theta) = \frac{1}{2} \left[P_\nu(\cos \theta) + i \frac{2}{\pi} Q_\nu(\cos \theta) \right]$$

and

$$Q_v^{(2)}(\cos \theta) = \frac{1}{2} \left[P_\nu(\cos \theta) - i \frac{2}{\pi} Q_\nu(\cos \theta) \right].$$

Their asymptotic forms are

$$\left. \begin{aligned} Q_v^{(1)}(\cos \theta) &\sim (2\pi\nu \sin \theta)^{-\frac{1}{2}} \\ &\times \exp \left[-i(\nu + \frac{1}{2})\theta + i \frac{\pi}{4} \right] \\ \text{and} \\ Q_v^{(2)}(\cos \theta) &\sim (2\pi\nu \sin \theta)^{-\frac{1}{2}} \\ &\times \exp \left[i(\nu + \frac{1}{2})\theta - i \frac{\pi}{4} \right] \end{aligned} \right\} 0 < \theta < \pi.$$

Local Characterization of Singularities in General Relativity*†

ROBERT GEROCH‡

Palmer Physical Laboratory, Princeton, New Jersey

(Received 5 July 1967)

We formulate a new approach to singularities: their local description. Given any incomplete space-time M , we define a topological space, the " g boundary," whose points consist of equivalence classes of incomplete geodesics of M . The points of the g boundary may be thought of as the "singular points" of M . Local properties of the singularity may now be described in a well-defined way in terms of local properties of the g boundary. For example, the notions: "dimensionality of a singularity," "past and future of a singular point," "neighborhood of a singular point," "spacelike or timelike character of a singularity," and "metric structure of a singularity" may all be expressed as properties of the g boundary. Two applications of the g boundary outside of the realm of singularities are discussed: (1) In the case in which the space-time M is extendable (for example, Taub space), the g boundary is shown to be that regular 3-surface across which M may be extended [in this case, the Misner boundary between Taub and Newman-Unti-Tamburino (NUT) space]. (2) With a slight modification of the definitions, the g boundary of an asymptotically simple space-time is shown to be Penrose's surface at "conformal infinity." The application of the g boundary technique to singularities is illustrated with a number of examples. The g -boundary structure of one particular example leads to our consideration of non-Hausdorff space-times.

I. INTRODUCTION

Singularities in general relativity are normally defined in terms of geodesic completeness.¹ Yet when one hears the word "singularity," he imagines that the temperature, the mass density, or perhaps the

curvature becomes infinite. Thus, one might ask of a space-time: "Does this or that physical quantity grow without bound in the vicinity of the singularity?" This very question presupposes that there is some "singular point" at which the singularity resides, and that we may ask questions about the vicinity of that point. But geodesic incompleteness does not at all provide us with the "singular points" or their "neighborhoods" that we should like to have in order to phrase our physical questions. We present an approach to the problem of bridging the gap between geodesic incompleteness on the one hand and our physical notions about singularities on the other. We

* This work was carried out under a National Science Foundation Graduate Fellowship.

† Submitted in partial fulfillment of the requirements of the degree of Doctor of Philosophy, Princeton University, Princeton, New Jersey. A brief summary of this work will appear in the *Proceedings of the Battelle Summer Recontres, Seattle, Washington, 1968* (W. A. Benjamin, Inc., New York) (to be published).

‡ Present address: Department of Mathematics, Birkbeck College, London, England.

¹ See, for example, C. W. Misner, *J. Math. Phys.* **4**, 924 (1963); R. Penrose, *Phys. Rev. Letters* **14**, 57 (1965); R. Geroch (unpublished).

and

$$H_v^{(2)'}(z) \sim -i \sin \tau H_v^{(2)}(z),$$

where

$$\tau = \cos^{-1}(\nu/z).$$

For argument and index both large, and $\nu \simeq z$, we have the Airy function representations for the Hankel functions as

$$H_v^{(2)}(z) \sim (2/\pi)ZA(t),$$

$$H_v^{(2)'}(z) \sim -(2/\pi)Z^2A'(t),$$

$$(\partial/\partial\nu)H_v^{(2)}(z) \sim (2/\pi)Z^2A'(t),$$

and

$$(\partial/\partial\nu)H_v^{(2)'}(z) \sim (2/\pi)(Z^3/3)tA(t),$$

where

$$Z = (6/z)^{1/3} e^{i(\pi/3)}, \quad t = Z(\nu - z).$$

In connection with Legendre functions, $Q_v^{(1)}(\cos \theta)$ and $Q_v^{(2)}(\cos \theta)$ are defined as

$$Q_v^{(1)}(\cos \theta) = \frac{1}{2} \left[P_\nu(\cos \theta) + i \frac{2}{\pi} Q_\nu(\cos \theta) \right]$$

and

$$Q_v^{(2)}(\cos \theta) = \frac{1}{2} \left[P_\nu(\cos \theta) - i \frac{2}{\pi} Q_\nu(\cos \theta) \right].$$

Their asymptotic forms are

$$\left. \begin{aligned} Q_v^{(1)}(\cos \theta) &\sim (2\pi\nu \sin \theta)^{-\frac{1}{2}} \\ &\times \exp \left[-i(\nu + \frac{1}{2})\theta + i \frac{\pi}{4} \right] \\ \text{and} \\ Q_v^{(2)}(\cos \theta) &\sim (2\pi\nu \sin \theta)^{-\frac{1}{2}} \\ &\times \exp \left[i(\nu + \frac{1}{2})\theta - i \frac{\pi}{4} \right] \end{aligned} \right\} 0 < \theta < \pi.$$

Local Characterization of Singularities in General Relativity*†

ROBERT GEROCH‡

Palmer Physical Laboratory, Princeton, New Jersey

(Received 5 July 1967)

We formulate a new approach to singularities: their local description. Given any incomplete space-time M , we define a topological space, the " g boundary," whose points consist of equivalence classes of incomplete geodesics of M . The points of the g boundary may be thought of as the "singular points" of M . Local properties of the singularity may now be described in a well-defined way in terms of local properties of the g boundary. For example, the notions: "dimensionality of a singularity," "past and future of a singular point," "neighborhood of a singular point," "spacelike or timelike character of a singularity," and "metric structure of a singularity" may all be expressed as properties of the g boundary. Two applications of the g boundary outside of the realm of singularities are discussed: (1) In the case in which the space-time M is extendable (for example, Taub space), the g boundary is shown to be that regular 3-surface across which M may be extended [in this case, the Misner boundary between Taub and Newman-Unti-Tamburino (NUT) space]. (2) With a slight modification of the definitions, the g boundary of an asymptotically simple space-time is shown to be Penrose's surface at "conformal infinity." The application of the g boundary technique to singularities is illustrated with a number of examples. The g -boundary structure of one particular example leads to our consideration of non-Hausdorff space-times.

I. INTRODUCTION

Singularities in general relativity are normally defined in terms of geodesic completeness.¹ Yet when one hears the word "singularity," he imagines that the temperature, the mass density, or perhaps the

curvature becomes infinite. Thus, one might ask of a space-time: "Does this or that physical quantity grow without bound in the vicinity of the singularity?" This very question presupposes that there is some "singular point" at which the singularity resides, and that we may ask questions about the vicinity of that point. But geodesic incompleteness does not at all provide us with the "singular points" or their "neighborhoods" that we should like to have in order to phrase our physical questions. We present an approach to the problem of bridging the gap between geodesic incompleteness on the one hand and our physical notions about singularities on the other. We

* This work was carried out under a National Science Foundation Graduate Fellowship.

† Submitted in partial fulfillment of the requirements of the degree of Doctor of Philosophy, Princeton University, Princeton, New Jersey. A brief summary of this work will appear in the *Proceedings of the Battelle Summer Recontres, Seattle, Washington, 1968* (W. A. Benjamin, Inc., New York) (to be published).

‡ Present address: Department of Mathematics, Birkbeck College, London, England.

¹ See, for example, C. W. Misner, *J. Math. Phys.* **4**, 924 (1963); R. Penrose, *Phys. Rev. Letters* **14**, 57 (1965); R. Geroch (unpublished).

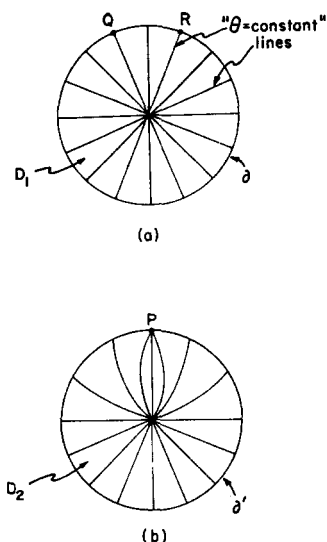


FIG. 1. An example illustrating that the boundary of a given manifold is not in general unique.

describe a construction by which, given any incomplete space-time, one may define the requisite singular points and describe their local properties.²

It may appear at first sight that in all practical cases one may, by inspection, define the appropriate singular points and attach them to form a boundary for the space-time. Thus, so it might appear, no complicated prescriptions are necessary. But appearances are notoriously coordinate dependent. We first give a simple example to show that the boundary of a manifold³ is not uniquely determined by the manifold itself. Consider an open disk, D_1 , in polar coordinates: $r < 1$, $\theta \in (0, 2\pi)$ [Fig. 1(a)]. A boundary ∂ may be appended to D_1 by inspection: the circle $r = 1$. We now construct a different boundary for the same manifold. Map our original disk, D_1 , diffeomorphically onto another identical disk, D_2 , where the mapping is shown in Fig. 1(b) by drawing the image in D_2 of the lines $\theta = \text{const}$ in D_1 . (The exact form of the mapping is not important here.) But now the natural boundary for D_2 [the circle $r' = 1$ in Fig. 1(b)] forms also a boundary, ∂' , for D_1 . With the boundary ∂ , each of the lines $\theta = \text{const}$ in D_1 end at a different boundary point, whereas with ∂' a number of these lines end at the same boundary point, namely P . Suppose that D_1 were a (two-dimensional) space-time and that the mass density became infinite near the point Q while remaining finite near R in Fig. 1(a). Assume we had "mistakenly" put the boundary ∂'

² A construction along similar lines has been discussed by S. W. Hawking, "Singularities and the Geometry of Space-Time." (Unpublished essay submitted for the Adams Prize, Cambridge University, December, 1966.)

³ For the definition of a manifold with boundary, see J. R. Munkres, *Elementary Differential Topology* (Princeton University Press, Princeton, N.J., 1963), p. 3.

on D_1 , and then asked whether or not the mass density becomes infinite near the "boundary point" P . It is clear that our question could not be meaningfully answered, for the distinct points Q and R of ∂ are collected into the single point P of ∂' . This example illustrates our point: the manifold alone will not tell one how a boundary of "singular points" should be affixed—one requires some additional structure such as a metric, and then a prescription for fixing the boundary given that metric.

In Sec. II we discuss the topological properties of singularities. Given any geodesically incomplete space-time M , we define equivalence classes of its incomplete geodesics, two geodesics being in the same equivalence class if they approach each other (in a sense to be defined precisely) as the geodesics approach the singularity. The space of equivalence classes is called the "g boundary." The g boundary, endowed with a topology, is to represent the "singular points" of the space-time. A prescription is given for attaching the g boundary to the original space-time M to form the "space-time with g boundary" \bar{M} . In the topological space \bar{M} the notions of a singular point and its neighborhood are well defined.

In Sec. III we consider the problem of assigning a causal, differentiable, and metric structure to the g boundary. The future and past of each g-boundary point are defined. The notions of future and past are used to define "spacelike" and "timelike" g boundaries. Conditions are given under which one may define a differentiable and metric structure on ∂ . The definitions in Sec. III are useful in particular in the analysis of the (highly symmetrical) cosmological models in general relativity.

In Sec. IV we discuss two applications of the g boundary outside of the field of singularities:

1. A prescription is given to determine whether or not a given space-time is extendable and, if so, to carry out an extension.

2. A technique is given to construct the surface at conformal infinity defined by Penrose⁴ for the study of asymptotically flat spacetimes.

Finally, in Sec. V, we consider a number of examples. The g boundaries of several well-known singular solutions of Einstein's equations are investigated.

II. CONSTRUCTION OF THE g BOUNDARY

The method we shall describe for identifying geodesics to form the g boundary is basically simple.

⁴ R. Penrose, in *Relativity, Groups, and Topology*, C. DeWitt and B. DeWitt, Eds. (Gordon and Breach, Science Publishers, Inc., New York, 1964), p. 565; R. Penrose, Proc. Roy. Soc. (London) **A284**, 159 (1965).

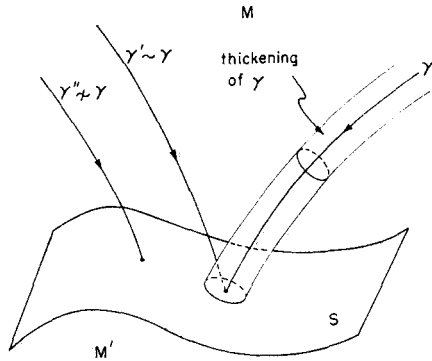


FIG. 2. The use of the "thickening" to restore the boundary of a given incomplete space-time.

Complications arise only in the course of generalizing the technique to accommodate more pathological space-times. We consider first, therefore, a special case. Let V be a geodesically complete space-time. Let S be a 3-submanifold of V which divides V into two disjoint parts, M and M' .⁵ We are to be given only the 4-geometry M . To what extent can we now recover the surface S ?

Consider the collection of all incomplete geodesics in M . These are precisely the geodesics of M which, when extended in V , pass through S . We wish to find a criterion, using only the geometry of M , to group these geodesics into equivalence classes, two geodesics being in the same equivalence class if they pass, when extended in V , through the same point of the surface S . Each equivalence class will then represent one point of S . Let γ be any incomplete geodesic in M . We may describe γ by its initial conditions: any one point on γ and the tangent vector to γ at this point. Consider now the family of geodesics which results from small variations in these initial conditions. This family traces out a four-dimensional tube in M which we call a *thickening* of γ (Fig. 2). If γ' is another incomplete geodesic in M , we write $\gamma' \sim \gamma$ if γ' enters and remains inside every thickening of γ . In this example, \sim is an equivalence relation. The equivalence classes reproduce the points of the surface S , i.e., $\gamma' \sim \gamma$ if and only if γ and γ' have the same end point on S .

In a general incomplete space-time, it will not be possible to find a regular 3-surface on which each incomplete geodesic terminates. We may, however, generalize the above construction so that it will be applicable to any space-time. The resulting equivalence classes of geodesics will then constitute the g boundary.

Let M be any space-time. Denote by G the collection

of all nonzero vectors in M , that is, a point of G represents a vector at *some given point* of M . Since the eight-dimensional manifold G is the tangent bundle⁶ of M with the zero vectors omitted, we call G the *reduced tangent bundle* of M . Each element of G may be written in the form (P, ξ^α) , where ξ^α is a vector at the point P of M . The point (P, ξ^α) of G uniquely determines that geodesic in M which begins at P and has its initial direction and affine parameter determined by the equation

$$\left. \frac{dx^\alpha}{d\lambda} \right|_P = \xi^\alpha.$$

It is convenient to use the word "geodesic" here to mean a parameterized curve having the properties:

1. The curve has one end point, and has been extended as far as possible in some direction from that end point.
2. The curve is a geodesic with the given parameter as affine parameter.
3. The affine parameter vanishes at the end point and is positive elsewhere on the curve.

There is a one-to-one correspondence between the points of the reduced tangent bundle and the geodesics of M . Define a scalar field φ on the manifold G as the total affine length of the corresponding geodesic in M . Thus, φ is infinite if and only if the geodesic in question is complete. Denote by G_I ("I" for "incomplete") that subset of G on which φ is finite.

Define the 9-manifold $H = G \times (0, \infty)$, and the following two subsets of H :

$$H_+ \equiv \{(P, \xi^\alpha, a) \in H \mid \varphi(P, \xi^\alpha) > a\},$$

$$H_0 \equiv \{(P, \xi^\alpha, a) \in H \mid \varphi(P, \xi^\alpha) = a\}.$$

There is a natural map⁷ $\Psi: H_+ \rightarrow M$ defined as follows: Given a point (P, ξ^α, a) of H_+ , let $\Psi(P, \xi^\alpha, a)$ be that point of M which results from traversing an affine distance a along the geodesic (P, ξ^α) .

We next define a topology on G_I . Let O be any open set of M . We associate with O a subset $S(O)$ of G_I as follows:

$$S(O) \equiv \{(P, \xi^\alpha) \in G_I \mid \text{there exists an open set } U \text{ in } H \text{ containing the point } (P, \xi^\alpha, \varphi(P, \xi^\alpha)) \text{ of } H_0 \text{ such that } \Psi(U \cap H_+) \subset O\}.$$

Thus, if O is a sufficiently small open set of M whose closure is compact, then $S(O)$ is the empty set. One

⁵ To simplify the discussion at this stage, let us also assume that every geodesic beginning tangent to S first enters M' rather than M .

⁶ See, for example, N. Steenrod, *The Topology of Fibre Bundles* (Princeton University Press, Princeton, N.J., 1951), p. 5.

⁷ See, for example, N. J. Hicks, *Notes on Differential Geometry* (D. Van Nostrand Co., Inc., Princeton, N.J., 1965), p. 131.

can easily verify that, given any two open sets O_1 and O_2 of M , the subsets $S(O_1)$ and $S(O_2)$ of G_I obey

$$S(O_1) \cap S(O_2) = S(O_1 \cap O_2).$$

It follows⁸ that the collection of sets $S(O)$, where O ranges over all open sets of M , serves as a basis for the open sets of a topology on G_I .

We use this topology on G_I to form equivalence classes of the elements of G_I as follows. If α and β are two elements of G_I , write $\alpha \approx \beta$ if every open set in G_I containing α also contains β , and every open set containing β also contains α . The relation \approx is an equivalence relation. The collection of equivalence classes will be denoted by ∂ and referred to as the *g boundary* ("g" for geodesic). The topology we have defined on G_I induces⁹ a topology on ∂ .¹⁰

In fact, there are several different ways to form the equivalence classes which define the *g boundary*. The identification described above is the weakest in the sense that points of G_I are identified only if they cannot be distinguished topologically, i.e., if they always appear in the same open sets.¹¹ The resulting *g boundary* satisfies the separation axiom¹² T_0 . A stronger identification is as follows: write $\alpha \approx \beta$ if for every continuous function f on G_I , $f(\alpha) = f(\beta)$.¹³ The *g boundary* as defined by this identification scheme satisfies the separation axiom T_3 . It is also possible to construct equivalence relations on G_I so that the topology of the *g boundary* satisfies separation axioms T_1 or T_2 .¹⁴ In many examples it makes no difference which type of identification is selected; the resulting *g boundaries* are identical. However, there are cases in which the *g boundary* depends on the choice of identification scheme. In the last example of

Sec. V, all the interesting properties of the *g boundary* would be destroyed if the T_3 identification were employed.

So far the *g boundary* exists only as an abstract topological space. We now attach ∂ to the space-time M . Define

$$\bar{M} \equiv M \cup \partial,$$

the disjoint union. A subset (O, U) of \bar{M} , where O is an open set of M and U is an open set of ∂ , will be called open in \bar{M} if $S(O) \supset U$. One can verify that the intersection of two such open sets is open. These open sets on \bar{M} are a basis for a topology on \bar{M} . The restriction of the topology on \bar{M} to M reproduces the manifold topology on M , while the restriction to ∂ reproduces the topology of the *g boundary*. \bar{M} will be called the *space-time with g boundary*. If $U = S(O)$, (O, U) will be called a *full open set* of \bar{M} .

On applying these definitions to the example at the beginning of this section, we see that the *g boundary* is just the surface S . The space-time with *g boundary* \bar{M} is a manifold with boundary: the space-time M with the boundary surface S attached.

The definitions of the *g boundary* and of the space-time with *g boundary* are the basic concepts of this paper. Given any incomplete space-time M , we may define the two topological spaces, ∂ and \bar{M} . The points of ∂ then represent the singular points, their neighborhoods in \bar{M} the neighborhoods of the singular points. Thus, the statement "the mass density becomes infinite at the singularity" may be expressed in precise terms as follows: "The point e of the *g boundary* has the property that for every ρ_0 , however large, there is an open neighborhood (O, U) of e in \bar{M} such that $\rho > \rho_0$ in O ." Many other notions about singularities may be similarly described. Note in particular that the *g boundary*, being a topological space, may be assigned a dimension¹⁵ if ∂ turns out to be a separable metric space.

It should be emphasized that though we have given a prescription that succeeds in principle in defining the *g boundary* for any space-time, in practice the construction may be quite difficult. In particular, one needs a considerable amount of information about the geodesics. In many of the known solutions of Einstein's equations, integrals for the geodesics may be obtained up to quadratures, and these suffice for the construction. But the known solutions are highly symmetric, and one cannot expect such a reduction to be possible in general. Since only the "asymptotic properties" of the geodesics are relevant, one might hope to find a method to carry out the *g-boundary*

⁸ See, for example, J. G. Hocking and G. S. Young, *Topology* (Addison-Wesley Publishing Co., Reading, Mass., 1961), p. 6.

⁹ That is, if we write π for the mapping $G_I \rightarrow \partial$, then a subset U of ∂ is open if $\pi^{-1}(U)$ is open in G_I .

¹⁰ There is another way to obtain a topology on ∂ . The manifold topology on G induces a topology on G_I , which is in general different from the one we have defined. We could define the topology on the *g boundary* to be that induced by this new topology on G_I .

¹¹ See W. J. Pervin, *General Topology* (Academic Press Inc., New York, 1964), p. 155.

¹² See Ref. 8, p. 37.

¹³ E. Chech, *Ann. Math.* **38**, 823 (1937); R. Vaidyanathaswamy, *Set Topology* (Chelsea Publishing Co., New York, 1960), p. 154.

¹⁴ Let S be any topological space, and let A be an equivalence relation on S . A may be represented as a subset \bar{A} of $S \times S$ in the following way. The point (α, β) of $S \times S$ is in \bar{A} if and only if α and β are identified under the equivalence relation A . The intersection of any collection of equivalence relations, defined as the intersection of the corresponding subsets of $S \times S$, gives a new equivalence relation. An equivalence relation will be said to be of type T_i ($i = 1, 2$) if the induced topology on the equivalence classes obeys separation axiom T_i . Define:

$A_1 =$ the intersection of all T_1 equivalence relations on S .

$A_2 =$ the intersection of all T_2 equivalence relations on S .

It is not difficult to show that A_1 and A_2 are, respectively, T_1 and T_2 equivalence relations on S .

¹⁵ See Ref. 8, p. 145.

construction without any detailed properties of the geodesics. Unfortunately, this hope has not yet been realized. We have seen, however, that *some* properties of the metric must be used to define a g boundary. Since the geodesics are among the simplest such properties, it is not clear where one should look for a different, more easily applied, construction which is still applicable to even the most pathological space-times.

Finally, we mention that no guarantee can be offered as to the uniqueness of the g boundary construction here. Our definitions were obtained by trying one definition, finding an example for which that definition did not perform as expected, then trying another definition, etc. It is certainly possible that a more natural definition for the g boundary exists, especially in light of the lack of an acceptable definition of a singularity.¹⁶ As one example of a possible modification, one might let the reduced tangent bundle G consist, not of all nonzero vectors in M , but only of the nonzero timelike vectors. One then forms the g boundary from equivalence classes of the incomplete timelike geodesics. Since there are examples¹⁷ of space-times which are timelike complete but neither spacelike nor null complete, the new g boundary would differ in general from the g boundary we have defined here.

III. FURTHER PROPERTIES OF THE g BOUNDARY

Use of the space-time with g boundary defined in Sec. II apparently solves the problem of describing local properties of singularities. Why, then, should one consider any further structures on the g boundary?

The first reason is that there are still several intuitive notions about singularities which cannot be expressed in terms of their topological features alone. Is the Schwarzschild singularity spacelike? Does the presence of a magnetic field in a collapsing star cause the collapse to proceed more quickly in directions parallel or perpendicular to the field lines? Which points of the Reissner-Nordström solution may be affected by information coming into the space-time from the singularity? One would like to give a meaning to the undefined notions expressed by these questions.

Secondly, the topological properties of the g boundary are not alone sufficient to begin to classify singularities. We shall discuss in Sec. V examples of space-times whose g boundaries are topologically identical, but whose singularities seem intuitively to be quite different.

A. Causal Structure

We consider first the following two problems:

1. Determine the "future" and "past" of each point of the g boundary.

2. Define "spacelike" and "timelike" g boundaries. These two aspects of the g boundary are considered together because they both involve only the causal structure of the space-time (i.e., "Can event A influence, by means of a signal, event B ?"). Kronheimer and Penrose¹⁸ have shown that the causal relations on a set can be discussed without reference to a metric, a differentiable structure, or even a topology on that set. Thus, one comes to think of the causal structure as the first and most basic relation to be specified on a collection of events, rather than a subsidiary property derived as an afterthought from the metric.

Let C be any directed curve in M . We say C has an *end point* at the point $e \in \partial$ if for every open neighborhood (O, U) of e in \bar{M} , C enters and remains in O . Of course, it may be that C has no end point in a given direction. It is easy to show, however, that if C does have an end point, then that point is unique if \bar{M} is Hausdorff. A curve in M having at least one end point in each direction is said to *connect* those points. Let e be any point of ∂ . Define the *future*¹⁹ of e :

$$I^+(e) \equiv \{Q \in \bar{M} \mid \text{there exists a timelike curve } C \text{ in } M \text{ with a future end point } Q \text{ and a past end point } e\}.$$

The *past*, $I^-(e)$, is similarly defined. Note that $I^+(e)$ and $I^-(e)$ are well defined for every g -boundary point e without restriction on the space-time.

Let us give an example to show the utility of this definition. Consider the Reissner-Nordström solution²⁰ (Fig. 3). One often hears the remark: "The event P may be influenced by information coming into the space-time from the singularity." What does that mean? Does it mean that there are incomplete past-directed timelike or null geodesics from the point P ? That is false. Does it mean that there are inextendable past-directed timelike curves from P of finite total length (i.e., the curve C in the figure)? This will not do either, for there are timelike curves of finite total length from any point of any spacetime. In fact, it is hard to give any meaning at all to the words "The singularity may influence P ." except to say that Fig. 3 leaves one with an impression of this sort. The remark $P \in I^+(\partial)$ is a precise statement of this impression.

¹⁸ E. H. Kronheimer and R. Penrose, Proc. Cambridge Phil. Soc. 63, 481 (1967).

¹⁹ Throughout this discussion, we must assume that the space-time is isochronous, at least in the region under consideration. Only under this assumption can we meaningfully distinguish future from past.

²⁰ See B. Carter, Phys. Letters 21, 423 (1966).

¹⁶ W. Kundt, Z. Physik 172, 488 (1963); L. Shepley, thesis, Princeton University, 1965; R. Geroch (unpublished).

¹⁷ W. Kundt, Ref. 16.

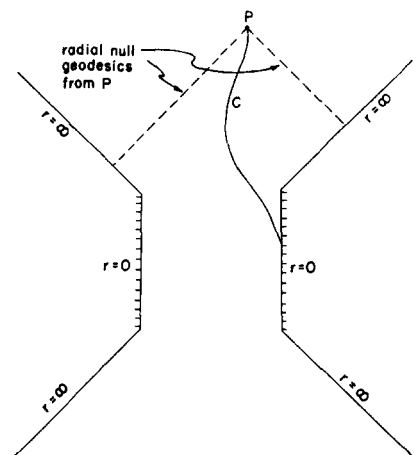


FIG. 3. The Reissner-Nordström solution. The point P is to the future of the g boundary ($r = 0$).

We next come to the question of determining whether a given g boundary should be considered spacelike, timelike, or perhaps neither. Unfortunately, a number of possible definitions come to mind, all of which yield the expected answers in common examples. We therefore present definitions of spacelike and timelike g boundaries which are fairly restrictive, so that many g boundaries will be in neither class. Perhaps a better understanding of this issue in the future will indicate the appropriate modifications or generalizations of the definitions given here.

We wish to have the spacelike or timelike character of a g boundary be a local property, i.e., a property unaffected by changes in the geometry in regions away from the g boundary. But the ordinary notion of the future and past of a point involves global considerations. For example, there are space-times having the property that the future of each point is the entire 4-geometry. We must define, therefore, the *local future* of a point. Let $e \in \partial$, and let (O, U) be a full open neighborhood of e in \bar{M} . Define

$$I^+(e; O, U) = \{\text{the future of } e \text{ in } (O, U)\},$$

and similarly for the *local past*, $I^-(e; O, U)$. We say that ∂ is *spacelike* at e if there exists a full open

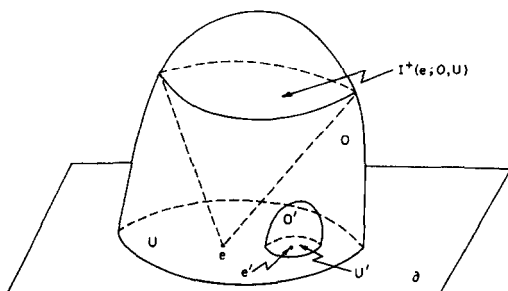


FIG. 4. The definition of a spacelike g boundary.

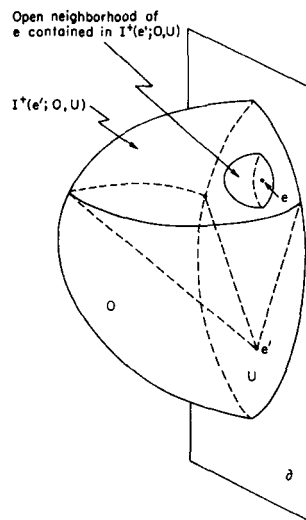


FIG. 5. The definition of a timelike g boundary. The point e'' and its local past have been suppressed.

neighborhood (O, U) of e in \bar{M} such that, for all $e' \in U$, there is a full open neighborhood (O', U') of e' with

$$[I^+(e; O, U) \cup I^-(e; O, U)] \cap (O', U') = \phi,$$

where ϕ is the empty set (Fig. 4). We say that ∂ is *timelike* at e if for every full open neighborhood (O, U) of e in \bar{M} there exist two points, e' and e'' , of U such that $I^+(e'; O, U) \cap I^-(e''; O, U)$ contains an open neighborhood of e in \bar{M} (Fig. 5).²¹

We have not defined a null g boundary because the null character of a surface is a causal property only when it obtains over a region. For example, the surface $x = t + t^3$ in Minkowski 3-space is null only on the line $x = t = 0$, yet this surface cannot be distinguished by its causal properties from an everywhere timelike surface. In fact, it is difficult to see why one would wish to characterize a surface as null unless it is also regular (in some sense). But when the g boundary is sufficiently regular, a causal definition of null is unnecessary as the metric properties we discuss next will provide the appropriate characterization.

B. Differentiable and Metric Structure

We consider next the construction of a differentiable and metric structure on the g boundary. In order to define these structures, it is necessary to impose very strong conditions on the space-time. One would not expect these conditions to hold in general in solutions

²¹ The following are examples of weaker definitions. ∂ is spacelike at $e \in \partial$ if there exists a full open neighborhood (O, U) of e such that $[I^+(e; O, U) \cup I^-(e; O, U)] \cap U = \phi$. ∂ is timelike at e if for every full open neighborhood (O, U) of e , $I^+(e; O, U) \cap U \neq \phi$ and $I^-(e; O, U) \cap U \neq \phi$.

of Einstein's equations. In particular, no differentiable structure can be defined on a region of ∂ unless that region is a topological manifold.

Let O be any open set of the space-time M , and let (O, U) be the corresponding full open set in \bar{M} . Since the points of ∂ are equivalence classes of elements of G_I , we have a natural map $\pi: G_I \rightarrow \partial$. Define, for each open set O of M , the following subset of G_I :

$$\bar{J}(O) \equiv \{(P, \xi^\alpha) \in \pi^{-1}(U) \mid \varphi(P, \xi^\alpha) = 1, \text{ and the geodesic in } M \text{ associated with the element } (P, \xi^\alpha) \text{ of } G_I \text{ lies entirely within the open set } O\}.$$

We say ∂ has a differentiable structure at the point e of ∂ if there exists a full open neighborhood (O, U) of e in \bar{M} and a subset J of $\bar{J}(O)$ such that the following four conditions are satisfied:

1. J includes almost all points of $\bar{J}(O)$.²²
2. J is a differentiable submanifold of the reduced tangent bundle G .
3. For each element e of U , $A_e \equiv \pi^{-1}(e) \cap J$ is a differentiable submanifold of J .
4. J may be written diffeomorphically as a cross product, $A \times B$, of two differentiable manifolds, where for each element of B the submanifold A is one of the A_e . We have, therefore, a one-to-one mapping $\Lambda: U \rightarrow B$. We require further that this mapping be a homeomorphism.

The differentiable structure of B and the mapping $\Lambda: U \rightarrow B$ automatically give us a differentiable structure on the open region U of ∂ .

From the definition we see that the differentiable structure, where defined, is unique. Note also that if a differentiable structure exists at a point $e \in \partial$, then the structure is defined in a neighborhood U of e . The collection of all points of ∂ at which a differentiable structure is defined forms a differentiable manifold. Returning to the example at the beginning of Sec. II, we see that the differentiable structure defined here reproduces the differentiable structure of the surface S .

A differentiable structure on the g boundary is useful only insofar as it leads to the introduction of a metric structure. We next discuss this problem.

Suppose that at the point e , ∂ has a differentiable structure characterized by the open neighborhood (O, U) of e in \bar{M} , the subset J of $\bar{J}(O)$, and the identification $J = A \times B$. Let V be a (contravariant) vector in the g boundary at the point e . We have a diffeomorphism $\Lambda: U \rightarrow B$. Therefore, the vector V in U determines a vector V' at the point $\Lambda(e)$ of B . Since

$J = A \times B$, the vector V' at $\Lambda(e)$ determines a vector field v on the submanifold A_e of J , v determined only up to the addition of an arbitrary vector tangent to the surface A_e . Make a particular choice of the arbitrary vector field tangent to A_e . The resulting vector field v on A_e will be called a representation of V . We shall describe the properties of vectors in ∂ in terms of their representations in J .

Let V be a vector at $e \in \partial$, and let v be a representation of V . The vector field v evaluated at the point q of A_e will be denoted by v_q . Since $A_e \subset J \subset G_I$, the point q of A_e is a point (P, ξ^α) of G_I . Let $x^\alpha(\lambda)$ denote the geodesic in M associated with the element (P, ξ^α) of G_I , where for each value of the affine parameter λ , $x^\alpha(\lambda)$ is a point of M . From the way we have defined J and A_e , it follows that $x(0) = P$, that λ has the range $[0, 1)$, and that the geodesic $x^\alpha(\lambda)$ is a member of that equivalence class of geodesics denoted by e . The vector v_q at the point q of J determines a Jacobi field²³ $v_q^\alpha(\lambda)$ on the geodesic $x^\alpha(\lambda)$ as follows. Let ϵ be infinitesimal, so ϵv_q is an infinitesimal vector in J at the point q . This vector may be thought of as joining the point q to a nearby point q' of J . We may write the geodesic corresponding to the point q' of J in the form $x^\alpha(\lambda) + \epsilon v_q^\alpha(\lambda)$. $v_q^\alpha(\lambda)$ is the Jacobi field in M associated with v_q . Choose a fixed number $\lambda_0 \in [0, 1)$. Then $v_q^\alpha(\lambda_0)$ is a vector at the point $x^\alpha(\lambda_0)$ of M . Parallel transport this vector along the geodesic $x^\alpha(\lambda)$ to the endpoint P of this geodesic. We thus determine a vector at P , which we denote by $\bar{v}_q^\alpha(\lambda_0)$. For each value of λ_0 in the interval $[0, 1)$, $\bar{v}_q^\alpha(\lambda_0)$ is a vector at P . We shall call $\bar{v}_q^\alpha(\lambda_0)$ the vector function associated with v_q .

Intuitively speaking, our construction is as follows. An infinitesimal vector V at the point e of ∂ may be considered as joining e to a nearby point e' . Let γ be any geodesic in M with endpoint e . (The freedom in choosing γ corresponds to the freedom in choosing the point q on A_e .) The vector V determines a Jacobi field on γ only up to the addition of an arbitrary Jacobi field which joins γ to another geodesic which ends at e . [The freedom in choosing the Jacobi field corresponds to the freedom to add to any representation of V an arbitrary vector tangent to A_e (Fig. 6)]. The limit of the norm of the associated vector function, if independent of the point q and of the representation chosen, will now give us the norm of the vector V in ∂ .

Let T_q be the tangent space to the manifold J at the point $q \in J$. A stratification of T_q is a sequence $T_q^1, T_q^2, \dots, T_q^m$ of subspaces of the vector space T_q

²² Our permitting some points of $\bar{J}(O)$ to lie outside of J is not a minor formality. With "almost all" replaced by "all," the g boundary of the region $x > r^2$ in two-dimensional Minkowski space could not be given a differentiable structure.

²³ See, for example, J. Milnor, *Morse Theory* (Princeton University Press, Princeton, N.J., 1963), p. 77. A Jacobi field is a solution of the equation of geodesic deviation.

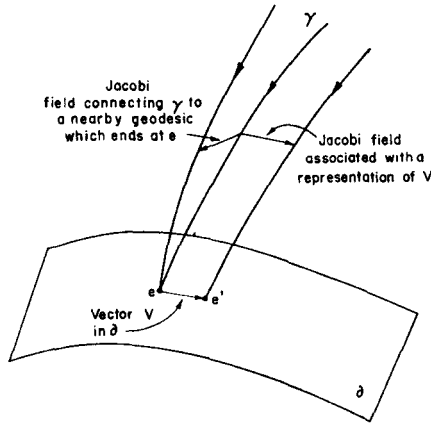


FIG. 6. The representation of vectors in the g boundary as Jacobi fields in the space-time M .

such that:

1. For each i , $T_q^i \supset T_q^{i-1}$ and $T_q^i \neq T_q^{i-1}$.
2. $T_q^m = T_q$.
3. There is a collection of positive functions $f^i(\lambda)$ ($i = 1, 2, \dots, m$), with $\lim_{\lambda \rightarrow 1} f^i/f^{i+1} = 0$, having the property $T_q^i = \{v_q \in T_q \mid \lim_{\lambda \rightarrow 1} \bar{v}_q^{(k)}(\lambda)/f^i(\lambda) \text{ exists (it may be zero)}\}$, where $\bar{v}_q^{(k)}(\lambda)$ is the vector function associated with v_q .

There always exists at least one stratification,²⁴ and in general there will exist many stratifications of T_q . We see that a stratification of T_q is a separation of T_q into vector subspaces according to the behavior of the corresponding vector functions. Those elements of T_q whose vector functions approach zero most rapidly as $\lambda \rightarrow 1$ are in the subspace T_q^1 , while those elements whose vector functions approach infinity most rapidly are in T_q^m but not in T_q^{m-1} .

We return now to the manifold J . As before, define $A_e \equiv \pi^{-1}(e) \cap J$. A global stratification of the surface A_e is a stratification (T_q^i, f^i) at each point q of A_e such that:

1. For each i , the subspace T_q^i is a differentiable function of q .
2. The $f^i(\lambda)$ can be chosen to be independent of q .
3. There is an integer p , $1 \leq p \leq m$, such that $T_q^p = T_{A_e}$, where T_{A_e} is the subspace of T_e tangent to the surface A_e at q .
4. If v_q is any representation of a vector V in ∂ ,

²⁴ Choose any collection of functions $f^i(\lambda)$ such that

$$\lim_{\lambda \rightarrow 1} f^i/f^{i+1} = 0.$$

Define the T_q^i as in property 3 above. By choosing f^m to go to infinity sufficiently quickly as $\lambda \rightarrow 1$, we may have $T_q^m = T_q$. We have now satisfied all the properties of a stratification except possibly $T_q^i \neq T_q^{i-1}$. If any two consecutive subspaces T_q^i are equal, say T_q^2 and T_q^3 , we simply drop the function f^3 from our list. Then T_q^2 is no longer among the T_q^i . Repeating this procedure until all the T_q^i are distinct, we obtain a stratification of T_q .

and if $v_q \in T_q^i$ at any one point q of A_e , then $v_q \in T_q^i$ at all points q of A_e .

Whereas a stratification is a separation of the tangent space T_q at a single point q of A_e , a global stratification is a stratification at every point of A_e . Condition 3 states that the behavior of the vector function is independent of the representation chosen. Conditions 2 and 4 state that the behavior of the vector function is independent of the point q on A_e .

A global stratification (T_q^i, f^i) of A_e will be called a metric stratification if there is a basis, $v_q^{(1)}, v_q^{(2)}, \dots, v_q^{(n)}$, for T_q at each point q of A_e , and a set of m integers, $1 \leq a_1 < a_2 < \dots < a_m = n$ (n is the dimension of the manifold J) such that:

1. The vectors $v_q^{(k)}$ for $k \leq a_i$ span T_q^i .
2. For each $k > a_p$ (p is the integer defined in condition 3 of the definition of the global stratification), $v_q^{(k)}$ is a representation of a vector in ∂ .
3. If k is an integer in the interval $a_{i-1} < k \leq a_i$, then $\lim_{\lambda \rightarrow 1} \bar{v}_q^{(k)\alpha}(\lambda)/f^i(\lambda)$ exists and is nonzero.
4. If k and j are integers in the intervals $a_{i-1} < k \leq a_i, a_{i-1} < j \leq a_i$, then

$$\lim_{\lambda \rightarrow 1} g_{\alpha\beta} \bar{v}_q^{(k)\alpha}(\lambda) \bar{v}_q^{(j)\beta}(\lambda) / [f^i(\lambda)]^2 \quad (i \text{ not summed})$$

is independent of the point q .

It is not difficult to show that if a metric stratification exists at all, then it is unique. Further, any global stratification can be obtained by merely combining several of the subspaces T_q^i of a metric stratification.

We shall be concerned here exclusively with the metric stratification because the g boundary of every space-time we have analyzed has turned out to have either no differentiable structure at all, or else a metric stratification. We could have proceeded directly to the metric stratification without introducing all the intermediate concepts. We have not done this, however, in order to emphasize the existence of an almost infinite variety of stratifications, i.e., of types of metric structures of the g boundary. If one is to take the metric properties of the g boundary seriously, it will be necessary to design the stratification to fit the g boundary under study. The basic idea is straightforward and appealing. The norm of a vector on the g boundary is given by the limit of the norm of its associated vector function. Our mathematical expression of this idea gives the appearance of being arbitrary and complicated. Perhaps there is some simple way to discuss this great variety of metric properties of the g boundary.

Suppose now that the g boundary ∂ of a space-time M admits a metric stratification. With each contravariant vector V in ∂ at e , we may associate an integer,

the *index*, defined as the smallest i such that $v_q \in T_q^i$, where v_q is a representation of V . From condition 3 of the definition of a global stratification, it follows that the index is independent of the representation chosen. We may define a metric on each subspace T_q^i by the limit in condition 4 of the definition of a metric stratification. In the metric of T_q^i , every vector has zero inner product with a vector whose index is less than i . In particular, the metric of T_q^i , restricted to the subspace T_q^{i-1} of T_q^i , is zero. If the limit $f^i(\lambda)$ is zero as $\lambda \rightarrow 1$, then each vector with index i is associated with a vector function which approaches zero as $\lambda \rightarrow 1$. In this case the metric in T_q^i , while providing a metric with which vectors with index i may be compared, must be multiplied by zero to obtain the metric in ∂ . [The limit in condition 4 above is nonzero only because $f^i(\lambda) \rightarrow 0$ as $\lambda \rightarrow 1$. Without the factors f^i , the limit itself would be zero.] Similarly, if $f^i(\lambda) \rightarrow \infty$ as $\lambda \rightarrow 1$, then the metric on T_q^i must be multiplied by "infinity" to obtain the metric on ∂ . Finally, if $f^i(\lambda)$ approaches a finite value as $\lambda \rightarrow 1$ [in which case we may, without loss of generality, set $f^i(\lambda) = 1$], then the metric on T_q^i defines the metric on ∂ . Thus, the metric on the g boundary is much richer than the ordinary notion of a metric. One admits components having the form of functions of the coordinates on ∂ multiplied by "zero" or "infinity."

As an example, consider the Schwarzschild solution. On each incomplete geodesic, the coordinates θ , φ , and t approach finite values as $r \rightarrow 0$. It turns out that two geodesics are in the same equivalence class if and only if their limiting values of θ , φ , and t are identical. Thus, the equivalence classes, i.e., the points of the g boundary, may be labeled by these three coordinates. These are also differentiable coordinates on ∂ . There are three subspaces T_q^i in the Schwarzschild solution:

T_q^1 : Jacobi fields connecting neighboring geodesics with the same limiting values of θ , φ , and t . $f^1(\lambda) \rightarrow 0$ as $\lambda \rightarrow 1$.

T_q^2 : Jacobi fields connecting neighboring geodesics with the same limiting value of t . $f^2(\lambda) \rightarrow 0$ as $\lambda \rightarrow 1$, but more slowly than $f^1(\lambda)$.

T_q^3 : All Jacobi fields. $f^3(\lambda) \rightarrow \infty$ as $\lambda \rightarrow 1$. Note that $T_q^3 \supset T_q^2 \supset T_q^1$. The integer p , defined in condition 3 of the definition of a global stratification, is one. Nonzero vectors in ∂ of the form $a d\theta + b d\varphi$ have index 2. The metric in T_q^2 is $d\theta^2 + \sin^2 \theta d\varphi^2$. Vectors in ∂ of the form $a d\theta + b d\varphi + c dt$ with $c \neq 0$ have index 3. The metric in T_q^3 is dt^2 . We may summarize this information about the metric structure of the g boundary of the Schwarzschild solution by writing

$$ds^2 = \infty dt^2 + 0(d\theta^2 + \sin^2 \theta d\varphi^2).$$

This is just the result one would have expected by examining the form of the Schwarzschild metric for $r > 0$.

Further examples will be considered in Sec. V.

IV. TWO APPLICATIONS

We shall discuss two applications of the g boundary.

1. A prescription is given to determine whether or not a given space-time is extendable, and, if so, to carry out an extension.

2. A technique is developed to construct the surface at conformal infinity defined by Penrose⁴ for the study of asymptotically flat space-times.

The constructions described in 1 and 2 above are of interest not only for their direct applications, but also because they demonstrate that our "equivalence classes of geodesics" idea, while developed to study singularities, also has applications to other problems. One is led to feel that the g boundary has some physical content.

A. Extending Spacetimes

Let M' be a space-time with boundary ∂'_3 such that each geodesic in M' is either complete, or else strikes ∂'_3 . Let ∂'_2 , ∂'_1 , and ∂'_0 be, respectively, closed 2-, 1-, and 0-dimensional (not necessarily connected) submanifolds of M' such that ∂'_3 , ∂'_2 , ∂'_1 , and ∂'_0 are disjoint by pairs. Define $\partial' \equiv \partial'_3 \cup \partial'_2 \cup \partial'_1 \cup \partial'_0$. Since ∂' is closed, $M \equiv M' - \partial'$ is a space-time. Any space-time M constructed in this manner will be said to be *extendable*.²⁵ By *extending* M we mean restoring, *given only* M , the space-time with boundary M' .

These definitions coincide with the intuitive notion of "eliminating a coordinate singularity." For example, let M be the Schwarzschild solution for $r > 2m$. Then M' is the Schwarzschild solution for $r \geq 2m$. Given M' , we may find a coordinate system in which the metric is regular also on the surface $r = 2m$, for each point of this surface is a regular point of M' . The Kruskal coordinates²⁶ are one such coordinate system.

It should be emphasized that we are not concerned here with the more difficult problem of extending an arbitrary space-time. Consider, for example, the interior M of a closed, nowhere differentiable curve drawn in Minkowski 2-space. Here, M can be extended to include the whole of Minkowski 2-space, yet our analysis would fail to bring to light such an extension. It would be preferable, of course, to treat the most general problem: the extension of an arbitrary incomplete space-time. An approach more general than that

²⁵ To simplify the exposition, we have given an unnecessarily restrictive definition of an extendable space-time. Three important generalizations, not treated in detail here, will be mentioned later.

²⁶ M. D. Kruskal, Phys. Rev. 119, 1743 (1960).

described here could probably be developed, based on the g boundary and the use of normal coordinate systems as suggested by Szekeres.²⁷ However, since the extension of space-times is essentially a practical problem, it may be preferable to wait until pathological examples arise in practice before considering sophisticated techniques to deal with them.

We now give a prescription for extending an arbitrary extendable space-time, as these terms have been defined here. The method is straightforward. One applies the definitions of Secs. II and III to the space-time M , but examines the consequences in M' . Since each point of ∂' is a regular point of M' , the effect of each definition is transparent in M' .

Consider first the g boundary of an extendable space-time M . Let γ be an incomplete geodesic in M . In M' , γ must have an end point at some point e of ∂' . Suppose $e \in \partial'_2$. Since the ∂'_i are closed and disjoint, we may find an open neighborhood O of e in M' such that O does not intersect ∂'_3 , ∂'_1 , or ∂'_0 . It is clear that if $\tilde{\gamma}$ is a second incomplete geodesic in M , then γ and $\tilde{\gamma}$ will appear in the same equivalence class if and only if $\tilde{\gamma}$ also has an endpoint at e . That is, two incomplete geodesics in M appear in the same equivalence class if and only if they have the same endpoints on ∂' . We have, therefore, a one-to-one onto mapping $\Delta: \partial \rightarrow \partial'$. Using similar arguments, one may verify that:

1. The mapping $\Delta: \partial \rightarrow \partial'$ is a homeomorphism.
2. ∂ has a differentiable structure. With this structure, Δ is a diffeomorphism.²⁸
3. Each point of ∂ has a metric stratification. The metric on ∂ is everywhere finite. With this metric, Δ becomes an isometry.²⁹
4. The space-time with g boundary \bar{M} is homeomorphic to M' . That is, Δ may be extended to a homeomorphism $\Delta: \bar{M} \rightarrow M'$.

We have been able to restore the topological, differentiable, and metric structures of ∂' and the topological structure of M' using only the space-time M . We must now determine the differentiable structure of M' .

Define the following collection of functions on \bar{M} :

$$F \equiv \{f \mid f \text{ is a function on } \bar{M} \text{ such that the scalar field } f\Psi \text{ on } H_0 \cup H_+ \text{ is a differentiable field}\}.$$

We now define the differentiable functions on \bar{M} to be the collection F of functions. It follows from the differentiability of the exponential map³⁰ Ψ that we

thereby fix a differentiable structure on \bar{M} and that $\Delta: \bar{M} \rightarrow M'$ becomes a diffeomorphism.

Let us summarize the steps involved in extending an incomplete space-time M .

1. Define the space-time with g boundary $\bar{M} \equiv M \cup \partial$. \bar{M} will be a topological manifold,³¹ possibly with boundary.
2. Define the collection F of differentiable functions on \bar{M} as above. We thus impose a differentiable structure on \bar{M} .
3. The given metric on the subset M of \bar{M} determines a unique metric on all of \bar{M} by continuity.

We have seen that if M is in fact extendable, i.e., if M results from cutting the submanifolds ∂'_i from a space-time M' , then \bar{M} will be just M' .

But what if one of these three steps should fail? Suppose that \bar{M} is not a topological manifold, or that the functions F do not define a differentiable structure on \bar{M} . Then, we are assured, M is not extendable.

We mention three generalizations of this construction.

1. It is not necessary to require that the submanifolds ∂'_i deleted from M' be differentiable (C^∞), but only C^0 , piecewise C^2 .
2. We could also have deleted from M' a 3-submanifold S not a part of the boundary of M' . This situation would be treated in a similar way except that locally S would be two-sided, and so the g boundary of M would include each point of S twice, once for each side. It would then be necessary to identify the appropriate pairs of points to restore the surface S .
3. M' might also contain a "real singularity," i.e., incomplete geodesics which do not terminate on ∂' . It would then be necessary, in order to extend M , to isolate and ignore those points of the g boundary of M whose equivalence classes consist of incomplete geodesics which do not terminate on ∂' .

B. The Construction of Conformal Infinity

A space-time M is called *asymptotically simple*⁴ if there exists a space-time M' with boundary \mathfrak{J} , and a diffeomorphism Δ of M onto the interior of M' such that:

1. Δ maps the metric $g_{\alpha\beta}$ of M into the metric $g'_{\alpha\beta} = \Omega^2 g_{\alpha\beta}$ of the interior of M' .
2. Ω is a differentiable scalar field on M' . On the boundary \mathfrak{J} of M' , $\Omega = 0$ and $\partial_a \Omega \neq 0$.
3. Every null geodesic in M' has two end points on \mathfrak{J} .

It is known⁴ that if the Ricci tensor vanishes in a

²⁷ G. Szekeres, Publ. Math. Deb. (Hungary) 7, 285 (1960).

²⁸ One must use the fact that the exponential mapping, $\Psi: H_+ \rightarrow M$, is differentiable (Ref. 7, p. 131).

²⁹ One must use the properties of Jacobi fields (Ref. 23).

³⁰ See Ref. 7.

³¹ For the definition of a topological manifold, see Ref. 3, p. 3.

neighborhood of \mathfrak{J} , then \mathfrak{J} is a null surface. We shall assume throughout this discussion that \mathfrak{J} is null.

We direct our attention to the following question: "How is it possible to determine whether or not a given space-time is asymptotically simple, and, if so, to find M' ?" We shall show how, given the space-time M , M' may be constructed using a sequence of steps very similar to those used in extending space-times.

In fact, it will be possible to carry out the entire construction using only the conformal properties of M . In particular, we deal with only the null geodesics. We must therefore slightly modify our previous definitions. Define

$$G \equiv \{(P, \xi^\alpha) \mid P \text{ is a point of } M, \text{ and } \xi^\alpha \text{ is a nonzero null vector at } P\}.$$

G is a 7-manifold. Each point of G determines a curve in M : the directed null geodesic with tangent vector ξ^α at the point P . Define, for each open set O of M ,

$$S(O) \equiv \text{interior } \{(P, \xi^\alpha) \in G \mid \text{the null geodesic associated with } (P, \xi^\alpha) \text{ enters and remains in } O\},$$

where the interior is taken in the manifold topology of G . The $S(O)$, where O ranges over all open sets of M , form a basis for the open sets of a new topology on G . We define equivalence classes in G : $(P, \xi^\alpha) \approx (P', \xi'^\alpha)$ if every open set (in the new topology on G) containing (P, ξ^α) contains also (P', ξ'^α) , and every open set containing (P', ξ'^α) contains also (P, ξ^α) . The collection of equivalence classes, denoted by ∂ , has a topology induced on it from the topology we have defined on G . The definitions of \bar{M} , of its topology, and of the differentiable structure of ∂ are the same as in Secs. II and III.

The following properties follow immediately from the asymptotic simplicity of M :

1. Two null geodesics in M are in the same equivalence class if and only if they have the same endpoint on \mathfrak{J} . We have, therefore, a one-to-one onto mapping $\Delta: \partial \rightarrow \mathfrak{J}$.
2. $\Delta: \partial \rightarrow \mathfrak{J}$ is a homeomorphism.
3. ∂ has a differentiable structure. With this structure, $\Delta: \partial \rightarrow \mathfrak{J}$ is a diffeomorphism.
4. Δ may be extended to a homeomorphism $\Delta: \bar{M} \rightarrow M'$.

We next use the conformal structure of M to define a conformal structure on ∂ . We require the transformation properties of Jacobi fields under conformal transformations. Let $\gamma(\lambda)$ be a null geodesic in a space-time M , where λ is an affine parameter. Define the tangent vector to this geodesic, $\xi \equiv d\gamma(\lambda)/d\lambda$. Let η^α be a Jacobi field on γ which connects γ to a second, nearby, null geodesic. Consider now the conformally

transformed metric, $g'_{\alpha\beta} = \Omega^2 g_{\alpha\beta}$. Write the Jacobi field connecting the *same* two null geodesics (but now with the metric $g'_{\alpha\beta}$) in the form

$$\eta'^\alpha = \eta^\alpha + \theta \xi^\alpha. \tag{1}$$

Imposing the Jacobi equation on η'^α (using the metric $g'_{\alpha\beta}$), we obtain an expression for θ :

$$\theta = -\frac{2}{\Omega} \int^\lambda (\eta^\alpha \partial_\alpha \Omega) d\lambda$$

As before, we define a vector function $\bar{\eta}^\alpha(\lambda)$ by parallel transporting the Jacobi field η^α , evaluated at $\gamma(\lambda)$, along γ to some fixed point P on γ .

Let (P, ξ^α) be a point of G . Let $T'_{P\xi^\alpha}$ be the tangent space to G at (P, ξ^α) , and let γ be the null geodesic in M determined by (P, ξ^α) . Define the following subspace of $T'_{P\xi^\alpha}$:

$$T'_{P\xi^\alpha} \equiv \{\eta^\beta \in T'_{P\xi^\alpha} \mid \eta^\alpha \xi_\alpha = 0\}.$$

We see from Eq. (1) that the subspace $T'_{P\xi^\alpha}$ is a conformal invariant of M . $T'_{P\xi^\alpha}$ is a 6-dimensional vector space. In M' , let the end point of γ be $e \in \mathfrak{J}$. Since \mathfrak{J} is null, the normal vector k to \mathfrak{J} at e lies in \mathfrak{J} . Since \mathfrak{J} is a regular surface in M' , we see that no representation of k is contained in $T'_{P\xi^\alpha}$. However, if h is any other vector tangent to \mathfrak{J} at e , then there exists a point (P, ξ^α) of G such that (i) the null geodesic associated with (P, ξ^α) has endpoint e , and (ii) a representation of h is contained in $T'_{P\xi^\alpha}$. These statements about M' are conformally invariant, and must therefore hold also in M . We have established the following: If M is asymptotically simple, and if $e \in \partial$, then there is precisely one direction k in ∂ at e having the property that for each point (P, ξ^α) in the equivalence class e , no representation of k is contained in $T'_{P\xi^\alpha}$.

Let (P, ξ^α) be an element of G whose null geodesic is in the equivalence class e . Since $T'_{P\xi^\alpha}$ is 6-dimensional, and since no representation of k is in $T'_{P\xi^\alpha}$, every vector η^α in $T'_{P\xi^\alpha}$ can be written in the form

$$\eta^\alpha = \eta_{(1)}^\alpha + \eta_{(2)}^\alpha,$$

where $\eta_{(1)}^\alpha$ is in $T'_{P\xi^\alpha}$, and $\eta_{(2)}^\alpha$ is a representation of k . There must exist³² a function $f(\lambda)$ such that

$$\lim_{\lambda \rightarrow \infty} g_{\alpha\beta} \bar{\eta}_{(1)}^\alpha(\lambda) \bar{\eta}_{(1)}^\beta(\lambda) / [f(\lambda)]^2 \tag{2}$$

exists for all η^α and is nonzero for at least one η^α . Define the norm of η^α by the limit (2). The norm of any vector in ∂ at e is given by the norm of one of its representations.³³ We thus define a metric on ∂ at e ,

³² This statement is conformally invariant because of Eq. (1) and the definition of $T'_{P\xi^\alpha}$. Therefore, since the statement is true of M' (with $f(\lambda) = \text{const}$), it is true of M .

³³ The norm thus defined is independent of the representation chosen, as one may verify in a neighborhood of $\Delta(e)$ in M' .

this metric determined only up to a factor because of the freedom to multiply $f(\lambda)$ by an arbitrary constant factor. From Eq. (1), we see that a conformal transformation on M induces no more than a conformal transformation of the metric we have defined at the point e of ∂ . Therefore, the null geodesic $\gamma(\lambda)$, ending at the point e of ∂ , determines a unique conformal metric at e . In M' , we see that this metric is precisely the induced conformal metric on \mathcal{J} , independently of the geodesic $\gamma(\lambda)$ chosen. Therefore, the conformal metric defined on ∂ is independent of the geodesic $\gamma(\lambda)$.

We must finally define a differentiable structure on \bar{M} . Let γ be any null geodesic in M with end point $e \in \partial$. Let η^a be a Jacobi field on γ such that the limit (2) is nonzero. Define the following function on γ :

$$w \equiv g_{\alpha\beta} \eta_{(1)}^\alpha(\lambda) \eta_{(1)}^\beta(\lambda).$$

The function w is not a conformal invariant. The transformation $g'_{\alpha\beta} = \Omega^2 g_{\alpha\beta}$ induces the transformation $w' = \Omega^2 w$. There will exist a point on γ beyond which w is strictly positive.³⁴ We may therefore conformally transform M so that $w = 1$ beyond some point on γ . Let t be any affine parameter on γ (in this conformally transformed metric on M) which is zero on ∂ and positive elsewhere on γ . We shall call t a *preferred* parameter on γ . A preferred parameter is not in general defined along an entire null geodesic, but only beyond a certain point on that geodesic.

Let F denote the collection of all functions f on \bar{M} such that

1. f is a continuous function on \bar{M} .
2. f , restricted to M , is differentiable.
3. f , restricted to ∂ , is differentiable.
4. For each null geodesic $\gamma(t)$ with preferred parameter t , $f[\gamma(t)]$ is a differentiable function of t for $t \geq 0$.

Since a function f on M' is differentiable if and only if it is in the collection F , the functions F define a differentiable structure on \bar{M} .

This completes the construction. We have shown how, given any asymptotically simple space-time M , M' may be constructed. We may now reverse the argument. Define a space-time M to be asymptotically simple if our construction succeeds at every step (i.e., if \bar{M} is a topological manifold, if the functions F define a differentiable structure on \bar{M} , etc.) and if the conformal factor Ω which relates the metric of M to that of \bar{M} has the properties $\Omega = 0$ and $\partial_\alpha \Omega \neq 0$ on ∂ . Thus, we may express the asymptotic simplicity of

³⁴ In M' , η^a is a nonzero spacelike vector at $\Delta(e)$. Therefore, $w > 0$ at $\Delta(e)$. It follows that $w > 0$ on some interval of γ about $\Delta(e)$.

M in terms of only the intrinsic properties of the space-time M .

V. EXAMPLES AND CONCLUSION

It is necessary to have a certain amount of detailed information about the geodesics in a space-time in order to construct its g boundary. Unfortunately, in many of the exact solutions of Einstein's equations, the geodesic equations are too complicated to permit the derivation, in any simple way, of the structure of the g boundary. This problem represents one of the most significant disadvantages of our method for practical use. There are, however, a number of space-times with sufficient Killing vectors so that the geodesics can be obtained up to quadratures. The results of g boundary analyses of several of these solutions follow.

1. *Schwarzschild solution.* We take the metric in the form

$$ds^2 = -[1 - (2m/r)] dt^2 + [1 - (2m/r)]^{-1} dr^2 + r^2(d\theta^2 + \sin^2 \theta d\varphi^2)$$

for $0 < r < 2m$. An extension beyond $r = 2m$ requires a change of coordinates.²⁶ The g boundary consists of two disjoint cylinders, each topologically $S^2 \times R$ (one cylinder for each of the two " $r = 0$ " singularities). The space-time with g boundary \bar{M} may be obtained by allowing the coordinate r to assume the value zero. θ , φ , and t are continuous coordinates on ∂ . Thus, the g boundary is three-dimensional, not one-dimensional³⁵ (Figs. 7 and 8).

The g boundary may be given a differentiable structure. It turns out that the coordinates t , θ , and φ are differentiable coordinates on ∂ (aside from the usual inadequacy of spherical coordinates at $\varphi = 0$, $\varphi = 2\pi$, $\theta = 0$, and $\theta = \pi$). The g boundary has a metric structure (see Sec. III). The metric may be written in the shorthand form

$$ds^2 = \infty dt^2 + 0(d\theta^2 + \sin^2 \theta d\varphi^2).$$

It is meaningful, in particular, to say that collapse takes place in the directions tangent to the 2-spheres, while infinite expansion occurs in the t direction.

The g boundary is spacelike.

³⁵ The 3-dimensional character of ∂ is not so unreasonable as it might appear at first sight. Consider Minkowski space with the t axis removed. Let γ be a geodesic which intersects the t axis at the origin O . On slightly varying the initial conditions of γ , we obtain a family of geodesics which trace out a thickening of γ containing an open neighborhood of the point O . This neighborhood completely surrounds the t axis in the vicinity of O (Fig. 7). On the other hand, let γ be a geodesic in the Schwarzschild solution, and let $(\bar{r}, \bar{\varphi}, \bar{\theta})$ be the limiting values of these three coordinates on γ as $r \rightarrow 0$. If we slightly vary the initial conditions of γ , we obtain a family of geodesics which strike the singularity $r = 0$ in a small neighborhood of $(\bar{r}, \bar{\varphi}, \bar{\theta})$ (Fig. 8). If the Schwarzschild singularity were to be described as 1-dimensional, we would expect all values of θ and φ to be represented in this neighborhood.

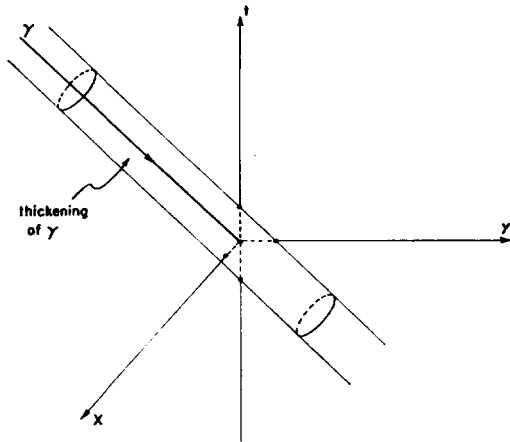


FIG. 7. The 1-dimensional character of the t axis in Minkowski 3-space.

In this case (and in all but one of the examples we shall discuss), the metric structure of the g boundary turns out to be exactly what one would have expected by looking at the metric components in the familiar coordinate system. It should be emphasized, however, that the results quoted here are not based on guesswork, but on application of the definitions given in Secs. II and III.

2. *The Reissner-Nordström solution.* Write the metric in the form³⁶

$$ds^2 = -\left(1 - \frac{2m}{r} + \frac{q^2}{r^2}\right) dt^2 + \left(1 - \frac{2m}{r} + \frac{q^2}{r^2}\right)^{-1} dr^2 + r^2(d\theta^2 + \sin^2 \theta d\varphi^2)$$

in the region $0 < r < m - (m^2 - q^2)^{1/2}$ (we assume $m > q$). Again, the g boundary consists of two disjoint cylinders, each topologically $S^2 \times R$ (compare, Hawking³⁷). The space-time with g boundary M is obtained by allowing the r coordinate to assume the value zero. The g boundary becomes, therefore, the surface $r = 0$, described by the differentiable coordinates t, θ , and φ .

The metric on the g boundary is

$$ds^2 = -\infty dt^2 + 0(d\theta^2 + \sin^2 \theta d\varphi^2)$$

The g boundary of the Reissner-Nordström solution is timelike.

There is one feature of the Reissner-Nordström solution which distinguishes it from the others. Whereas G_I is eight-dimensional for all the other space-times discussed here, G_I is only six-dimensional for the Reissner-Nordström solution. That is, only under very special circumstances does a geodesic

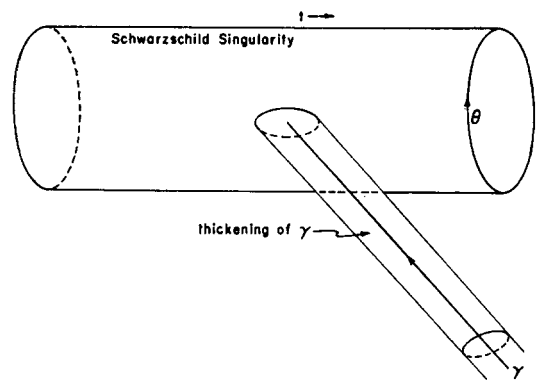


FIG. 8. The three-dimensional character of the Schwarzschild singularity. One angular variable has been suppressed.

strike the Reissner-Nordström singularity. (In fact, no timelike geodesics, and no spacelike or null geodesics with any motion around the (θ, φ) 2-spheres, reach $r = 0$.)

3. *Closed Friedmann and Tolman universes.* Take the metric in the form

$$ds^2 = -dt^2 + R^2(t)[d\chi^2 + \sin^2 \chi(d\theta^2 + \sin^2 \theta d\varphi^2)],$$

where

$$R = a(1 - \cos \eta), \quad t = a(\eta - \sin \eta)$$

for the Friedmann universe,³⁸ and

$$R = a \sin \eta, \quad t = a(\eta - \sin \eta)$$

for the Tolman universe.³⁹ The two singularities in each solution occur when $R(t) = 0$. Each singularity is diffeomorphic to the 3-sphere, with χ, θ , and φ differentiable coordinates. The metric on the g boundary is zero. The g boundary is spacelike.

4. *The spatially homogeneous, anisotropic, dust-filled solutions of Kantowski and Sachs.*⁴⁰ Here we find an example in which the g boundary would be incorrectly predicted by a quick glance at the metric in the given coordinate system. Write the metric in the form

$$ds^2 = -dt^2 + X^2(t) dr^2 + Y^2(t)(d\theta^2 + \sin^2 \theta d\varphi^2).$$

The functions $X(t)$ and $Y(t)$ are given in parametric form

$$X = 1 + (\eta + b) \tan \eta \quad t = a(\eta + \frac{1}{2} \sin 2\eta), \\ Y = a \cos^2 \eta$$

where $a > 0$ and $-\frac{1}{2}\pi \leq b < 0$ are constants.

Three different types of singularities can occur in

³⁶ See, for example, L. D. Landau and E. M. Lifshitz, *The Classical Theory of Fields* (Addison-Wesley Publ. Co., Reading, Mass., 1962), p. 380.

³⁷ R. C. Tolman, *Phys. Rev.* **37**, 1639 (1931).

⁴⁰ R. Kantowski and R. K. Sachs, *J. Math. Phys.* **7**, 443 (1966).

³⁸ J. C. Graves and D. R. Brill, *Phys. Rev.* **120**, 1507 (1960).

³⁷ Reference 2, p. 163.

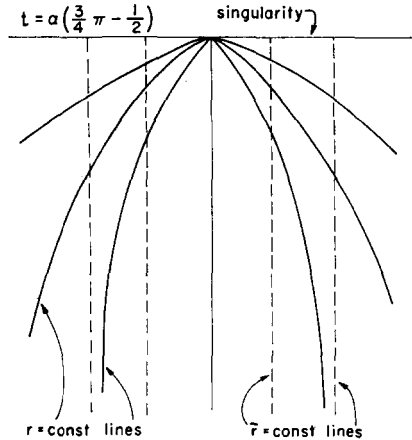


FIG. 9. The structure of the singularity of a spatially homogeneous, anisotropic cosmological model of Kantowski-Sachs.

these solutions, but we restrict consideration to one of particular interest. Set $b = 1 - (\frac{3}{4})\pi$. We restrict η to the range $(\frac{3}{4}\pi, \frac{3}{2}\pi)$. The two singularities occur at $\eta = \frac{3}{4}\pi$ (where $X = 0$) and at $\eta = \frac{3}{2}\pi$ (where $Y = 0$). We consider the singularity at $\eta = \frac{3}{4}\pi$, at which point $t = a(\frac{3}{4}\pi - \frac{1}{2})$.

The g boundary is topologically the cylinder, $S^2 \times R$, but we cannot use the coordinates r, θ, φ as continuous coordinates on ∂ . A schematic diagram of the (r, t) coordinates is shown in Fig. 9. Each point of Fig. 9 represents a 2-sphere. We see that in the Kantowski-Sachs coordinates, all of the lines $r = \text{const}$ approach the same point on the g boundary (compare Fig. 1). However, if we replace the coordinate r by the new coordinate

$$\bar{r} \equiv r[\ln(t - a(\frac{3}{4}\pi - \frac{1}{2}))]^{-1},$$

then the g boundary is properly characterized by the continuous coordinates \bar{r}, θ, φ . The lines $\bar{r} = \text{const}$ are shown in Fig. 9.

The g boundary is spacelike. The metric on the g boundary is given by

$$ds^2 = 0 d\bar{r}^2 + \frac{1}{2}a^2(d\theta^2 + \sin^2 \theta d\varphi^2).$$

Note that the metric on ∂ has signature $(0, +, +)$. Thus, it appears possible that ∂ is a regular null surface, and that the space-time may be extended through ∂ . However, the density of dust, a scalar invariant of the metric, becomes infinite as $\eta \rightarrow \frac{3}{4}\pi$, and so no extension is possible.

5. *Kasner metrics.*⁴¹ These solutions are of interest because they were used by Lifshitz and Khalatnikov⁴²

as the prototype of the general singular solution of Einstein's (source-free) equations.⁴³

Write the metric in the form

$$ds^2 = -dt^2 + t^{2p_1} dx^2 + t^{2p_2} dy^2 + t^{2p_3} dz^2.$$

The p_i are given by

$$\begin{aligned} p_1 &= -s/(1 + s + s^2), \\ p_2 &= s(1 + s)/(1 + s + s^2), \\ p_3 &= (1 + s)/(1 + s + s^2), \end{aligned}$$

where s is an arbitrary real number in the interval $0 < s < 1$. The Kasner metrics are a 1-parameter family of source-free solutions.

The g boundary of the Kasner solution has a topology and differentiable structure given by the coordinates $x, y,$ and z at $t = 0$. The metric on ∂ may be written

$$ds^2 = \infty dx^2 + 0 dy^2 + 0 dz^2.$$

The g boundary is spacelike.

6. *Taub space.*⁴⁴ In this space-time, there is a sufficient number of Killing vectors⁴⁵ to calculate the geodesics. From the information that there are two possible extensions of Taub into NUT space,⁴⁶ one can deduce immediately that the g boundary includes two disjoint, null 3-spheres. It is necessary to investigate only the equivalence classes of geodesics which are incomplete in both extensions of Taub space.

A 2-geometry has been given by Misner⁴⁶ which is not the solution of any particular field equations, but which displays several interesting features. Write the metric in the form

$$ds^2 = -\cos t dt^2 + 2 \sin t d\theta dt + \cos t d\theta^2,$$

where θ is an angular coordinate, $0 \leq \theta < 2\pi$ [Fig. 10(a)]. Let M denote the lower half ($t < \frac{1}{2}\pi$) of the cylinder in Fig. 10(a). There is a second extension of the space M . This may be seen intuitively as follows. Fix the bottom of the cylinder in Fig. 10(b), and rotate the top of the cylinder clockwise until the light cones are facing the other way, as in Fig. 10(c). Now extend the space M as shown in Fig. 10(d). The geodesic γ_1 in Fig. 10(a) passes from M into the extension of M . However, we see that, in Fig. 10(d), γ_1

⁴³ It is perhaps not widely known that the Kasner metrics represent a very specialized class of solutions. Consider the general static, cylindrically symmetric solution of Einstein's equations. That is, we consider metrics which are both (1) a special case of the Weyl-Levi-Civita static, axially symmetric solutions and (2) a special case of the Einstein-Rosen cylindrical wave solutions. By writing such a metric in cylindrical coordinates and performing a complex coordinate transformation, we recover the general Kasner metric.

⁴⁴ A. H. Taub, *Ann. Math.* **53**, 472 (1951).

⁴⁵ E. T. Newman, L. Tamburino, and T. Unti, *J. Math. Phys.* **4**, 915 (1963).

⁴⁶ C. W. Misner, *J. Math. Phys.* **4**, 924 (1963).

⁴¹ E. Kasner, *Am. J. Math.* **43**, 217 (1921).

⁴² E. M. Lifshitz and I. M. Khalatnikov, *Advan. Phys.* **12**, 185, (1963).

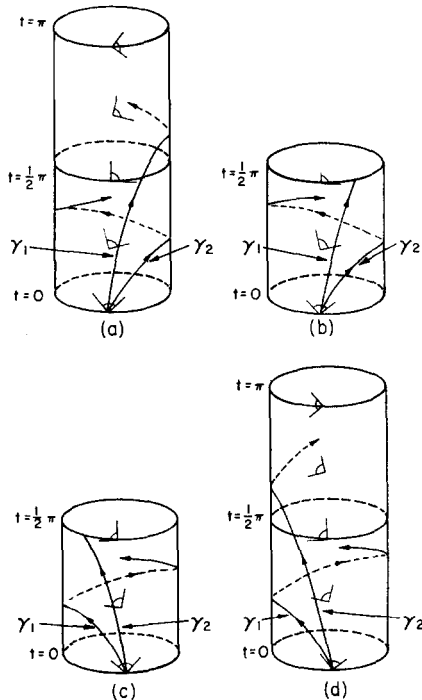


FIG. 10. The g boundary analysis of a 2-geometry which has two inequivalent extensions. Typical light cones are shown.

has become wrapped around the cylinder, asymptotically approaching the circle $t = \frac{1}{2}\pi$. The geodesic γ_2 wraps around the cylinder in Fig. 10(a), but has been “unwound” in Fig. 10(c), and so γ_2 may be extended in Fig. 10(d). From the standpoint of the space-time M , the two geodesics γ_1 and γ_2 are on equal footing. Without specifying one of the two extensions it is not possible to say which geodesic is “spiraling” around M , and which is not. M displays many of the properties of Taub space, with the two extensions playing the role of the two extensions of Taub into NUT space.

The g boundary of M may be described as the union of three sets: two circles, C and C' , and one abstract point α . A point on the circle C is the equivalence class of all those geodesics which, in Fig. 10(b), strike the circle $t = \frac{1}{2}\pi$ at one fixed point. Thus C represents the regular boundary $t = \frac{1}{2}\pi$ in Fig. 10(b). The circle C' represents the regular boundary $t = \frac{1}{2}\pi$ in Fig. 10(c). The abstract point α is the equivalence class consisting of all geodesics which spiral around the cylinder in both extensions.⁴⁷

⁴⁷ It is not difficult to see that such geodesics must exist. In Fig. 10(a), draw the collection of all timelike geodesics from some fixed point on the circle $t = 0$. Some, such as γ_1 , will pass through the circle $t = \frac{1}{2}\pi$ while others, such as γ_2 , will spiral around the cylinder on approaching $t = \frac{1}{2}\pi$. If we continuously change our geodesic from γ_1 to γ_2 , we must find one geodesic which is on the boundary between the two classes. It follows from continuity arguments that this geodesic must spiral around the cylinder in both Fig. 10(a) and 10(d).

The topology of ∂ is as follows: The open sets consist of the ordinary open sets of the circle C and the ordinary open sets of the circle C' . However, every open set containing the point α contains the whole g boundary ∂ . In this case, the g boundary is a T_0 topological space.

If we attach just the one circle C to M , a differentiable and metric structure can be placed on C , and we find the extension of Fig. 10(a). Similarly, if we attach the circle C' , the extension of Fig. 10(d) is indicated. The following idea was suggested by Brill: “Why must we carry out only one or the other extension? Why not extend M in both ways simultaneously?” The space which results was pointed out by Penrose (see Fig. 11). Cylinders A and D are copies of Fig. 10(b). Cylinders B and C are copies of the upper half of the cylinder in Fig. 10(a). Geodesic γ_1 in Fig. 11 does not spiral around cylinder A, and so passes into cylinder B. γ_1 then turns around in B but, on approaching the boundary $t = \frac{1}{2}\pi$ of B, γ_1 spirals around cylinder B. This is the signal that γ_1 must be continued into cylinder D. Geodesic γ_2 spirals around cylinder A, and is therefore extended into cylinder C. On turning around in C, γ_2 continues into cylinder D. Note that the observers on the geodesics γ_1 and γ_2 do not detect any anomalous effects on crossing the boundaries from one cylinder to another. An observer on γ_1 , for example, does not concern himself with cylinder C, for he never enters that region. In his local neighborhood, the geometry is differentiable and of the proper signature at all times.

Unfortunately, there is a bit of difficulty with regard to the geodesics in the equivalence class α . How shall these be extended? If $\gamma_3 \in \alpha$, one would like to continue γ_3 directly into cylinder D without entering

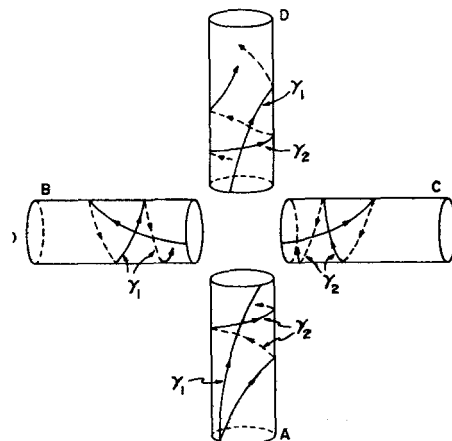


FIG. 11. Penrose's extension of the 2-geometry in Fig. 10(b). Each of the four cylinders is either the upper or the lower half of the cylinder in Fig. 10(a).

either of the cylinders B or C . However, there is a 1-parameter family of geodesics in D which strike α , and we have no way of choosing which one will become the extension of γ_3 . Of course, we could invent a consistent assignment of how each such geodesic is extended into D , but such an assignment would involve an arbitrary choice, thus destroying the unique character of the construction. Furthermore, it would not be possible to make the 2-geometry regular in a neighborhood of α . We therefore drop the point α and accept the consequence that our 2-geometry is incomplete.

Figure 11 represents a non-Hausdorff manifold. Despite this fact, there is no difficulty in writing a metric, in writing Einstein's equations or Maxwell's equations; that is, in doing all the physics to which one is accustomed. The book of Hicks⁷ deals in large part with the differential geometry of non-Hausdorff manifolds.

We have seen that the natural extension of the space-time M above is a non-Hausdorff manifold. But how are we to know in general whether or not a given space-time should be extended so as to be Hausdorff? The g boundary takes care of that. If the g boundary of a space-time consists of several disjoint parts, and if each of these parts, when individually attached to the space-time, allows an extension, then we are free to carry out all of the extensions simultaneously. The resulting space-time may or may not be Hausdorff. (In fact, the extension is Hausdorff if and only if the space-time with g boundary is Hausdorff.) The concept of a non-Hausdorff space-time has not been arbitrarily introduced, but enters here in a natural way.

One is tempted at this point to discard the usual requirement that space-times be Hausdorff. However, this step may be too drastic, as the following example (based on a suggestion of Svetlichney) shows. Let A and B each represent the open interval $(0, 1)$. Let C denote the collection of all pairs $\{x, y\}$, of rational real numbers with $0 < x < 1$, $0 < y < 1$. Define:

$$M \equiv A \cup B \cup C,$$

the disjoint unions. Consider the following basis for the open sets on M :

1. The ordinary open sets of A .
2. The ordinary open sets of B .
3. A point $\{x, y\}$ of C , along with open sets (x, x') (with $x < x' < 1$) of A and (y, y') (with $y < y' < 1$) of B . M is a non-Hausdorff 1-manifold. Every neighborhood of the point $\{x, y\}$ of C intersects both

- (a) every neighborhood of the point x of A , and
- (b) every neighborhood of the point y of B . A continuous curve drawn in the direction of decreasing x along the line A may, at any rational point, jump continuously to any rational point of B . A continuous curve drawn in the direction of increasing x in A may not jump continuously to the line B .

This example is too pathological to be of physical interest. It would seem that some restriction must be imposed on those non-Hausdorff manifolds which are to be deemed acceptable candidates for a space-time manifold. Two possible restrictions immediately come to mind:

1. Only those non-Hausdorff space-times are permitted in which every geodesic has a unique extension.
2. Only those non-Hausdorff space-times are permitted in which every curve has no more than one end point.

It seems, however, that many more examples will be needed before one will be able to formulate the proper criteria for admitting non-Hausdorff cosmological models.

We have described a construction by which one may define and describe the "singular points" associated with any geodesically incomplete space-time. In view of the fact⁴⁸ that singularities occur as a characteristic feature of solutions of Einstein's equations, a description of the singularities would be a useful tool in understanding how to deal with them. The approach outlined here suffers from two significant disadvantages:

1. The construction is difficult to carry out in practice for some space-times.
2. Because of the freedom available in the formulation of many of the definitions, the construction is less natural than one would like.

ACKNOWLEDGMENT

I am deeply grateful to many individuals for stimulating discussions during the course of this work. Special thanks go to J. A. Wheeler, my thesis advisor, for his encouragement and numerous suggestions. Particularly useful comments and suggestions were offered by S. W. Hawking, A. Staruszkiewicz, D. Brill, and especially by R. Penrose and G. Svetlichney.

⁴⁸ R. Penrose, *Phys. Rev. Letters* **14**, 57 (1965); S. W. Hawking, *Phys. Rev. Letters* **17**, 445 (1966); *Proc. Roy. Soc. (London)* **294A**, 511 (1966); *Proc. Roy. Soc. (London)* **295A**, 490 (1966); R. Geroch, *Phys. Rev. Letters* **17**, 445 (1966).

General Griffiths' Inequalities on Correlations in Ising Ferromagnets*

D. G. KELLY† AND S. SHERMAN
 Department of Mathematics, Indiana University, Bloomington, Indiana

(Received 17 May 1967)

Let $N = (1, 2, \dots, n)$. For each subset A of N , let $J_A \geq 0$. For each $i \in N$, let $\sigma_i \pm 1$. For each subset A of N , define $\sigma^A = \prod_{i \in A} \sigma_i$. Let the Hamiltonian be $-\sum_{A \subset N} J_A \sigma^A$. Then for each $A, B \subset N$, $\langle \sigma^A \rangle \geq 0$ and $\langle \sigma^A \sigma^B \rangle - \langle \sigma^A \rangle \langle \sigma^B \rangle \geq 0$. This weakens the hypothesis and widens the conclusion of a result due to Griffiths.

INTRODUCTION

The inequality of Griffiths, which exhibits the monotonic behavior of the moments ("correlations") in a ferromagnetic Ising system as functions of the interactions, has consequences which are sufficiently interesting to warrant further study of the inequality. (For example, the appearance of long-range order in the three-dimensional Ising ferromagnet follows immediately.) In particular, one would like to know what makes the inequality work mathematically.

In this paper we generalize the result and give alternate proofs. An interesting problem is to obtain conditions on the moments of a distribution of a system of Ising spins which are necessary and sufficient for the distribution to have come from a ferromagnetic Hamiltonian. This problem is not solved here. In fact, a list of unsolved problems appears in the Appendix.

1. ISING MODEL WITH LONG-RANGE INTERACTION

Let N denote the index set $\{1, 2, \dots, n\}$; consider the space of all 2^n spin configurations $(\sigma_1, \sigma_2, \dots, \sigma_n)$, where each σ_i is allowed the values $+1$ ("up") or -1 ("down"). We customarily denote a general configuration by γ , and $(\sigma_i)_\gamma$ is the number (± 1) which appears as the i th spin (component) in γ .

Suppose that for each pair (i, j) of distinct indices in N the extended real number

$$J_{ij} = J_{ji} \geq 0 \tag{1.1}$$

is given ($J_{ij} = \infty$ is permitted). The requirement $J_{ij} \geq 0$ is that the system be *ferromagnetic*. The *Hamiltonian* of the system is a real-valued function on configurations, whose value at the configuration γ is

$$\mathcal{H}_\gamma = - \sum_{1 \leq i < j \leq n} J_{ij} (\sigma_i \sigma_j)_\gamma. \tag{1.2}$$

[Note that we have written $(\sigma_i \sigma_j)_\gamma$ for $(\sigma_i)_\gamma (\sigma_j)_\gamma$; we

also employ such abbreviations as $(\sigma_i + \sigma_j)_\gamma$ for $(\sigma_i)_\gamma + (\sigma_j)_\gamma$, etc.]

The Gibbs probability on the space of configurations is defined by

$$P(\gamma) = Z^{-1} \exp(-\beta \mathcal{H}_\gamma), \tag{1.3}$$

where

$$\beta = (kT)^{-1} > 0, \tag{1.4}$$

k being Boltzmann's constant and T the (absolute) temperature, and where the *partition function* Z is defined by

$$Z = \sum_\gamma \exp(-\beta \mathcal{H}_\gamma). \tag{1.5}$$

The expected value of a random variable X on this probability space is called its *thermal average* and is denoted by angular brackets:

$$\langle X \rangle = E(X) = \sum_\gamma X(\gamma) \exp(-\beta \mathcal{H}_\gamma) \left[\sum_\gamma \exp(-\beta \mathcal{H}_\gamma) \right]^{-1}. \tag{1.6}$$

As an example, we note that if σ_i denotes the random variable whose value at γ is $(\sigma_i)_\gamma$, then

$$\langle \sigma_i \rangle = 0. \tag{1.7}$$

This is true because $\mathcal{H}_\gamma = \mathcal{H}_{-\gamma}$ (where $-\gamma$ denotes the configuration obtained from γ by multiplying every spin by -1); hence $P(\gamma) = P(-\gamma)$, whereas $(\sigma_i)_\gamma = -(\sigma_i)_{-\gamma}$.

In fact, for the same reason,

$$\left\langle \prod_{i \in A} \sigma_i \right\rangle = 0 \tag{1.8}$$

whenever A contains an odd number of indices.

Griffiths¹ proved the following sets of inequalities:

$$\langle \sigma_k \sigma_l \rangle \geq 0 \quad \text{for all } k, l \in N. \tag{1.9}$$

$$\langle \sigma_k \sigma_l \sigma_m \sigma_n \rangle - \langle \sigma_k \sigma_l \rangle \langle \sigma_m \sigma_n \rangle \geq 0 \quad \text{for all } k, l, m, n \in N. \tag{1.10}$$

(Note that k, l, m , and n need not be distinct.)

In view of the fact that $\langle \sigma_k \rangle = 0$ and $\text{var}(\sigma_k) = 1$ for all k , Eq. (1.9) says that the covariance of any pair

* Partially supported by NSF GP 3941 and NSF GP 7469.
 † Present address: Jet Propulsion Laboratory 238-420, 4800 Oak Grove Drive, Pasadena, California 91103.

¹ R. B. Griffiths, *J. Math. Phys.* **8**, 478 (1967).

of spins is nonnegative; that is, the probability that they have the same sign is at least $\frac{1}{2}$.

On account of the identity

$$\frac{1}{\beta} \frac{\partial \langle \sigma_k \sigma_l \rangle}{\partial J_{mn}} = \langle \sigma_k \sigma_l \sigma_m \sigma_n \rangle - \langle \sigma_k \sigma_l \rangle \langle \sigma_m \sigma_n \rangle \quad (1.11)$$

(which may be verified by substituting (1.2) into (1.6) and differentiating), Eq. (1.10) says that increasing any one of the interaction strengths J_{mn} can never decrease the covariance of an arbitrary pair of spins. The physical consequences of this are rather far-reaching; (1.10) may be used, for example, for a simple proof of the existence of a phase transition in the three-dimensional Ising model (see Ref. 1).

We remark that the hypothesis $J_{ij} \geq 0$ is essential; both (1.9) and (1.10) fail without this assumption (see Sec. 6).

We shall refer to (1.9) and (1.10) as the *Griffiths inequalities*.

Classical Cases

In what we shall call the *classical Ising model* of a ferromagnet, the spins are considered to be situated at the vertices of a one-, two-, or three-dimensional square (or cubic) lattice, and the interactions are given by

$$J_{ij} = \begin{cases} J \geq 0 & \text{if spins } i \text{ and } j \text{ are} \\ & \text{nearest neighbors on} \\ & \text{the lattice,} \\ 0 & \text{otherwise.} \end{cases} \quad (1.12)$$

This model has been studied extensively (see Newell and Montroll,² Onsager,³ Griffiths,⁴ and their bibliographies).

A generalization of the classical Ising model which includes the long-range model is the model of a magnet in an *external magnetic field*, in which the Hamiltonian is given by

$$\mathcal{H}_\gamma^{\text{ext}} = \sum_{i < j} J_{ij} (\sigma_i \sigma_j)_\gamma - H \sum_{i=1}^n (\sigma_i)_\gamma, \quad (1.13)$$

where $H \geq 0$ is the external field strength.

The Griffiths inequalities hold for this model also; they are proved in Griffiths.⁵ They will appear also as a consequence of the extension in the next section.

2. GENERALIZATIONS; MAIN THEOREM

One extension of Griffiths' inequalities is as follows. For each subset A of the index set N , define

$$\sigma^A = \prod_{i \in A} \sigma_i \quad (\sigma^\emptyset \equiv 1). \quad (2.1)$$

Then under the same hypothesis as in Sec. 1, viz. (1.1), we have

$$\langle \sigma^R \rangle \geq 0 \quad \text{for all } R \subset N, \quad (2.2)$$

$$\langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle \geq 0 \quad \text{for all } R, S \subset N. \quad (2.3)$$

We shall not prove these now, as they are included in the main theorem below.

As an example of (2.3), let $R = \{k, l, m, n\}$ and $S = \{k, l\}$. We then have

$$\langle \sigma_m \sigma_n \rangle - \langle \sigma_k \sigma_l \sigma_m \sigma_n \rangle \langle \sigma_k \sigma_l \rangle \geq 0. \quad (2.4)$$

Here we have upper bounds on the fourth-order moments of the spins in terms of ratios of second-order moments. In contrast to this, (1.10) gives lower bounds on the fourth-order moments in terms of products of second-order moments.

We note in passing that

$$\sigma^R \sigma^S = \sigma^{R \Delta S}, \quad (2.5)$$

where $R \Delta S$ denotes the set-theoretic *symmetric difference*:

$$R \Delta S = (R \cup S) - (R \cap S). \quad (2.6)$$

An interpretation of (2.2) is similar to that for (1.9): The probability that a given set of spins contains an even number of "down" spins is at least $\frac{1}{2}$. However, no identity such as (1.11) exists (at this level of generality) to relate (2.2) and (2.3).

Main Theorem

We now extend the probability on the space of spin configurations. For each nonempty subset A of N , let the number

$$J_A > 0 \quad (J_\emptyset = 0) \quad (2.7)$$

be given, and define the Hamiltonian by

$$\mathcal{H}_\gamma = - \sum_{A \subset N} J_A (\sigma^A)_\gamma. \quad (2.8)$$

Probabilities are again defined by (1.3), (1.4), and (1.5), with (2.8) replacing (1.2) for \mathcal{H}_γ . [Since the only distinction between (1.2) and (2.8) is one of specialization, we use the same symbol for both.]

Theorem: In the probability space defined by (2.7), (2.8), (1.3), (1.4), and (1.5), we have

$$(I) \quad \langle \sigma^R \rangle \geq 0 \quad \text{for all } R \subset N,$$

$$(II) \quad \langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle \geq 0 \quad \text{for all } R, S \subset N.$$

With our new Hamiltonian we again have

$$\frac{1}{\beta} \frac{\partial \langle \sigma^R \rangle}{\partial J_S} = \langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle. \quad (2.9)$$

² G. Newell and E. W. Montroll, Rev. Mod. Phys. 25, 353 (1953).

³ L. Onsager, Phys. Rev. 65, 117 (1944).

⁴ R. B. Griffiths, Phys. Rev. 136, A437 (1964).

⁵ R. B. Griffiths, J. Math. Phys. 8, 484 (1967).

Hence II admits an interpretation like that of (1.10): The tendency of each set of spins to contain an even number of "down" spins is increased when any of the interaction constants J_S is increased.

The next three sections are devoted to the proofs of I and II, which we shall call the *generalized Griffiths' inequalities*.

3. NOTATION AND CONVENTIONS

Before embarking on the proofs, we introduce certain notational conveniences.

Observe first the effect of adding a constant C to the Hamiltonian: The partition function Z is multiplied by $\exp(-\beta C)$, but so is the numerator of (1.3), so that no probabilities, and hence no thermal averages, are changed. Hence, if we can prove I and II with a constant added to \mathcal{K}_γ , we shall have proved our theorem.

We in fact add the constant $C = \sum_{A \subset N} J_A$ to \mathcal{K} , using

$$\mathcal{K}_\gamma^* = - \sum_{A \subset N} J_A (\sigma^A - 1)_\gamma \tag{3.1}$$

instead of (2.8). In place of (1.5) we then have

$$Z^* = \sum_\gamma \exp(-\beta \mathcal{K}_\gamma^*). \tag{3.2}$$

In place of (1.6) we have

$$\langle X \rangle = \sum_\gamma X(\gamma) \exp(-\beta \mathcal{K}_\gamma^*) \left[\sum_\gamma \exp(-\beta \mathcal{K}_\gamma^*) \right]^{-1}. \tag{3.3}$$

Now

$$\sigma^A - 1 = \begin{cases} -2 & \text{if } \sigma^A = -1, \\ 0 & \text{if } \sigma^A = +1, \end{cases}$$

so that

$$\mathcal{K}_\gamma^* = 2 \sum_{\substack{A \subset N \\ (\sigma^A)_\gamma = -1}} J_A. \tag{3.4}$$

If, for each $A \subset N$, we let

$$x_A = \exp(-2\beta J_A), \tag{3.5}$$

we then have

$$\exp(-\beta \mathcal{K}_\gamma^*) = \prod_{\substack{A \subset N \\ (\sigma^A)_\gamma = -1}} x_A \geq 0. \tag{3.6}$$

Defining $Z_\gamma = \exp(-\beta \mathcal{K}_\gamma^*)$, we get

$$Z^* = \sum_\gamma Z_\gamma \geq 0, \tag{3.7}$$

and

$$\langle \sigma^R \rangle = (Z^*)^{-1} \sum_\gamma (\sigma^R)_\gamma Z_\gamma. \tag{3.8}$$

Note that $0 \leq J_A \leq \infty$ is equivalent to

$$0 \leq x_A \leq 1 \text{ for all } A \subset N. \tag{3.9}$$

The following "lemma" will be used only once, in the initial step of the inductive proof of I. We state it

as a lemma to avoid annoying details in that proof, and to establish a notation for what follows.

Lemma: If A and B are nonempty, unequal subsets of N , then the two sets of configurations

$$A^+ = \{\gamma: \sigma^A = +1\}, \quad A^- = \{\gamma: \sigma^A = -1\}$$

each contain 2^{n-1} configurations; and the four sets

$$A^+B^+ = \{\gamma: \sigma^A = +1, \sigma^B = +1\},$$

(A^+B^- , A^-B^+ , A^-B^- defined analogously) each contain 2^{n-2} configurations.

Proof: Choose an index in A , say i . A one-to-one correspondence between A^+ and A^- is given by $\gamma \leftrightarrow \gamma'$, where γ' is obtained from γ by multiplying σ_i by -1 (and γ is obtained from γ' in the same way).

Now suppose (without loss of generality) that B contains an index j not in A . Then we have a one-to-one correspondence between A^+B^+ and A^+B^- , and one between A^-B^+ and A^-B^- , both given by $\gamma \leftrightarrow \gamma''$, where γ'' is obtained from γ by multiplying σ_j by -1 .

Since A^+B^+ and A^+B^- are subsets of A^+ , and A^-B^+ and A^-B^- are subsets of A^- , each of the four has 2^{n-2} configurations. End of proof of lemma.

4. PROOF OF I

Since $Z^* > 0$, it suffices to prove that

$$A(R) = Z^* \langle \sigma^R \rangle = \sum_\gamma (\sigma^R)_\gamma Z_\gamma \geq 0 \text{ for all } R \subset N. \tag{4.1}$$

Lemma 1: In a system where $\langle \sigma^R \rangle \geq 0$ for all $R \subset N$, we have

$$\sum_{\gamma \in B^+} (\sigma^R)_\gamma Z_\gamma \geq 0 \text{ for all } R, B \subset N. \tag{4.2}$$

Proof: Notice that

$$\sigma^B \sigma^R + \sigma^R = \begin{cases} 2\sigma^R & \text{if } \sigma^B = +1 \\ 0 & \text{if } \sigma^B = -1. \end{cases}$$

Hence,

$$\begin{aligned} \sum_{\gamma \in B^+} (\sigma^R)_\gamma Z_\gamma &= \frac{1}{2} \sum_\gamma (\sigma^B \sigma^R + \sigma^R)_\gamma Z_\gamma \\ &= \frac{Z^*}{2} [\langle \sigma^{B \Delta R} \rangle + \langle \sigma^R \rangle], \end{aligned}$$

which is nonnegative by the hypothesis. End of proof of Lemma 1.

Now we prove (4.1) by induction on the number of J_A that are nonzero in the Hamiltonian.

Note that

$$A(R) = \sum_{\gamma \in R^+} Z_\gamma - \sum_{\gamma \in R^-} Z_\gamma. \tag{4.3}$$

If all J_A are equal to zero, then all x_A are 1, and so $Z_\gamma = 1$ for every γ . Hence,

$$A(R) = \sum_{\gamma \in R^+} 1 - \sum_{\gamma \in R^-} 1 = 2^{n-1} - 2^{n-1} = 0.$$

If $J_B > 0$ and all other J_A are 0, we have two cases: $B = R$ and $B \neq R$.

If $B = R$, then

$$Z_\gamma = \begin{cases} x_R < 1 & \text{if } (\sigma^R)_\gamma = -1, \\ 1 & \text{if } (\sigma^R)_\gamma = 1. \end{cases}$$

So

$$A(R) = \sum_{\gamma \in R^+} 1 - \sum_{\gamma \in R^-} x_R = 2^{n-1}(1 - x_R) > 0.$$

If $B \neq R$, then

$$Z_\gamma = \begin{cases} x_B < 1 & \text{if } (\sigma^B)_\gamma = -1, \\ 1 & \text{if } (\sigma^B)_\gamma = 1. \end{cases}$$

So

$$\begin{aligned} A(R) &= \sum_{\gamma \in R^+ B^+} 1 + \sum_{\gamma \in R^+ B^-} x_B - \sum_{\gamma \in R^- B^+} 1 - \sum_{\gamma \in R^- B^-} x_B \\ &= 2^{n-2}(1 + x_B - 1 - x_B) = 0. \end{aligned}$$

So $A(R) \geq 0$ for all R when at most one of the J_A is nonzero.

Suppose $A(R) \geq 0$ for all R in every system in which k (≥ 1) of the J_A are nonzero; consider a system in which $k + 1$ of the J_A , among them J_B , are nonzero.

Now $A(R)$ is a polynomial in the variables x_A which is an inhomogeneous linear function of each variable, and in particular, of x_B . So if we can prove $A(R) \geq 0$ for the end-point values $x_B = 0$ and $x_B = 1$ (with fixed but arbitrary values of the other x_A , only k of them being unequal to 1), the proof of I will be complete.

(Here is where the hypothesis $J_A \geq 0$ for all $A \subset N$ is used; if $J_B < 0$ is permitted, then $x_B > 1$, and this method of proof will fail.)

When $x_B = 1$, we have $J_B = 0$, and we are back to the case in which only k of the J_A are nonzero; the induction hypothesis gives $A(R) \geq 0$.

When $x_B = 0$, we have

$$Z_\gamma = \begin{cases} 0 & \text{when } (\sigma^B)_\gamma = -1 \\ Z'_\gamma & \text{when } (\sigma^B)_\gamma = 1, \end{cases}$$

where Z'_γ does not involve x_B . For this reason, Z'_γ can be considered to have come from a Hamiltonian in which $J_B = 0$, i.e., $x_B = 1$. We have

$$A(R) = \sum_{\gamma \in B^+} (\sigma^R)_\gamma Z'_\gamma.$$

But by the induction hypothesis and Lemma 1, $A(R) \geq 0$ for all R in such a system (only k of the J_A are nonzero). End of proof of I.

We note in passing that the hypothesis of Lemma 1 is just I, which has now been proved; hence we have (4.2) for all ferromagnetic systems. In probabilistic terms,

$$E(\sigma^R \mid \sigma^B = +1) \geq 0 \quad \text{for all } R, B \subset N.$$

5. PROOF OF II

First we derive two expressions for F , defined by

$$F = \frac{(Z^*)^2}{4} [\langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle]. \quad (5.1)$$

$$\begin{aligned} \langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle &= (Z^*)^{-2} \left(\sum_\gamma (\sigma^R \sigma^S)_\gamma Z_\gamma \right) \left(\sum_{\gamma'} Z_{\gamma'} \right) \\ &\quad - (Z^*)^{-2} \left(\sum_\gamma (\sigma^R)_\gamma Z_\gamma \right) \left(\sum_{\gamma'} (\sigma^S)_{\gamma'} Z_{\gamma'} \right), \end{aligned}$$

so that

$$F = \frac{1}{4} \sum_\gamma \sum_{\gamma'} [(\sigma^R \sigma^S)_\gamma - (\sigma^R)_\gamma (\sigma^S)_{\gamma'}] Z_\gamma Z_{\gamma'}. \quad (5.2)$$

For the second expression let

$$\begin{aligned} p &= \sum_{\gamma \in R^+ S^+} Z_\gamma, \\ q &= \sum_{\gamma \in R^+ S^-} Z_\gamma, \\ r &= \sum_{\gamma \in R^- S^+} Z_\gamma, \\ s &= \sum_{\gamma \in R^- S^-} Z_\gamma. \end{aligned} \quad (5.3)$$

Then we have

$$Z^* = p + q + r + s$$

and

$$Z^* \langle \sigma^R \sigma^S \rangle = p - q - r + s,$$

so

$$(Z^*)^2 \langle \sigma^R \sigma^S \rangle = (p + s)^2 - (q + r)^2. \quad (5.4)$$

Similarly,

$$(Z^*)^2 \langle \sigma^R \rangle \langle \sigma^S \rangle = (p - s)^2 - (q - r)^2. \quad (5.5)$$

Combining (5.4) and (5.5) yields

$$(Z^*)^2 [\langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle] = 4(ps - qr),$$

whence we get

$$F = ps - qr. \quad (5.6)$$

We next make some observations about the form of F . It is first of all a polynomial in the variables x_A , in which each x_A occurs in any term to power 0, 1, or 2. (For this proof we assume that all the variables x_A appear, that is, $0 < J_A < \infty$ for all $A \subset N$, even including x_ϕ . If $F \geq 0$ can be proved for such cases, it will follow by continuity of F for the cases in which some x_A are 1 or 0.)

Let us define the *Griffiths' form* of a polynomial, which is at most quadratic in each of its variables, as

follows. For each term, the *linear part* is the product of those variables appearing with exponent 1. (For example, the linear part of $-2x^2yzw^2$ is yz ; $-3y^2$ has no linear part.)

Now label the distinct linear parts of all terms, say g_1, g_2, \dots (some linear parts may serve for more than one term; there will be at most 2^n different linear parts). Combine all terms having the same linear part g_i , calling the sum $g_i G_i$, where G_i is called the *nonlinear part*. In addition, there may be terms without linear parts; combine them to form G_0 . The Griffiths' form is then

$$G_0 + \sum_i g_i G_i.$$

[For example, the Griffiths' form of

$$x^2 + wx^2y^2z - 2wxy^2 + wy^2z + wxz^2 + 3wy^2$$

is

$$(x^2) + wz(x^2y^2 + y^2) + wx(-2y^2 + z^2) + w(3y^2),$$

where the G_i are enclosed in parentheses.]

Now when F is put into Griffiths' form, there will be no G_0 , for $x_R x_S$ is a linear factor of every term. (Proof: x_R appears in every term of r and s and in no term of p and q ; x_S appears in every term of q and s and in no term of p and r .) We thus have

$$F = \sum_i g_i G_i. \tag{5.7}$$

Since $0 < x_A < 1$ for all $A \subset N$, we have $g_i \geq 0$ for all i . Fix an index $k \geq 1$; to prove II it suffices to show $G_k \geq 0$.

This is done by proving the three following statements:

(1) There exists a nonempty subset B_k of N such that

- (i) x_A appears in g_k iff $\#(A \cap B_k)$ is odd.
- (ii) $\#(R \cap B_k)$ and $\#(S \cap B_k)$ are odd.

(2) Consider the *reduced system*, obtained by setting all the x_A appearing in g_k equal to 1. For this system the Griffiths' form of F becomes

$$F^{(r)} = G'_0 + \sum_i g'_i G'_i, \tag{5.8}$$

and we have $G'_0 = G_k$.

(3) In the reduced system, $G'_0 \geq 0$.

II will of course follow when these three statements are proved.

Proof of Statement 1. First, (i) implies (ii), for x_R and x_S appear in g_i for every $i \geq 1$.

On account of (5.6), each term in F is of the form $\pm Z_\gamma Z_{\gamma'}$ for some pair (γ, γ') of *distinct* configurations.

So

$$g_k G_k = \pm Z_{\gamma_1} Z_{\gamma_2} \pm Z_{\gamma_3} Z_{\gamma_4} \pm \dots, \tag{5.9}$$

where the terms in the sum all have the same linear part g_k . We fix our attention on γ_1 and γ_2 only.

Let

$$B_k = \{i: (\sigma_i)_{\gamma_1} \neq (\sigma_i)_{\gamma_2}\} \neq \phi, \tag{5.10}$$

and for $j = 1$ or 2 let

$$D_j = \{i: (\sigma_i)_{\gamma_j} = -1\}, \tag{5.11}$$

which is the set of "down spins" in γ_j .

Now x_A appears in Z_{γ_j} iff $\#(A \cap D_j)$ is odd ($j = 1, 2$), and x_A appears in g_k iff x_A appears to power 1 in $Z_{\gamma_1} Z_{\gamma_2}$, i.e., iff one of $\#(A \cap D_1)$, $\#(A \cap D_2)$ is odd and the other even, i.e., iff $\#(A \cap D_1) + \#(A \cap D_2)$ is odd.

But

$$\begin{aligned} \#(A \cap D_1) + \#(A \cap D_2) &= \#(A \cap D_1 \cap B_k) + \#(A \cap D_1 \cap B_k^c) \\ &\quad + \#(A \cap D_2 \cap B_k) + \#(A \cap D_2 \cap B_k^c). \end{aligned}$$

However,

$$\begin{aligned} A \cap D_2 \cap B_k &= \{i \in A: (\sigma_i)_{\gamma_2} = -1 \text{ and } (\sigma_i)_{\gamma_1} \neq (\sigma_i)_{\gamma_2}\} \\ &= \{i \in A: (\sigma_i)_{\gamma_1} = 1 \text{ and } (\sigma_i)_{\gamma_1} \neq (\sigma_i)_{\gamma_2}\} \\ &= A \cap D_1^c \cap B_k, \end{aligned}$$

and

$$\begin{aligned} A \cap D_2 \cap B_k^c &= \{i \in A: (\sigma_i)_{\gamma_2} = -1 \text{ and } (\sigma_i)_{\gamma_1} = (\sigma_i)_{\gamma_2}\} \\ &= \{i \in A: (\sigma_i)_{\gamma_1} = -1 \text{ and } (\sigma_i)_{\gamma_1} = (\sigma_i)_{\gamma_2}\} \\ &= A \cap D_1 \cap B_k^c. \end{aligned}$$

Hence,

$$\begin{aligned} \#(A \cap D_1) + \#(A \cap D_2) &= \#(A \cap D_1 \cap B_k) + 2\#(A \cap D_1 \cap B_k^c) \\ &\quad + \#(A \cap D_1^c \cap B_k), \end{aligned}$$

which is odd iff $\#(A \cap D_1 \cap B_k) + \#(A \cap D_1^c \cap B_k)$ is odd. End of proof of Statement 1.

Proof of Statement 2. It is clear that

$$G'_0 = G_k + G_{k_1} + G_{k_2} + \dots,$$

where G_{k_1}, G_{k_2}, \dots are those G_i whose associated g_i contain only variables appearing also in g_k . But we shall show that for any two distinct linear terms, each always contains a variable not appearing in the other; it will follow that there is no g_i whose variables all appear in g_k , so that Statement 2 will follow.

In fact, each g_i contains precisely half of the variables x_A (including x_ϕ). [Proof: x_A appears in g_i

iff $\#(A \cap B_i)$ is odd, where B_i is as in Statement 1. But the sets having an odd intersection with a fixed nonempty set comprise half of all sets, for a one-to-one correspondence between these sets and the others is given by fixing an index $j \in B_i$ and associating a set C with $C \Delta \{j\}$.

But distinct sets having the same finite number of elements must each contain an element not in the other. End of proof of Statement 2.

Proof of Statement 3. On account of Statement 1, we have, in the reduced system,

$$Z_\gamma = \prod [x_A : \#(A \cap B_k) \text{ is even, } (\sigma^A)_\gamma = -1]. \quad (5.12)$$

For each configuration γ , let γ^* be the configuration obtained from γ by multiplying $(\sigma_i)_\gamma$ by -1 for each $i \in B_k$, and having all other $(\sigma_i)_\gamma$ unaltered.

Note that

$$Z_{\gamma^*}^{(r)} = Z_\gamma^{(r)}, \quad (5.13)$$

where the superscript r indicates that we are in the reduced system, i.e., the variables in g_k are being ignored.

$$Z_{\gamma^*}^{(r)} = \prod [x_A : \#(A \cap B_k) \text{ is even, } (\sigma^A)_{\gamma^*} = -1],$$

while in this product

$$\begin{aligned} (\sigma^A)_{\gamma^*} &= (\sigma^A \cap B_k)_{\gamma^*} (\sigma^A \cap B_k^c)_{\gamma^*} \\ &= (\sigma^A \cap B_k)_\gamma (\sigma^A \cap B_k^c)_\gamma, \end{aligned}$$

since $\#(A \cap B_k)$ is even, and γ and γ^* agree on $A \cap B_k^c$. So

$$(\sigma^A)_{\gamma^*} = (\sigma^A)_\gamma,$$

and hence

$$\begin{aligned} Z_{\gamma^*}^{(r)} &= \prod [x_A : \#(A \cap B_k) \text{ is even, } (\sigma^A)_\gamma = -1] \\ &= Z_\gamma^{(r)}. \end{aligned}$$

Also,

$$(\sigma^R)_{\gamma^*} = -(\sigma^R)_\gamma, \quad (5.14)$$

for the transformation $\gamma \rightarrow \gamma^*$ involves changing the spins in B_k , and $\#(B_k \cap R)$ is odd. Similarly,

$$(\sigma^S)_{\gamma^*} = -(\sigma^S)_\gamma. \quad (5.15)$$

Hence $p = s$ in the reduced system, for if $Z_\gamma^{(r)}$ appears in p , then $Z_{\gamma^*}^{(r)} = Z_\gamma^{(r)}$ appears in s . Similarly, $q = r$. So in the reduced system,

$$\begin{aligned} F^{(r)} &= p^2 - q^2 = \left(\sum_{\gamma \in R^+ S^+} Z_\gamma^{(r)} \right)^2 - \left(\sum_{\gamma \in R^+ S^-} Z_\gamma^{(r)} \right)^2; \\ F^{(r)} &= \sum_{\gamma \in R^+ S^+} (Z_\gamma^{(r)})^2 - \sum_{\gamma \in R^+ S^-} (Z_\gamma^{(r)})^2 + \text{cross terms.} \end{aligned} \quad (5.16)$$

But a cross term is of the form $Z_\gamma^{(r)} Z_{\gamma'}^{(r)}$, where γ and γ' are distinct configurations with $(\sigma^R)_\gamma = (\sigma^R)_{\gamma'}$ and $(\sigma^S)_\gamma = (\sigma^S)_{\gamma'}$. (Note that $Z_\gamma Z_{\gamma'}$ does not appear in the original F ; nevertheless, its reduction $Z_\gamma^{(r)} Z_{\gamma'}^{(r)}$ appears in $F^{(r)}$.) So the linear part of $Z_\gamma Z_{\gamma'}$ (before reduction) does not contain x_R and x_S , and hence is not g_k . Therefore, as in the proof of Statement 1, it contains a variable not in g_k . Hence $Z_\gamma^{(r)} Z_{\gamma'}^{(r)}$ has a linear part which remains after reduction.

Clearly, however, the two sums in (5.16) do not contain any linear terms; hence the nonlinear part of $F^{(r)}$ is

$$\begin{aligned} G'_0 &= \sum_{\gamma \in R^+ S^+} (Z_\gamma^{(r)})^2 - \sum_{\gamma \in R^+ S^-} (Z_\gamma^{(r)})^2 \\ &= \sum_{\gamma \in R^+} (\sigma^S)_\gamma (Z_\gamma^{(r)})^2. \end{aligned}$$

But this is precisely the left-hand side of (4.2), except that we are in a system where certain J_A have been set to zero and the remainder of the J_A have been doubled (thus causing each x_A and Z_γ to be squared). Such a system is still ferromagnetic, so that (4.2) holds. End of proof of Statement 3; end of proof of II.

We remark here the similarity in form between the above proofs and Griffiths' proofs. In both cases the proof of I is by induction on the number of non-vanishing interactions, and the proof of II uses the Griffiths' form and shows that each G_k is the G_0 of a reduced system, which is always nonnegative.

We conclude this section with a remark on the effect of increasing J_S . As observed earlier, II and (2.9) imply that none of the moments $\langle \sigma^R \rangle$ can decrease when J_S is increased; moreover, one has a strong feeling that $\langle \sigma^S \rangle$ should increase at a faster rate than any of the other $\langle \sigma^R \rangle$. That this is so is confirmed by the following:

Proposition:

$$\frac{\partial \langle \sigma^S \rangle}{\partial J_S} - \frac{\partial \langle \sigma^R \rangle}{\partial J_S} \geq 0 \quad \text{for all } R, S \subset N, \quad (5.17)$$

or

$$Q = 1 - \langle \sigma^S \rangle^2 - \langle \sigma^R \sigma^S \rangle + \langle \sigma^R \rangle \langle \sigma^S \rangle \geq 0. \quad (5.18)$$

Proof:

$$\begin{aligned} (Z^*)^2 Q &= (Z^*)^2 - (Z^* \langle \sigma^S \rangle)^2 - Z^* (Z^* \langle \sigma^R \sigma^S \rangle) \\ &\quad + (Z^* \langle \sigma^R \rangle) (Z^* \langle \sigma^S \rangle). \end{aligned}$$

With p, q, r , and s as in (5.3) we have

$$\begin{aligned} Z^* &= p + q + r + s, \\ Z^* \langle \sigma^R \rangle &= p + q - r - s, \\ Z^* \langle \sigma^S \rangle &= p - q + r - s, \\ Z^* \langle \sigma^R \sigma^S \rangle &= p - q - r + s. \end{aligned}$$

Hence

$$\begin{aligned} (Z^*)^2 Q &= (p + q + r + s)^2 - (p - q + r - s)^2 \\ &\quad - (p + q + r + s)(p - q - r + s) \\ &\quad + (p + q - r - s)(p - q + r - s) \\ &= 4(p + r)(q + s) - (p + s)^2 + (q + r)^2 \\ &\quad + (p - s)^2 - (q - r)^2 \\ &= 4(pq + 2qr + rs), \end{aligned}$$

which is nonnegative because $p, q, r,$ and s are nonnegative.

End of proof of proposition.

6. COUNTEREXAMPLE TO I AND II WHEN $J_A < 0$

First we give a general form for counterexamples with $N = \{1, 2, 3\}$; this will shorten the descriptions of the several counterexamples that will follow in this and later sections.

If $A = \{i_1, \dots, i_k\}$, we will denote J_A by J_{i_1, i_2, \dots, i_k} and x_A by x_{i_1, i_2, \dots, i_k} ; for example, $x_{\{1, 2\}}$ will be called x_{12} , and $x_{\{1\}}$ will be x_1 . So for $N = \{1, 2, 3\}$ a system is completely specified by giving the values of the seven numbers.

$$x_1, x_2, x_3, x_{12}, x_{13}, x_{23}, x_{123} \quad (x_\phi = 1 \text{ always}).$$

The configurations γ and their relative probabilities Z_γ will then be displayed as in Table I.

TABLE I. Configurations γ and their relative probabilities Z_γ .

σ_1	σ_2	σ_3	Z_γ
+1	+1	+1	1
+1	+1	-1	$x_{13}x_{23}x_3x_{123}$
+1	-1	+1	$x_{12}x_{23}x_3x_{123}$
+1	-1	-1	$x_{12}x_{13}x_3x_3$
-1	+1	+1	$x_{12}x_{13}x_1x_{123}$
-1	+1	-1	$x_{12}x_{23}x_1x_3$
-1	-1	+1	$x_{13}x_{23}x_1x_2$
-1	-1	-1	$x_1x_2x_3x_{123}$

Example 6.1: Here we allow some of the J_A to be negative, and show that I and II fail. $J_A < 0$ is equivalent to $x_A > 1$.

$$\text{Let } x_1 = x_2 = x_3 = x_{123} = x_{13} = 1,$$

$$x_{12} = x > 1,$$

$$x_{23} = y > 1.$$

$$Z^* = 2(1 + x)(1 + y),$$

$$Z^* \langle \sigma_1 \sigma_2 \rangle = 2(1 - x)(1 + y) < 0,$$

$$Z^* \langle \sigma_2 \sigma_3 \rangle = 2(1 + x)(1 - y) < 0,$$

$$Z^* \langle \sigma_1 \sigma_3 \rangle = 2(1 - x)(1 - y) > 0.$$

TABLE II. Relative probabilities for Example 6.1.

σ_1	σ_2	σ_3	Z_γ
+1	+1	+1	1
+1	+1	-1	y
+1	-1	+1	xy
+1	-1	-1	x
-1	+1	+1	x
-1	+1	-1	xy
-1	-1	+1	y
-1	-1	-1	1

Now

$$\frac{1}{\beta} \frac{\partial \langle \sigma_1 \sigma_3 \rangle}{\partial J_{23}} = \langle \sigma_1 \sigma_2 \rangle - \langle \sigma_1 \sigma_3 \rangle \langle \sigma_2 \sigma_3 \rangle,$$

and

$$\begin{aligned} \frac{(Z^*)^2}{4} [\langle \sigma_1 \sigma_2 \rangle - \langle \sigma_1 \sigma_3 \rangle \langle \sigma_2 \sigma_3 \rangle] \\ &= (1 - x)(1 + y)(1 + x)(1 + y) \\ &\quad - (1 - x)(1 - y)(1 + x)(1 - y) \\ &= (1 - x^2)4y < 0. \end{aligned}$$

Note that this counterexample serves for Griffiths' model also, since all interactions are binary.

7. THIRD-ORDER INEQUALITIES

In the system of Sec. 2 notice that

$$\langle \sigma^R \rangle = \frac{1}{\beta} \frac{\partial \ln Z}{\partial J_R}, \tag{7.1}$$

and so

$$\langle \sigma^R \sigma^S \rangle - \langle \sigma^R \rangle \langle \sigma^S \rangle = \frac{1}{\beta^2} \frac{\partial^2 \ln Z}{\partial J_R \partial J_S}. \tag{7.2}$$

It is natural to consider

$$\begin{aligned} \frac{1}{\beta^3} \frac{\partial^3 \ln Z}{\partial J_R \partial J_S \partial J_T} &= \langle \sigma^R \sigma^S \sigma^T \rangle - \langle \sigma^R \rangle \langle \sigma^S \sigma^T \rangle \\ &\quad - \langle \sigma^S \rangle \langle \sigma^R \sigma^T \rangle - \langle \sigma^T \rangle \langle \sigma^R \sigma^S \rangle + 2 \langle \sigma^R \rangle \langle \sigma^S \rangle \langle \sigma^T \rangle. \end{aligned} \tag{7.3}$$

We give two examples to show that this may be either positive or negative.

Example 7.1: $R = S$. The right-hand side of (7.3) is then

$$\begin{aligned} -2 \langle \sigma^R \rangle \langle \sigma^R \sigma^T \rangle + 2 \langle \sigma^R \rangle^2 \langle \sigma^T \rangle \\ = -2 \langle \sigma^R \rangle [\langle \sigma^R \sigma^T \rangle - \langle \sigma^R \rangle \langle \sigma^T \rangle]. \end{aligned}$$

By I and II, this is never positive; there are easy examples in which the inequalities I and II are both strict, so that the right side of (7.3) may be negative. (Such an easy example is $N = \{1, 2, 3\}$; $R = T = \{1, 2\}$; $J_{\{1, 2\}} = J > 0$, $J_A = 0$ for all other $A \subset N$.) End of Example 7.1.

Example 7.2: In the long-range model of Sec. 1, let $R = \{k\}$, $S = \{l\}$, $T = \{k, l\}$. Then the right-hand side of (7.3) is

$$1 - \langle \sigma_k \rangle^2 - \langle \sigma_l \rangle^2 - \langle \sigma_k \sigma_l \rangle^2 + 2 \langle \sigma_k \rangle \langle \sigma_l \rangle \langle \sigma_l \sigma_k \rangle.$$

But by (1.7), this is just

$$1 - \langle \sigma_k \sigma_l \rangle^2 \geq 0.$$

Again, this is strict in some cases; the same "easy example" given above is such a case (where $k = 1$, $l = 2$). End of Example 7.2.

Concavity of Magnetization

The above questions are interesting partly because of the following considerations.

Consider the *external magnetic field* model, given by (1.13) above. Define

$$s = \sigma_1 + \dots + \sigma_n, \tag{7.4}$$

and

$$M = \frac{\langle s \rangle}{n},$$

called the "average magnetization per spin." Since $0 \leq \langle \sigma_i \rangle \leq 1$, we have

$$0 \leq M \leq 1. \tag{7.5}$$

Applying II and increasing the $J_{(i)}$ one at a time from H to $H + k$ will show that M is an increasing function of H :

$$M(H) \leq M(H + k) \text{ for } k \geq 0. \tag{7.6}$$

Furthermore, $\langle \sigma_i \rangle \rightarrow 1$ as $H \rightarrow \infty$. This is true because the numerator and denominator of (3.8) are jointly continuous in the $x_{(j)}$, and tend to 1 as the $x_{(j)}$ tend to zero.

The open question is: Is M a concave function of H ? One feels that the situation must be rather pathological if this is not the case, but we know of no proof of concavity.

Since M is a rational function of e^H , the other variables J_{ij} being fixed, concavity of $M(H)$ is equivalent to

$$\frac{\partial^2 M}{\partial H^2} \leq 0. \tag{7.7}$$

Performing the differentiation as above, we find this to be equivalent to

$$D = \langle s^3 \rangle - 3 \langle s \rangle \langle s^2 \rangle + 2 \langle s \rangle^2 \leq 0. \tag{7.8}$$

Mere manipulation shows that

$$D = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n D_{ijk}, \tag{7.9}$$

where

$$D_{ijk} = \langle \sigma_i \sigma_j \sigma_k \rangle - \langle \sigma_i \rangle \langle \sigma_j \sigma_k \rangle - \langle \sigma_j \rangle \langle \sigma_i \sigma_k \rangle - \langle \sigma_k \rangle \langle \sigma_i \sigma_j \rangle + 2 \langle \sigma_i \rangle \langle \sigma_j \rangle \langle \sigma_k \rangle. \tag{7.10}$$

So by (7.3), a sufficient condition for concavity of $M(H)$ is that, in the general system (2.8),

$$\left[\frac{1}{\beta^3} \frac{\partial^3 \ln Z}{\partial J_i \partial J_j \partial J_k} \right]_{J_1 = \dots = J_n = H} \leq 0 \text{ for all } H > 0 \tag{7.11}$$

(where we have abbreviated $J_{(i)}$ by J_i).

Notice that Example 7.2 does not contradict this proposition, for in Example 7.2 the sets R, S, T are not distinct sets of one spin each, as they are in (7.11). However, the following example shows that (7.11) is not true in general when J_A for other than one- and two-element sets are permitted.

Example 7.3: $N = \{1, 2, 3\}$; $x_{12} = x_{13} = x_{23} = \sqrt{\alpha}$, $x_1 = x_2 = x_3 = x_{123} = \sqrt{\beta}$, $x_\phi = 1$. Table III shows configurations and their relative probabilities.

$$Z^* = 1 + 6\alpha\beta + \beta^2,$$

$$Z^* \langle \sigma_1 \rangle = Z^* \langle \sigma_2 \rangle = Z^* \langle \sigma_3 \rangle = Z^* \langle \sigma_1 \sigma_2 \sigma_3 \rangle = 1 - \beta^2,$$

$$Z^* \langle \sigma_1 \sigma_2 \rangle = Z^* \langle \sigma_1 \sigma_3 \rangle = Z^* \langle \sigma_2 \sigma_3 \rangle = 1 - 2\alpha\beta + \beta^2.$$

Then

$$\frac{Z^3 D_{123}}{(1 - \beta^2)^3} = +8\beta^2(9\alpha^2 - 1),$$

which is positive when $\alpha > \frac{1}{9}$ and negative when $\alpha < \frac{1}{9}$. End of Example 7.3.

TABLE III. Relative probabilities for Example 7.3.

σ_1	σ_2	σ_3	Z_γ
+1	+1	+1	1
+1	+1	-1	$\alpha\beta$
+1	-1	+1	$\alpha\beta$
+1	-1	-1	$\alpha\beta$
-1	+1	+1	$\alpha\beta$
-1	+1	-1	$\alpha\beta$
-1	-1	+1	$\alpha\beta$
-1	-1	-1	β^2

However, the following question is still open, and an affirmative answer would provide a proof of the concavity of $M(H)$: In the external magnetic-field system, is $D_{ijk} \geq 0$ for all triplets (i, j, k) of distinct indices?

Notice that $D_{ijk} \geq 0$ whenever any two of i, j, k are equal, as seen in Example 7.1.

A further consequence of the result in Example 7.1, namely

$$\frac{\partial^2 \ln Z}{\partial J_R \partial J_R \partial J_T} \leq 0,$$

is that the second partial derivative of any σ^T with

respect to any J_R is nonpositive, for

$$\frac{\partial^2 \langle \sigma^T \rangle}{\partial J_R^2} = \frac{\partial^3 \ln Z}{\partial J_T \partial J_R \partial J_R} = \frac{\partial^3 \ln Z}{\partial J_R \partial J_R \partial J_T}.$$

[Compare this with the result of II and (2.9) that the first partial is nonnegative.]

8. CONVERSE PROBLEM

Given any system of n random variables X_1, \dots, X_n , each assuming values ± 1 , we can solve the 2^n equations

$$Z^{-1} \exp(-\beta J_\gamma) = \Pr [(X_1, \dots, X_n) = \gamma] \quad (8.1)$$

for the 2^n constants J_A . The question is: Under what conditions are all $J_A \geq 0$? In other words, we have in I and II a necessary condition for ferromagnetism; what is a sufficient condition?

That I and II alone are not sufficient is shown by the following example:

Example 8.1: Let $n = 4$, and consider random variables X_1, X_2, X_3, X_4 . Let the probabilities of the configurations be as follows:

$$\begin{aligned} P(++++) &= p_0 & P(-+++ &= p_8 \\ P(+++-) &= p_1 & P(-++-) &= p_9 \\ P(++-+) &= p_2 & P(-+-+) &= p_{10} \\ P(++--) &= p_3 & P(-+--) &= p_{11} \\ P(+--+ &= p_4 & P(--++) &= p_{12} \\ P(+--+) &= p_5 & P(--+-) &= p_{13} \\ P(+---) &= p_6 & P(----) &= p_{14} \\ P(+----) &= p_7 & P(----) &= p_{15}. \end{aligned}$$

[Here $(+-+-)$, for example, is the event $(X_1 = +1, X_2 = -1, X_3 = +1, X_4 = -1)$.]

If these probabilities come from a Hamiltonian as in (2.8), (1.3), and (1.5), then it can be seen that

$$x_{1234} = \frac{p_1 p_2 p_4 p_7 p_8 p_{11} p_{13} p_{14}}{p_0 p_3 p_5 p_6 p_9 p_{10} p_{12} p_{15}},$$

where

$$x_{1234} = e^{-2\beta J_{1234}}.$$

We show values of p_0, \dots, p_{15} such that $x_{1234} > 1$, but such that I and II hold, thus contradicting the sufficiency of I and II for ferromagnetism.

$$\text{Let } p_0 = p_{15} = (1 + \alpha^4)K^{-1},$$

$$\begin{aligned} p_1 = p_2 = p_4 = p_7 = p_8 = p_{11} = p_{13} \\ = p_{14} = (\alpha + \alpha^3)K^{-1}, \end{aligned}$$

$$p_3 = p_5 = p_6 = p_9 = p_{10} = p_{12} = (2\alpha^2)K^{-1},$$

where α will be some small positive number and K is the proper norming constant to make $\sum_0^{15} p_i = 1$.

Then we have

$$x_{1234} = \frac{(\alpha + \alpha^3)^8}{(1 + \alpha^4)^2 (2\alpha^2)^6} = \frac{\alpha^8 + o(\alpha^8)}{2^6 \alpha^{12} + o(\alpha^{12})},$$

where $o(\alpha^K)$ divided by α^K tends to 0 as α tends to 0.

$$x_{1234} = \frac{\alpha^8 + o(\alpha^8)}{o(\alpha^8)} \rightarrow \infty \text{ as } \alpha \rightarrow 0.$$

So x_{1234} certainly exceeds 1 for some small positive value of α .

To show that the random variables X_1, X_2, X_3 , and X_4 satisfy I and II, consider the following Ising model:

$$N = \{1, 2, 3, 4, 5\},$$

$$x_{15} = x_{25} = x_{35} = x_{45} = \alpha,$$

$$x_A = 1 \text{ unless } A = \{1, 5\}, \{2, 5\}, \{3, 5\}, \{4, 5\}.$$

This is in fact a Griffiths' model, so I and II are certainly satisfied. Hence I and II are satisfied for the system $\sigma_1, \sigma_2, \sigma_3, \sigma_4$ only; for certainly if I and II hold for all $R, S \subset \{1, 2, 3, 4, 5\}$, then they hold for $R, S \subset \{1, 2, 3, 4\}$. But the marginal probabilities for $\sigma_1, \sigma_2, \sigma_3, \sigma_4$ are seen to be precisely the values we have prescribed above. End of Example 8.1.

In fact, the question of finding sufficient conditions for ferromagnetism is, at the time of writing, open. The remainder of this paper will consist of different approaches to the problem, which are motivated by this question.

9. EXPONENTIAL-EXPANSION PROOF OF I AND II

In what follows we shall ignore the positive constant β ; to make this consistent with what we have above, we replace each J_A by βJ_A .

From (1.6) and (2.8) we have, for any $R \subset N$,

$$\begin{aligned} Z \langle \sigma^R \rangle &= \sum_\gamma (\sigma^R)_\gamma \exp \left[\sum_{A \subset N} J_A (\sigma^A)_\gamma \right] \\ &= \sum_\gamma (\sigma^R)_\gamma \sum_{k=0}^\infty \frac{1}{k!} \left[\sum_{A \subset N} J_A (\sigma^A)_\gamma \right]^k \\ &= \sum_{k=0}^\infty \frac{1}{k!} \sum_\gamma (\sigma^R)_\gamma \left[\sum_{A \subset N} J_A (\sigma^A)_\gamma \right]^k \\ &= \sum_{k=0}^\infty \frac{1}{k!} \sum_\gamma (\sigma^R)_\gamma \\ &\quad \times \sum_{A_1 \subset N} \dots \sum_{A_k \subset N} J_{A_1} \dots J_{A_k} (\sigma^{A_1 \Delta \dots \Delta A_k})_\gamma \\ &= \sum_{k=0}^\infty \frac{1}{k!} \sum_{A_1 \subset N} \dots \sum_{A_k \subset N} J_{A_1} \dots J_{A_k} \\ &\quad \times \sum_\gamma (\sigma^{A_1 \Delta \dots \Delta A_k \Delta R})_\gamma. \end{aligned}$$

Now for any set $B \subset N$,

$$\sum_{\gamma} (\sigma^B)_{\gamma} = \sum_{\gamma \in B^+} 1 - \sum_{\gamma \in B^-} 1 = \begin{cases} 0 & \text{if } B \neq \phi \\ 2^n & \text{if } B = \phi. \end{cases} \quad (9.1)$$

So

$$Z\langle \sigma^R \rangle = 2^n \sum_{k=0}^{\infty} \frac{1}{k!} \sum J_{A_1} \cdots J_{A_k}, \quad (9.2)$$

the sum being taken over all ordered k -tuples (A_1, \dots, A_k) of (not necessarily distinct) subsets of N satisfying $A_1 \Delta \cdots \Delta A_k = R$.

Now consider *multiplicity functions* μ on the subsets of N , that is, functions assigning a nonnegative integer $\mu(A)$ to each $A \subset N$. Define

$$J_{\mu} = \prod_{A \subset N} J_A^{\mu(A)}, \quad (9.3)$$

$$\mu! = \prod_{A \subset N} (\mu(A)!), \quad (9.4)$$

$$\Delta \mu = \Delta \prod_{A \subset N} A^{\mu(A)}; \quad (9.5)$$

the last expression denotes the symmetric difference of all the subsets of N , each subset A being taken $\mu(A)$ times.

For example, suppose A , B , and C are subsets of N , and suppose $\mu(A) = 3$, $\mu(B) = 1$, $\mu(C) = 2$, and $\mu(R) = 0$ unless R is A or B or C . Then

$$J_{\mu} = J_A^3 J_B J_C^2,$$

$$\mu = 3!1!2! = 12,$$

$$\Delta \mu = A \Delta A \Delta A \Delta B \Delta C \Delta C = A \Delta B.$$

For simplicity, we will sometimes denote a multiplicity function μ by listing the subsets of N , each set A appearing $\mu(A)$ times in the list. For example, the μ considered above would be listed

$$\mu = (AAABCC).$$

Using this notation, we shall list the μ which assigns multiplicity 0 to every set, simply as (0) . Note that $\Delta 0 = \phi$.

Now choose an arbitrary μ and suppose

$$\sum_{A \subset N} \mu(A) = k, \quad \Delta \mu = R;$$

that is, there are k items on the list for μ (possibly including some repetitions), and the symmetric difference of these items is R . Then there are exactly $k!/|\mu|$ ordered k -tuples (A_1, \dots, A_k) of subsets of N having the property that

$$J_{\mu} = J_{A_1} \cdots J_{A_k}.$$

Each of these k -tuples has in addition the property

$$A_1 \Delta \cdots \Delta A_k = R.$$

Because of this, we can rewrite (9.2) as

$$Z\langle \sigma^R \rangle = 2^n \sum_{\mu: \Delta \mu = R} \frac{1}{\mu!} J_{\mu}. \quad (9.6)$$

Finally, we define, for each subset R of N ,

$$\Sigma(R) = Z\langle \sigma^R \rangle 2^{-n} = \sum_{\mu: \Delta \mu = R} \frac{1}{\mu!} J_{\mu}. \quad (9.7)$$

It is clear that

$$\Sigma(\phi) = 2^{-n} Z. \quad (9.8)$$

With this notation, I and II can be restated:

$$(I') \quad \Sigma(R) \geq 0 \quad \text{for all } R \subset N;$$

$$(II') \quad \Sigma(\phi) \Sigma(R \Delta S) \geq \Sigma(R) \Sigma(S) \quad \text{for all } R, S \subset N.$$

Now notice that I' is obvious, in view of the definition (9.7), since $J_A \geq 0$ for all $A \subset N$.

Proof of II': We show that for each μ the coefficient of J_{μ} on the right of II' is either 0 or else equal to the coefficient of J_{μ} on the left (which of course is nonnegative).

Since $J_{v_1} J_{v_2} = J_{v_1 \cup v_2}$, we have

$$\begin{aligned} \Sigma(\phi) \Sigma(R \Delta S) &= \sum_{v_1: \Delta v_1 = \phi} \sum_{v_2: \Delta v_2 = R \Delta S} \frac{1}{v_1! v_2!} J_{v_1 \cup v_2} \\ &= \sum_{\mu: \Delta \mu = R \Delta S} \sum_{\substack{v \leq \mu \\ \Delta v = \phi}} \frac{1}{v!(\mu - v)!} J_{\mu}. \end{aligned} \quad (9.9)$$

Similarly,

$$\Sigma(R) \Sigma(S) = \sum_{\mu: \Delta \mu = R \Delta S} \sum_{\substack{v \leq \mu \\ \Delta v = R}} \frac{1}{v!(\mu - v)!} J_{\mu}. \quad (9.10)$$

So J_{μ} will appear in II' if and only if $\Delta \mu = R \Delta S$; in this case the coefficients of J_{μ} on the left and right of II' are, respectively,

$$L_{\mu}(\phi) = \sum_{\substack{v \leq \mu \\ \Delta v = \phi}} \frac{1}{v!(\mu - v)!} \quad (9.11)$$

and

$$L_{\mu}(R) = \sum_{\substack{v \leq \mu \\ \Delta v = R}} \frac{1}{v!(\mu - v)!}.$$

Defining, for each $A \subset N$,

$$B_{\mu}(A) = \{v \leq \mu: \Delta v = A\}, \quad (9.12)$$

we have

$$L_{\mu}(A) = \sum_{v \in B_{\mu}(A)} \frac{1}{v!(\mu - v)!}, \quad (9.13)$$

and $L_{\mu}(A) = 0$ iff $B_{\mu}(A)$ is empty. Clearly, $B_{\mu}(\phi)$ is not empty, for it contains the multiplicity function 0. We prove the following proposition:

Proposition 1: If $B_\mu(R)$ is not empty, then $L_\mu(R) = L_\mu(\phi)$.

Now observe that every multiplicity function μ can be built up, starting with some μ_0 satisfying $\mu_0(A) \leq 2$ for all $A \subset N$, by successively adding 2 to the values of $\mu(B)$ for B which have $\mu_0(B) > 0$. For example,

$$\mu = (AAAAABBBBCDD)$$

can be obtained from

$$\mu_0 = (ABBCDD)$$

in stages as follows:

$$\begin{aligned} \mu_1 &= (AAABBCDD) = \mu_0 + 2\chi_{\{A\}}, \\ \mu_2 &= (AAAAABBCDD) = \mu_1 + 2\chi_{\{A\}}, \\ \mu &= \mu_2 + 2\chi_{\{B\}}, \end{aligned}$$

where $\chi_{\{A\}}$ is the characteristic function of the class $\{A\}$, i.e., the function assigning multiplicity 1 to the set A and 0 to all other sets.

Notice also that $\Delta\mu = \Delta\mu_0$ in the above discussion. We shall prove Proposition 1 in two steps:

- (i) Proposition 1 is true provided $\mu(A) \leq 2$ for all $A \subset N$.
- (ii) If Proposition 1 is true for μ , then it is true for $\mu + 2\chi_{\{A\}}$, where A is any set for which $\mu(A) > 0$.

The following example is for illustrative purposes only; we include it to demonstrate how we will compare $L_\mu(R)$ and $L_\mu(\phi)$.

Example 9.1: Let $N = \{1, 2, 3, 4\}$, and let

$$\begin{aligned} A &= \{1, 2\}, & B &= \{1, 3\}, & C &= \{1, 4\}, \\ D &= \{2, 3\}, & E &= \{2, 4\}, & F &= \{3, 4\}. \end{aligned}$$

Let $R = A$ and $S = F$, so that $R\Delta S = N$. Fix $\mu = (AABDDE)$. We list $B_\mu(\phi)$ and $B_\mu(R)$, and next

TABLE IV

$B_\mu(\phi)$	$\frac{1}{v!(\mu - v)!}$	$B_\mu(R)$	$\frac{1}{v!(\mu - v)!}$
(0)	$\frac{1}{2!2!}$	(A)	$\frac{1}{2!}$
(AA)	$\frac{1}{2!2!}$	(ADD)	$\frac{1}{2!}$
(DD)	$\frac{1}{2!2!}$	(BD)	$\frac{1}{2!}$
(AADD)	$\frac{1}{2!2!}$	(AABD)	$\frac{1}{2!}$
(ABD)	1		

to each v in these classes we write $[1/v!(\mu - v)!]$. (See Table IV.)

Note that here $\mu(T) \leq 2$ for all $T \subset N$, that $B_\mu(R)$ is nonempty, and also that $L_\mu(R) = L_\mu(\phi) = 2$. So in this example, (i) is confirmed. End of Example 9.1.

Now we prove (i).

Motivated by Example 9.1, we define an equivalence relation on the class of multiplicity functions v for which $v(A) \leq 2$ for all $A \subset N$, as follows: v_1 and v_2 are equivalent if $|v_1 - v_2|$ assigns the value 0 or 2 to every set. [For example, (B), (BDD), (AABDD), etc. are all equivalent.]

Now $B_\mu(\phi)$ and $B_\mu(R)$ split into equivalence classes; each class contains one and only one characteristic function (that is, a function whose values are 0 and 1). [For example, in Example 9.1, $B_\mu(\phi)$ and $B_\mu(R)$ each contain two equivalence classes; the four representatives are (0), (ABD), (A), and (BD). We have listed the functions above to display the classes.]

Now each equivalence class will contain exactly 2^k members, where k is the number of sets A for which $\mu(A) = 2$ and which are assigned multiplicity 0 by the class representative. [For example, the first listed equivalence class in $B_\mu(\phi)$ in Example 9.1 contains 2^2 members, as there are 2 sets, A and D , which are given multiplicities 2 by μ and 0 by the representative (0).]

However, it is easily seen that $[1/v!(\mu - v)!]$ is $1/(2!)^k$ for each of these members. [This is where the property $\mu(A) \leq 2$ for all $A \subset N$ is essential.] Hence the contribution of every equivalence class to $L_\mu(\phi)$ or $L_\mu(R)$ is just 1.

To prove (i), it thus remains to show only that $B_\mu(\phi)$ and $B_\mu(R)$ contain the same number of equivalence classes, i.e., the same number of characteristic functions.

A one-to-one correspondence between the characteristic functions in $B_\mu(\phi)$ and those in $B_\mu(R)$ is established as follows. For two characteristic functions v and η , define $v\Delta\eta$ to be the Boolean sum of v and η :

$$(v\Delta\eta)(A) = \begin{cases} 0 & \text{if } v(A) = \eta(A), \\ 1 & \text{if } v(A) \neq \eta(A). \end{cases}$$

Characteristic functions form a group under Boolean sum; the identity is 0, and $v^{-1} = v$. Furthermore, we can see that if $\Delta v = A$ and $\Delta\eta = B$, then $\Delta(v\Delta\eta) = A\Delta B$.

Now $B_\mu(R)$ is nonempty by hypothesis, so there exists $v_0 \leq \mu$ with $\Delta v_0 = R$. So if $\eta \in B_\mu(\phi)$, then

$\eta\Delta v_0 \in B_\mu(R)$, since $\Delta(\eta\Delta v_0) = (\Delta\eta)\Delta(\Delta v_0) = \phi\Delta R = R$. Hence the mapping

$$\eta \rightarrow \eta\Delta v_0$$

is a map from $B_\mu(\phi)$ into $B_\mu(R)$; because of the group property above, it is one-to-one. And it is onto because the same operation maps $B_\mu(R)$ into $B_\mu(\phi)$. End of proof of (i).

We proceed now with the proof of (ii).

Let μ be an arbitrary multiplicity function satisfying $\Delta_\mu = R\Delta S$ and such that $B_\mu(R)$ is nonempty. By the hypothesis of (ii), $L_\mu(\phi) = L_\mu(R)$; that is,

$$\sum_{\eta \in B_\mu(\phi)} \frac{1}{\eta!(\mu - \eta)!} = \sum_{\eta \in B_\mu(R)} \frac{1}{\eta!(\mu - \eta)!}.$$

Suppose A is a particular set with $\mu(A) = k > 0$, and let

$$\mu_1 = \mu + 2\chi_{\{A\}}.$$

[For example, if $\mu = (A A A B B C D)$, then $\mu_1 = (A A A A A B B C D)$.]

We need to show that $L_{\mu_1}(\phi) = L_{\mu_1}(R)$.

Define a new notion of equivalence as follows: v_1 and v_2 are equivalent if $|v_1 - v_2|$ assigns an even multiplicity to A , and 0 to all other sets. [For example, (BC) , $(AABC)$, and $(A A A A B C)$ are equivalent, while (BC) and $(B B B C)$ are not.] Again, this is indeed an equivalence.

Now the equivalence classes in $B_\mu(\phi)$, or in $B_\mu(R)$, will have one of the following forms:

$$\begin{aligned} (B_1 B_2 \cdots) & \quad \frac{1}{0!k!\alpha} \\ (A A B_1 B_2 \cdots) & \quad \frac{1}{2!(k-2)!\alpha} \\ (A A A A B_1 B_2 \cdots) & \quad \frac{1}{4!(k-4)!\alpha} \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ (\underbrace{A \cdots A}_{k \text{ or } k-1} B_1 B_2 \cdots) & \quad \frac{1}{k!0!\alpha} \text{ or } \frac{1}{(k-1)!1!\alpha} \end{aligned}$$

or

$$\begin{aligned} (A B_1 B_2 \cdots) & \quad \frac{1}{1!(k-1)!\alpha} \\ (A A A B_1 B_2 \cdots) & \quad \frac{1}{3!(k-3)!\alpha} \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ (\underbrace{A \cdots A}_{k \text{ or } k-1} B_1 B_2 \cdots) & \quad \frac{1}{k!0!\alpha} \text{ or } \frac{1}{(k-1)!1!\alpha}. \end{aligned}$$

Here we have again listed $[1/\nu!(\mu - \nu)!]$ next to each entry; α represents the factorials arising from multiplicities among B_1, B_2, \dots .

So the contribution of any equivalence class to $L_\mu(\phi)$ or $L_\mu(R)$ is either

$$\frac{1}{\alpha k!} \left[\binom{k}{0} + \binom{k}{2} + \cdots \right]$$

or

$$\frac{1}{\alpha k!} \left[\binom{k}{1} + \binom{k}{3} + \cdots \right].$$

(We have used the convention that $\binom{k}{j} = 0$ when $j > k$.) But the above two expressions are both equal to $2^{k-1}/\alpha k!$. That is, the contribution of any equivalence class in $B_\mu(\phi)$ [or $B_\mu(R)$] to $L_\mu(\phi)$ [or $L_\mu(R)$] is $2^{k-1}/\alpha k!$, where only α depends on which equivalence class is being considered.

Now corresponding to each such equivalence class in $B_\mu(\phi)$ or $B_\mu(R)$ is an equivalence class in $B_{\mu_1}(\phi)$ or $B_{\mu_1}(R)$ of one of the following forms:

$$\begin{aligned} (B_1 B_2 \cdots) & \quad \frac{1}{0!(k+2)!\alpha} \\ (A A B_1 B_2 \cdots) & \quad \frac{1}{2!k!\alpha} \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ (\underbrace{A \cdots A}_{k+2 \text{ or } k+1} B_1 B_2 \cdots) & \quad \frac{1}{(k+2)!0!\alpha} \text{ or } \frac{1}{(k+1)!1!\alpha} \end{aligned}$$

or

$$\begin{aligned} (A B_1 B_2 \cdots) & \quad \frac{1}{1!(k+1)!\alpha} \\ (A A A B_1 B_2 \cdots) & \quad \frac{1}{3!(k-1)!\alpha} \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ \vdots & \quad \vdots \\ (\underbrace{A \cdots A}_{k+2 \text{ or } k+1} B_1 B_2 \cdots) & \quad \frac{1}{(k+2)!0!\alpha} \text{ or } \frac{1}{(k+1)!1!\alpha}. \end{aligned}$$

Here α is the same constant, depending on the equivalence class, as before. The contribution of this equivalence class to $L_{\mu_1}(\phi)$ or $L_{\mu_1}(R)$ is

$$\begin{aligned} & \frac{1}{\alpha(k+2)!} \left[\binom{k+2}{0} + \binom{k+2}{2} + \cdots \right] \\ & = \frac{1}{\alpha(k+2)!} \left[\binom{k+2}{1} + \binom{k+2}{3} + \cdots \right] \\ & = \frac{2^{k+1}}{\alpha(k+2)!} = \frac{2^{k-1}}{\alpha k!} \frac{4}{(k+1)(k+2)}. \end{aligned}$$

That is, given any equivalence class in $B_\mu(\phi)$ or $B_\mu(R)$, there exists an equivalence class in $B_{\mu_1}(\phi)$ or $B_{\mu_1}(R)$, whose contribution to $L_{\mu_1}(\phi)$ or $L_{\mu_1}(R)$ is just $[4/(k+1)(k+2)]$ times the contribution of the original equivalence class to $L_\mu(\phi)$ or $L_\mu(R)$.

But there are no equivalence classes in $B_{\mu_1}(\phi)$ or $B_{\mu_1}(R)$ which are *not* obtained in this way, for all sets with nonzero multiplicity under μ_1 also have nonzero multiplicity under μ .

Hence

$$L_{\mu_1}(\phi) = L_\mu(\phi) \frac{4}{(k+1)(k+2)}$$

and

$$L_{\mu_1}(R) = L_\mu(R) \frac{4}{(k+1)(k+2)}$$

Since $L_\mu(\phi) = L_\mu(R)$ by hypothesis, the result is proved. End of proof of (ii); end of proof of II'.

10. FINITE-SUM EXPANSION

The method used in the last section is effective as a means of proving the generalized Griffiths' inequalities; but it is somewhat ill-adapted to the examination of examples. The reason for this is that there are infinitely many multiplicity functions μ on the subsets of N , and so any $\Sigma(R)$ is an infinite series. The following approach to the problem overcomes this difficulty. The cost of the simplification, however, is that we are unable to discover a proof of II by this method.

We have, for each $R \subset N$,

$$\begin{aligned} Z\langle\sigma^R\rangle &= \sum_\gamma (\sigma^R)_\gamma \exp\left(\beta \sum_{A \subset N} J_A (\sigma^A)_\gamma\right) \\ &= \sum_\gamma (\sigma^R)_\gamma \prod_{A \subset N} \exp(\beta J_A (\sigma^A)_\gamma) \\ &= \sum_\gamma (\sigma^R)_\gamma \prod_{A \subset N} [\cosh \beta J_A + (\sigma^A)_\gamma \sinh \beta J_A] \\ &\hspace{15em} (\text{since } \sigma^A = \pm 1) \\ &= \left(\prod_{A \subset N} \cosh \beta J_A\right) \sum_\gamma (\sigma^R)_\gamma \\ &\hspace{10em} \times \prod_{A \subset N} [1 + (\sigma^A)_\gamma \tanh \beta J_A] \\ &= \left(\prod_{A \subset N} \cosh \beta J_A\right) \sum_\gamma (\sigma^R)_\gamma \sum_{k=0}^n \sum^* (\sigma^{A_1} \dots \sigma^{A_k})_\gamma \\ &\hspace{10em} \times \tanh \beta J_{A_1} \dots \tanh \beta J_{A_k}, \end{aligned}$$

where \sum^* is the sum over all unordered k -tuples $\{A_1, \dots, A_k\}$ of *distinct* subsets of N .

Writing, for convenience,

$$K^{-1} = \prod_{A \subset N} \cosh \beta J_A > 0 \tag{10.1}$$

and

$$\tau_A = \tanh \beta J_A \quad (0 \leq \tau_A \leq 1), \tag{10.2}$$

we have

$$\begin{aligned} KZ\langle\sigma^R\rangle &= \sum_\gamma \sum_{k=0}^n \sum^* (\sigma^{A_1 \Delta \dots \Delta A_k \Delta R})_\gamma \tau_{A_1} \dots \tau_{A_k} \\ &= \sum_{k=0}^n \sum^* \tau_{A_1} \dots \tau_{A_k} \sum_\gamma (\sigma^{A_1 \Delta \dots \Delta A_k \Delta R})_\gamma. \end{aligned}$$

Now on account of (9.1), the last sum is 2^n when $A_1 \Delta \dots \Delta A_k = R$, and 0 otherwise; so

$$KZ\langle\sigma^R\rangle = 2^n \sum_{k=0}^n \sum' \tau_{A_1} \dots \tau_{A_k}, \tag{10.3}$$

where \sum' is the sum over unordered k -tuples $\{A_1, \dots, A_k\}$ of distinct subsets of N having the property that $A_1 \Delta \dots \Delta A_k = R$.

Now for each subclass $\mathcal{G} = \{A_1, \dots, A_j\}$ of the class 2^N of all subsets of N , define

$$\tau^\mathcal{G} = \prod_{A \in \mathcal{G}} \tau_A, \tag{10.4}$$

$$\Delta \mathcal{G} = A_1 \Delta \dots \Delta A_j. \tag{10.5}$$

Then we can rewrite (10.3) as

$$KZ\langle\sigma^R\rangle = 2^n \sum_{\mathcal{G}: \Delta \mathcal{G} = R} \tau^\mathcal{G}, \tag{10.6}$$

where \mathcal{G} runs over all subclasses of 2^N satisfying $\Delta \mathcal{G} = R$, i.e., over all classes of subsets of N with this property.

Finally, define

$$\alpha_R = KZ2^{-n}\langle\sigma^R\rangle = \sum_{\mathcal{G}: \Delta \mathcal{G} = R} \tau^\mathcal{G}. \tag{10.7}$$

We have

$$\alpha_\phi = KZ2^{-n},$$

and I and II are equivalent to

$$(I'') \quad \alpha_R \geq 0 \quad \text{for all } R \subset N,$$

$$(II'') \quad \alpha_\phi \alpha_{R \Delta S} \geq \alpha_R \alpha_S \quad \text{for all } R, S \subset N.$$

Again, I'' is obvious, because of (10.2). However, as mentioned earlier, we have no independent proof of II''.

As an example of the usefulness of the finite-sum approach to the problem, we prove the following theorem. This theorem is a generalization of Theorem 6 of Griffiths' paper.⁶

Theorem 10.1: Consider a generalized Ising model as in the main theorem of Sec. 2, and let

$$\mathcal{A} = \{A \subset N: J_A > 0\}.$$

⁶ R. B. Griffiths, Commun. Math. Phys. 6, 121 (1967).

(So we have

$$\alpha_R = \sum_{\substack{\mathcal{G}: \mathcal{G} \subset \mathcal{A} \\ \Delta \mathcal{G} = R}} \tau^{\mathcal{G}} \quad \text{for all } R \subset N.)$$

Then for any nonempty $R \subset N$ and any $k \in R$,

$$\langle \sigma^R \rangle \leq \sum_{\substack{S: S \in \mathcal{A} \\ k \in S}} (\tanh \beta J_S) \langle \sigma^S \sigma^R \rangle. \quad (10.8)$$

Proof: It is necessary and sufficient to prove

$$\alpha_R \leq \sum_{\substack{S: S \in \mathcal{A} \\ k \in S}} \tau_S \alpha_{RAS}. \quad (10.9)$$

The left side of (10.9) is

$$\sum_{\substack{\mathcal{G}: \mathcal{G} \subset \mathcal{A} \\ \Delta \mathcal{G} = R}} \tau^{\mathcal{G}},$$

and the right side of (10.9) is

$$\sum_{\substack{S: S \in \mathcal{A} \\ k \in S}} \sum_{\substack{\mathcal{G}: \mathcal{G} \subset \mathcal{A} \\ \Delta \mathcal{G} = R \Delta S}} \tau_S \tau^{\mathcal{G}}.$$

Now every term $\tau^{\mathcal{G}}$ appearing on the left, appears only once on the left. [We may treat the expressions in (10.9) as polynomials in the τ_A ; then each term on the left has coefficient 1.] We will prove (10.9) by choosing an arbitrary \mathcal{G} for which $\tau^{\mathcal{G}}$ appears on the left, and showing that $\tau^{\mathcal{G}}$ is also a term in the sum on the right.

We have $\mathcal{G} \subset \mathcal{A}$, $\Delta \mathcal{G} = R$, $k \in R$. Hence $\exists S \in \mathcal{G}$ with $k \in S$. Let $\mathcal{G}' = \mathcal{G} - \{S\} = \mathcal{G} \Delta \{S\}$. Then $\mathcal{G}' \subset \mathcal{A}$ (since $\mathcal{G}' \subset \mathcal{G} \subset \mathcal{A}$), and

$$\Delta \mathcal{G}' = (\Delta \mathcal{G}) \Delta (\Delta \{S\}) = R \Delta S.$$

Hence on the right appears the term

$$\tau_S \tau^{\mathcal{G}'} = \tau^{\mathcal{G}}.$$

End of proof of Theorem 10.1.

In the special case of the *classical Ising model* (see Sec. 1), the finite-sum approach gives the well-known van der Waerden expansion for the partition function Z .

First of all, since J_A (and hence τ_A) is nonzero only when A has two elements, we get

$$\alpha_\phi = 2^{-n} KZ = \sum_{k=0}^{\infty} \sum^* \tau_{i_1 i_2} \tau_{i_3 i_4} \cdots \tau_{i_{2k-1} i_{2k}}, \quad (10.10)$$

where \sum^* denotes the sum over all classes

$$\{\{i_1, i_2\}, \{i_3, i_4\}, \dots, \{i_{2k-1}, i_{2k}\}\}$$

of (distinct) 2-element subsets of N satisfying

$$\{i_1, i_2\} \Delta \{i_3, i_4\} \Delta \cdots \Delta \{i_{2k-1}, i_{2k}\} = \phi. \quad (10.11)$$

Now (10.11) is equivalent to the requirement that each index i appear in an even number of the subsets

$\{i_1, i_2\}, \dots, \{i_{2k-1}, i_{2k}\}$. Now regard the 2-element subsets of N as the edges of the complete graph on n vertices (labeled $1, 2, \dots, n$), and view the classes of 2-element subsets as subgraphs. Then (10.11) is equivalent to the condition that each vertex $1, 2, \dots, n$ have an even number of edges adjacent to it. But this in turn is equivalent to the condition that the graph be *closed*; in fact, it may be taken as the definition of a closed graph. (See Berge.⁷)

Finally, in the classical Ising model, we have the spins on a square or cubic lattice (which we now regard as a graph), and the J_{ij} given by (1.12). Letting

$$\tau = \tanh J, \quad (10.12)$$

we have

$$2^{-n} KZ = \sum_{n=0}^{\infty} \mathcal{G}(n) \tau^n, \quad (10.13)$$

where $\mathcal{G}(n)$ is the number of closed subgraphs of the lattice which has n edges. This is the van der Waerden expansion for Z . (See van der Waerden.⁸)

A probabilistic interpretation of the finite-sum approach, and in particular of the α_R given by (10.7), is as follows.

Let $h = 2^n - 1$, and label the nonempty subsets of N as A_1, A_2, \dots, A_h . Let

$$\beta_k = \frac{\tau_{A_k}}{1 + \tau_{A_k}}, \quad k = 1, 2, \dots, h.$$

Note that $0 \leq \beta_k \leq \frac{1}{2}$.

Now consider h independent trials, with the probability of success on the k th trial being β_k , the probability of failure being $1 - \beta_k$.

Define, for $k = 1, 2, \dots, h$,

$$B_k = \begin{cases} A_k & \text{if the } k\text{th trial results in success,} \\ \phi & \text{if the } k\text{th trial results in failure.} \end{cases}$$

Finally, define

$$B = B_1 \Delta B_2 \Delta \cdots \Delta B_h.$$

Then for any subset R of N , we see that

$$\begin{aligned} P(B = R) &= \prod_{k=1}^h (1 + \tau_{A_k})^{-1} \sum_{\mathcal{G}: \Delta \mathcal{G} = R} \tau^{\mathcal{G}} \\ &= \alpha_R \prod_{k=1}^h (1 + \tau_{A_k})^{-1}. \end{aligned}$$

Thus the Griffiths' inequality II is equivalent to

$$P(B = \phi) P(B = R \Delta S) \geq P(B = R) P(B = S)$$

for all $R, S \subset N$.

⁷ C. Berge, *The Theory of Graphs and its Applications*, translated by Alison Doig (John Wiley & Sons, Inc., New York, 1962).

⁸ B. L. van der Waerden, *Z. Physik* **118**, 473 (1941).

11. GROUP ALGEBRA

The following method of approaching the problem, which is due to Lenard, unifies the methods of the preceding two sections and paves the way for the methods of the next section. For reference in group algebras, see Loomis.⁹

Let G denote the group $(2^N, \Delta)$ of subsets of N under the operation of symmetric difference. We write G multiplicatively:

$$AB = A\Delta B. \tag{11.1}$$

The identity in G is ϕ ; each element A is its own inverse.

The (real) group algebra $L^1(G)$ is the set of all functions on G with values in the real numbers \mathbb{R} . (We shall customarily express the values of the functions using subscripts; thus the function α has the value α_R at $R \subset N$.) Addition and multiplication in $L^1(G)$ are defined by

$$(\alpha + \beta)_R = \alpha_R + \beta_R \tag{11.2}$$

and

$$\alpha\beta = \alpha * \beta, \tag{11.3}$$

where the convolution $\alpha * \beta$ is given by

$$(\alpha * \beta)_A = \sum_{B \in G} \alpha_B \beta_{BA^{-1}} = \sum_{B \in G} \alpha_B \beta_{AB}. \tag{11.4}$$

The additive identity in $L^1(G)$ is of course the zero function; the multiplicative identity is the function I given by

$$I_A = \begin{cases} 1 & \text{if } A = \phi \\ 0 & \text{otherwise.} \end{cases} \tag{11.5}$$

We can view the elements of $L^1(G)$ alternatively as formal linear combinations

$$\alpha = \sum_{A \in G} \alpha_A A$$

of the elements of G , with real coefficients. The usual rules for addition and multiplication of linear combinations are seen to coincide with (11.2) and (11.3). We then have the linear combination 0 as additive identity and $I = 1 \cdot \phi$ as multiplicative identity.

We shall use the notation α^{k*} for the k -fold convolution of α with itself; that is, the k th power of α in $L^1(G)$. α_R^{k*} will denote the value of α^{k*} at R , and must be distinguished from α_R^k , which is the k th power of the number α_R .

Since G is finite, a sequence of elements $\alpha^{(k)}$ of $L^1(G)$ converges to $\alpha \in L^1(G)$ iff it converges pointwise; that is,

$$\alpha^{(k)} \rightarrow \alpha \quad \text{iff} \quad \alpha_A^{(k)} \rightarrow \alpha_A \quad \text{for each } A \in G.$$

⁹ L. H. Loomis, *An Introduction to Abstract Harmonic Analysis* (D. Van Nostrand Company, Inc., Princeton, N.J., 1953).

Now G has $h = 2^n$ elements, and every function $\alpha \in L^1(G)$ is bounded, so that there exists a real number $B_\alpha \geq 0$ such that

$$|\alpha_A| \leq B_\alpha \quad \text{for all } A \in G. \tag{11.6}$$

Hence

$$|\alpha_A^{2*}| = \left| \sum_{B \in G} \alpha_{AB} \alpha_B \right| \leq h B_\alpha^2, \tag{11.7}$$

and, by induction,

$$|\alpha_A^{k*}| = \sum_{B \in G} \alpha_{AB} \alpha_B^{k*} \leq h^{k-1} B_\alpha^k. \tag{11.8}$$

Hence, for each $A \in G$, and every $\alpha \in L^1(G)$, the series

$$\sum_{k=0}^{\infty} \frac{1}{k!} \alpha_A^{k*}$$

converges (we define $\alpha^{0*} = I$), and we can define a new element $(\exp \alpha)$ of $L^1(G)$ by

$$\exp \alpha = \sum_{k=0}^{\infty} \frac{1}{k!} \alpha^{k*}. \tag{11.9}$$

Now consider the function J on G , where J_A is the same as in Sec. 2. Define

$$\pi = \exp J. \tag{11.10}$$

Then we have

$$\pi_R = (\exp J)_R = \sum_{k=0}^{\infty} \frac{1}{k!} J_R^{k*}. \tag{11.11}$$

But it can be seen that

$$J_R^{k*} = \sum_{\substack{A_1, \dots, A_k \in G \\ A_1 \cdots A_k = R}} J_{A_1} \cdots J_{A_k}, \tag{11.12}$$

so that

$$\pi_R = 2^{-n} Z\langle \sigma^R \rangle = \sum (R). \tag{11.13}$$

[See (9.2) and (9.7).]

Hence the Griffiths inequalities I' and II' of Sec. 9 (exponential-expansion approach) can be stated in terms of the function $\pi = \exp J$ in $L^1(G)$:

$$(I''') \quad \pi_R \geq 0 \quad \text{for all } R \in G;$$

$$(II''') \quad \pi_\phi \pi_{RS} \geq \pi_R \pi_S \quad \text{for all } R, S \in G.$$

We can similarly express the finite-sum approach of Sec. 10 as follows. Consider the function $\tau \in L^1(G)$, given by (10.2), and view τ as a formal linear combination:

$$\tau = \sum_{A \in G} \tau_A A. \tag{11.14}$$

Define a new formal linear combination α by

$$\alpha = \sum_{A \in G} \alpha_A A = \prod_{A \neq \phi} (1 \cdot \phi + \tau_A A). \tag{11.15}$$

The coefficient α_R of R in the right-hand member of (11.15) is clearly given by (10.7), and hence α_R is the same α_R as in Sec. 10.

Returning to the function $\pi = \exp J$, one is led to consider this function on groups other than $(2^N, \Delta)$, and ask whether Π''' holds in these cases. The following proposition shows that Π''' characterizes the group $(2^N, \Delta)$ in the sense that Π''' fails in $L^1(G)$ for any group G other than $(2^N, \Delta)$.

Theorem 11.1: If G is a group with an element a of order $r > 2$ (or of infinite order), then there exists a nonnegative function $J \in L^1(G)$ such that, if $\pi = \exp J$, then

$$\pi_1 \pi_{a^2} < (\pi_a)^2.$$

Proof: Suppose a is of finite order $r > 2$. Define

$$J_g = \begin{cases} K > 0 & \text{if } g = a \\ 0 & \text{otherwise.} \end{cases}$$

Writing π and J as formal linear combinations of elements of G , we have

$$\begin{aligned} \sum_{g \in G} \pi_g g &= \exp(Ka) = 1 + Ka + \frac{K^2 a^2}{2!} \\ &+ \cdots + \frac{K^{r-1} a^{r-1}}{(r-1)!} + \frac{K^r}{r!} + \frac{K^{r+1} a}{(r+1)!} + \cdots. \end{aligned}$$

Clearly

$$\pi_1 = 1 + \frac{K^r}{r!} + \cdots = 1 + O(K^r),$$

$$\pi_a = K + \frac{K^{r+1}}{(r+1)!} + \cdots = K + O(K^{r+1}),$$

$$\pi_{a^2} = \frac{K^2}{2} + \frac{K^{r+2}}{(r+2)!} + \cdots = \frac{K^2}{2} + O(K^{r+2}).$$

Hence

$$\begin{aligned} \frac{\pi_1 \pi_{a^2}}{(\pi_a)^2} &= \frac{\frac{1}{2}K^2 + O(K^{2r+2})}{K^2 + O(K^{2r+2})} \\ &= \frac{\frac{1}{2} + O(K^r)}{1 + O(K^r)} \rightarrow \frac{1}{2} \text{ as } K \rightarrow 0. \end{aligned}$$

So there exists a small positive value of K for which the asserted inequality holds.

Now if a is of infinite order, we simply have $\pi_1 = 1$, $\pi_a = K$, $\pi_{a^2} = \frac{1}{2}K^2$, and the asserted inequality is obvious. End of proof of Theorem 11.1.

12. FOURIER TRANSFORMS IN THE GROUP ALGEBRA

In this section we give necessary and sufficient conditions for ferromagnetism in terms of the Fourier transform of the function π of Sec. 11.

A character of a group is a homomorphism of the group into the (multiplicative) unit circle $S = \{e^{i\theta} : 0 \leq \theta < 2\pi\}$. The characters of a group G

themselves form a group \hat{G} under (pointwise) multiplication; if G is finite, then \hat{G} is isomorphic to G , although in general there is no "canonical" isomorphism.

The Fourier transform of a function $f : G \rightarrow \mathbb{R}$ is the function $\hat{f} : \hat{G} \rightarrow \mathbb{R}$ defined by

$$\hat{f}(\hat{g}) = \sum_{g \in G} \hat{g}(g) f(g) \text{ for } \hat{g} \in \hat{G}. \tag{12.1}$$

Questions of convergence naturally arise; but in our special case $G = (2^N, \Delta)$ the sum in (12.1) is a finite sum, and the topology of G need not concern us.

Now in the case $G = (2^N, \Delta)$, every element B of G (i.e., subset B of N) induces a character t_B of G defined by

$$t_B(A) = (-1)^{\#(A \cap B)}, \tag{12.2}$$

and different sets B give different characters. Hence the characters t_B exhaust \hat{G} in this case, and the correspondence $B \leftrightarrow t_B$ gives a natural isomorphism of G and \hat{G} .

So the Fourier transform of a real function f on G [that is, an element f of $L^1(G)$] is the function \hat{f} defined on \hat{G} (which we now identify with G) by

$$\hat{f}(A) = \sum_{B \subset N} (-1)^{\#(A \cap B)} f(B). \tag{12.3}$$

The Fourier inversion formula is

$$f(B) = 2^{-n} \sum_{A \subset N} (-1)^{\#(A \cap B)} \hat{f}(A). \tag{12.4}$$

The central fact we use regarding Fourier transforms is that the transform of the convolution of two functions is the (pointwise) product of the transforms of the functions. In the notation of Sec. 11, if $f, g \in L^1(G)$, then for any $R \in G$ we have

$$\widehat{(f * g)}_R = (\hat{f}_R)(\hat{g}_R). \tag{12.5}$$

Furthermore, the transform of the sum of two functions is the sum of the transforms; in our case, since convergence is pointwise in $L^1(G)$, this extends to infinite series of functions.

Now we have the functions J and π in $L^1(G)$, related by

$$\pi = \exp J = I + J + \frac{1}{2!} J^{2*} + \frac{1}{3!} J^{3*} + \cdots. \tag{12.6}$$

[Note that we can no longer allow J to assume infinite values, since we need bounded functions for convergence of the right side of (12.6). We can either modify the theory to include ∞ as a permissible value for J , or treat $J_A = \infty$ as a limiting case. We choose the latter.]

Theorem 12.1: $J_A \geq 0$ for all $A \subset N$ if the following two sets of inequalities hold:

$$(III) \quad \hat{\pi}_R > 0 \quad \text{for all } R \subset N;$$

$$(IV) \quad \prod_{\substack{A \subset N \\ \#(A \cap R) \text{ even}}} \hat{\pi}_A \geq \prod_{\substack{A \subset N \\ \#(A \cap R) \text{ odd}}} \hat{\pi}_A \quad \text{for all non-empty } R \subset N.$$

Proof: It suffices to show that the condition $J_A \geq 0$ for all $A \subset N$ implies III and is equivalent to IV.

Now the Fourier transform of I is

$$\hat{f} = 1, \tag{12.7}$$

where 1 is the function which assumes the value 1 everywhere on G . And the transform of J^{k*} is the function whose value at R is

$$(J^{k*})_R = J_R^k \tag{12.8}$$

(the k th power of the number J_R).

So taking Fourier transforms on both sides of (12.6) gives, for each $R \in G$,

$$\hat{\pi}_R = 1 + J_R + \frac{1}{2!} J_R^2 + \frac{1}{3!} J_R^3 + \dots = \exp(J_R). \tag{12.9}$$

Now (12.9) is an equation between numbers; "exp" here represents the usual exponential function. So if $J_A \geq 0$ for all A , then J is a real function, and hence \hat{f} is a real function. That is, J_R is real for all R , and so by (12.9) $\hat{\pi}_R > 0$ for all R . Thus III is proved. (Note that we have not yet used $J_A \geq 0$; III depends only on J_A being real.)

Now using (12.4), we get

$$\begin{aligned} J_B &= 2^{-n} \sum_{A \subset N} (-1)^{\#(A \cap B)} J_A \\ &= 2^{-n} \left[\sum_{\substack{A \subset N \\ \#(A \cap B) \text{ even}}} J_A - \sum_{\substack{A \subset N \\ \#(A \cap B) \text{ odd}}} J_A \right]. \end{aligned}$$

However, $J_A = \ln \hat{\pi}_A$ by (12.9); so

$$2^n J_B = \ln \left(\prod_{\substack{A \subset N \\ \#(A \cap B) \text{ even}}} \hat{\pi}_A \right) - \ln \left(\prod_{\substack{A \subset N \\ \#(A \cap B) \text{ odd}}} \hat{\pi}_A \right).$$

Hence $J_B \geq 0$ for every $B \neq \phi$ iff the above expression is nonnegative for every $B \neq \phi$; i.e., iff IV holds. End of proof of Theorem 12.1.

Note that a stronger statement can be made: III and IV for a particular $R \subset N$ are equivalent to $J_R \geq 0$ for that R .

Next we show that $\hat{\pi}$ is essentially the probability function on the set of configurations. For each subset R of N , let γ_R denote the configuration in which every spin in R has the value -1 and all others have

the value $+1$. Define

$$P_R = P(\gamma_R). \tag{12.10}$$

Again, consider β to be a factor of every J_A .

Theorem 12.2:

$$\hat{\pi}_R = Z P_R. \tag{12.11}$$

Proof:

$$P_R = Z^{-1} \exp(-\mathcal{H}_{\gamma_R}) = Z^{-1} \exp \left[\sum_{A \subset N} J_A (\sigma^A)_{\gamma_R} \right].$$

But

$$(\sigma^A)_{\gamma_R} = \begin{cases} +1 & \text{if } \#(A \cap B) \text{ is even} \\ -1 & \text{if } \#(A \cap B) \text{ is odd.} \end{cases} \tag{12.12}$$

Hence

$$Z P_R = \exp \left[\sum_{A \subset N} (-1)^{\#(A \cap B)} J_A \right] = \exp J_R = \hat{\pi}_R.$$

End of proof of Theorem 12.2.

Note that condition III is again seen to hold regardless of the nonnegativity of the J_A ; whenever they are real, they define a set of probabilities, which are nonnegative. (And if they are finite, the probabilities are nonzero.)

We can now restate Theorem 12.1 in terms of the probability distribution as follows:

Theorem 12.1': $J_A \geq 0$ for all $A \subset N$,

$$\prod_{\substack{A \subset N \\ \#(A \cap B) \text{ even}}} P_A \geq \prod_{\substack{A \subset N \\ \#(A \cap B) \text{ odd}}} P_A \quad \text{for all } B \subset N. \tag{12.13}$$

Note that (12.13) holds even for $B = \phi$, as $P_R \leq 1$ for all R .

Thus we actually have an answer to the converse problem posed in Sec. 8; but the conditions given are not in terms of the expected values $\langle \sigma^A \rangle$. At the time of writing, the problem of interpreting (12.12) or IV in terms of the $\langle \sigma^A \rangle$ is open.

We conclude with an attempt to exploit the duality between a function and its Fourier transform, and a counterexample which serves to disprove many conjectures.

We have a distribution on the space of configurations given by the function $P = Z^{-1} \hat{\pi}$, where $\hat{\pi} > 0$ is the Fourier transform of the moment function $\pi \geq 0$. It seems reasonable to guess that $\hat{\pi}$ is the moment function of some other ferromagnetic distribution. If this were the case, the probability function of the new distribution would be $\hat{P} = \hat{Z}^{-1} \pi$, where \hat{Z} is Z times a suitable norming constant.

Then we would have a set of inequalities similar to IV for the function π . That is, we can conjecture that

for every ferromagnetic system given by $J_A \geq 0$, we have

$$\prod_{\substack{A \subset N \\ \#(A \cap B) \text{ even}}} \langle \sigma^A \rangle \geq \prod_{\substack{A \subset N \\ \#(A \cap B) \text{ odd}}} \langle \sigma^A \rangle \quad \text{for all } B \subset N.$$

(Recall that $\pi_A = 2^{-n} Z \langle \sigma^A \rangle$.)

In fact, this is *false*, as the following example shows:

Example 12.1: Let $N = \{1, 2, 3\}$ and (using the form for examples given in Sec. 6)

$$x_{12} = x_{13} = 1;$$

$$x_1 = x_2 = x_3 = x_{23} = x_{123} = 1 - \epsilon$$

for some small $\epsilon > 0$.

Let $B = N$; the conjectured inequality becomes

$$\frac{\langle \sigma^\phi \rangle \langle \sigma_1 \sigma_2 \rangle \langle \sigma_1 \sigma_3 \rangle \langle \sigma_2 \sigma_3 \rangle}{\langle \sigma_1 \rangle \langle \sigma_2 \rangle \langle \sigma_3 \rangle \langle \sigma_1 \sigma_2 \sigma_3 \rangle} \geq 1.$$

(Recall that $\sigma^\phi = 1$.)

Now

$$\begin{aligned} Z^* &= 1 + 2(1 + \epsilon)^2 + 4(1 + \epsilon)^3 + (1 - \epsilon)^4 \\ &= 8 + o(1), \end{aligned}$$

$$\begin{aligned} Z^* \langle \sigma_1 \sigma_2 \rangle &= Z^* \langle \sigma_1 \sigma_3 \rangle = 1 - 2(1 - \epsilon)^2 + (1 - \epsilon)^4 \\ &= 4\epsilon^2 + o(\epsilon^2), \end{aligned}$$

$$\begin{aligned} Z^* \langle \sigma_2 \sigma_3 \rangle &= 1 + 2(1 - \epsilon)^2 - 4(1 - \epsilon)^3 + (1 - \epsilon)^4 \\ &= 4\epsilon + o(\epsilon), \end{aligned}$$

$$\begin{aligned} Z^* \langle \sigma_1 \rangle &= Z^* \langle \sigma_2 \rangle = Z^* \langle \sigma_3 \rangle = Z^* \langle \sigma_1 \sigma_2 \sigma_3 \rangle \\ &= 1 - (1 - \epsilon)^4 = 4\epsilon + o(\epsilon). \end{aligned}$$

So the ratio in question, with numerator and denominator multiplied by $(Z^*)^4$, is

$$\begin{aligned} &= \frac{[8 + o(1)][4\epsilon^2 + o(\epsilon^2)][4\epsilon^2 + o(\epsilon^2)][4\epsilon + o(\epsilon)]}{[4\epsilon + o(\epsilon)]^4} \\ &= \frac{512\epsilon^5 + o(\epsilon^5)}{256\epsilon^4 + o(\epsilon^4)} = \frac{o(\epsilon^4)}{256\epsilon^4 + o(\epsilon^4)}, \end{aligned}$$

which tends to zero as ϵ tends to zero. So the ratio is certainly less than 1 for some small positive values of ϵ . End of Example 12.1.

TABLE V. Relative probabilities for Example 12.1.

σ_1	σ_2	σ_3	Z_γ
+1	+1	+1	1
+1	+1	-1	$(1 - \epsilon)^3$
+1	-1	+1	$(1 - \epsilon)^3$
+1	-1	-1	$(1 - \epsilon)^2$
-1	+1	+1	$(1 - \epsilon)^2$
-1	+1	-1	$(1 - \epsilon)^3$
-1	-1	+1	$(1 - \epsilon)^3$
-1	-1	-1	$(1 - \epsilon)^4$

Finally, we note here the consequence of Bochner's theorem in the algebra we have described. We can consider $\hat{\pi}$ as a positive measure on G ; Bochner's theorem then provides that

$$\pi_A = \sum_{B \subset N} (-1)^{\#(A \cap B)} \hat{\pi}_B$$

defines a positive semidefinite function on G .

But it can be seen that the positive semidefiniteness of the function π is equivalent to the positive semidefiniteness of the variance-covariance matrix of the 2^n random variables σ^A . This, however, is a property of all systems of random variables, and is thus not a consequence of ferromagnetism.

ACKNOWLEDGMENTS

Much of the work reported here grew out of a seminar at Indiana University during the academic year 1966-67, in which A. Lenard and the authors participated. Much is owed to Lenard for his enthusiastic heckling during this period. The gallant efforts of W. Gustin on the group algebra approach served as further inspiration. We also wish to acknowledge useful correspondence and conversations with R. B. Griffiths.

APPENDIX: LIST OF OPEN QUESTIONS

(1) The major unsolved problem given in this paper is the converse problem of Sec. 8: Find conditions on the moments $\langle \sigma^R \rangle$ which imply that $J_A \geq 0$ for all $A \subset N$.

(2) One way of solving the converse problem would be to interpret the conditions III and IV of Sec. 12 in terms of the moments $\langle \sigma^R \rangle$. Equivalently, find a condition on the moments which is equivalent to (12.12).

(3) Give a proof of II'' (see Sec. 10) using the finite-sum approach of Sec. 10.

(4) Prove (or disprove) that the average magnetization per spin is a concave function of H (see Sec. 7), for $H \geq 0$.

(5) As a sufficient condition for concavity of magnetization in the special case of the external magnetic field Ising model of (1.13), prove that D_{ijk} [see (7.10)] is nonpositive in that system. [Example 7.3 shows that it is not always nonnegative in more general systems.]

(6) Griffiths' inequality II''' can be rephrased as follows: If G_0 is any two-element subgroup of $(2^N, \Delta)$ and AG_0 is any coset, then

$$\prod_{B \in G_0} \pi_B \geq \prod_{B \in AG_0} \pi_B.$$

For what groups and subgroups does this hold? [We know it holds for the cyclic group of order 4; and Example 12.1 shows that it does not hold for the group $(2^N, \Delta)$ and the subgroup $G_B = \{A: \#(A \cap G)$

is even}. In addition, we know that it does *not* hold for the cyclic group $\{1, a, a^2, a^3, a^4, a^5\}$ of order 6 and the subgroup $\{1, a^3\}$; the question is still open for the same group and the subgroup $\{1, a^2, a^4\}$.]

Tail of a Gravitational Wave*

W. E. COUCH AND R. J. TORRENCE
Syracuse University, Syracuse, N. Y.

AND

A. I. JANIS AND E. T. NEWMAN
University of Pittsburgh, Pittsburgh, Pa.

(Received 4 May 1967).

A first-order quadrupole sandwich wave of gravitational radiation exploding from a first-order Schwarzschild mass is examined to second order. If the second-order field preceding the sandwich wave vanishes, it is shown that the region of space-time following the sandwich wave contains a second-order, imploding quadrupole wave. The rest of the second-order field in the space-time region following the sandwich wave is also given, and it is seen to consist of monopole, quadrupole, and 16-pole nonradiative motions.

1. INTRODUCTION

In the study of gravitational radiation in a space with isolated sources two formalisms have been used. Bondi and others, working directly with the metric tensor, have been successful in understanding some important aspects of retarded radiative solutions. Thus Bondi, Van der Burg, and Metzner,¹ having expanded the metric in inverse powers of a luminosity parameter, were able to show that a pulse of outgoing radiation carries off mass from the source. Recently Bonnor and Rotenberg,² by adding a perturbation approximation to the asymptotic expansion of Bondi, were able to calculate two additional interesting effects. They showed that the emission of gravitational radiation is sometimes accompanied by a recoil of the source and that, in the quadrupole radiation \times monopole interaction, "wave tails" must exist in the sense that if the space is pure Schwarzschild before the emission of the radiation, then after the radiation stops the space cannot be stationary.

A parallel treatment of the problem of radiation

from isolated sources was initiated by Newman and Penrose³ and by Newman and Unti.⁴ Using a null-tetrad formalism, they worked with an asymptotic expansion in inverse powers of an affine parameter along the retarded light cones of the isolated source; and instead of the metric tensor, they took the physical components of the Riemann tensor as the fundamental quantities. Subsequently Janis and Newman⁵ worked out the linear theory in the NP formalism and proposed on this basis definitions of multipole moments for the exact theory. Following this, Torrence and Janis⁶ developed the second-order perturbation theory in the NP formalism, with advanced radiation solutions explicitly included. In the same paper it was shown that, although the expansion in inverse powers of an affine parameter remains a convenience, at least some second-order corrections can be calculated without an asymptotic approximation.

In the present paper a systematic method of calculating second-order corrections to first-order solutions is used, and a physically interesting result

* Supported in part by Aerospace Research Laboratories, Office of Aerospace Research, United States Air Force, and Office of Scientific Research.

¹ H. Bondi, M. Van der Burg, and A. Metzner, *Proc. Roy. Soc. (London)* **A269**, 21 (1962).

² W. B. Bonnor and M. A. Rotenberg, *Proc. Roy. Soc. (London)* **A289**, 247 (1965).

³ E. Newman and R. Penrose, *J. Math. Phys.* **3**, 566 (1962). This work will be referred to in the text as NP.

⁴ E. Newman and T. Unti, *J. Math. Phys.* **3**, 891 (1962).

⁵ A. Janis and E. Newman, *J. Math. Phys.* **6**, 902 (1965). This work will be referred to in the text as JN.

⁶ R. Torrence and A. Janis, *J. Math. Phys.* **8**, 1355 (1967). This work will be referred to in the text as TJ.

For what groups and subgroups does this hold? [We know it holds for the cyclic group of order 4; and Example 12.1 shows that it does not hold for the group $(2^N, \Delta)$ and the subgroup $G_B = \{A: \#(A \cap G)$

is even}. In addition, we know that it does *not* hold for the cyclic group $\{1, a, a^2, a^3, a^4, a^5\}$ of order 6 and the subgroup $\{1, a^3\}$; the question is still open for the same group and the subgroup $\{1, a^2, a^4\}$.]

Tail of a Gravitational Wave*

W. E. COUCH AND R. J. TORRENCE
Syracuse University, Syracuse, N. Y.

AND

A. I. JANIS AND E. T. NEWMAN
University of Pittsburgh, Pittsburgh, Pa.

(Received 4 May 1967).

A first-order quadrupole sandwich wave of gravitational radiation exploding from a first-order Schwarzschild mass is examined to second order. If the second-order field preceding the sandwich wave vanishes, it is shown that the region of space-time following the sandwich wave contains a second-order, imploding quadrupole wave. The rest of the second-order field in the space-time region following the sandwich wave is also given, and it is seen to consist of monopole, quadrupole, and 16-pole nonradiative motions.

1. INTRODUCTION

In the study of gravitational radiation in a space with isolated sources two formalisms have been used. Bondi and others, working directly with the metric tensor, have been successful in understanding some important aspects of retarded radiative solutions. Thus Bondi, Van der Burg, and Metzner,¹ having expanded the metric in inverse powers of a luminosity parameter, were able to show that a pulse of outgoing radiation carries off mass from the source. Recently Bonnor and Rotenberg,² by adding a perturbation approximation to the asymptotic expansion of Bondi, were able to calculate two additional interesting effects. They showed that the emission of gravitational radiation is sometimes accompanied by a recoil of the source and that, in the quadrupole radiation \times monopole interaction, "wave tails" must exist in the sense that if the space is pure Schwarzschild before the emission of the radiation, then after the radiation stops the space cannot be stationary.

A parallel treatment of the problem of radiation

from isolated sources was initiated by Newman and Penrose³ and by Newman and Unti.⁴ Using a null-tetrad formalism, they worked with an asymptotic expansion in inverse powers of an affine parameter along the retarded light cones of the isolated source; and instead of the metric tensor, they took the physical components of the Riemann tensor as the fundamental quantities. Subsequently Janis and Newman⁵ worked out the linear theory in the NP formalism and proposed on this basis definitions of multipole moments for the exact theory. Following this, Torrence and Janis⁶ developed the second-order perturbation theory in the NP formalism, with advanced radiation solutions explicitly included. In the same paper it was shown that, although the expansion in inverse powers of an affine parameter remains a convenience, at least some second-order corrections can be calculated without an asymptotic approximation.

In the present paper a systematic method of calculating second-order corrections to first-order solutions is used, and a physically interesting result

* Supported in part by Aerospace Research Laboratories, Office of Aerospace Research, United States Air Force, and Office of Scientific Research.

¹ H. Bondi, M. Van der Burg, and A. Metzner, *Proc. Roy. Soc. (London)* **A269**, 21 (1962).

² W. B. Bonnor and M. A. Rotenberg, *Proc. Roy. Soc. (London)* **A289**, 247 (1965).

³ E. Newman and R. Penrose, *J. Math. Phys.* **3**, 566 (1962). This work will be referred to in the text as NP.

⁴ E. Newman and T. Unti, *J. Math. Phys.* **3**, 891 (1962).

⁵ A. Janis and E. Newman, *J. Math. Phys.* **6**, 902 (1965). This work will be referred to in the text as JN.

⁶ R. Torrence and A. Janis, *J. Math. Phys.* **8**, 1355 (1967). This work will be referred to in the text as TJ.

presented. We consider weak, gravitational, 2ℓ -pole radiation exploding from a small Schwarzschild mass. In the particular case where the retarded 2ℓ -pole radiation is a sandwich wave, interest has been expressed in the nature of the field after the original wave has passed—that is, in the existence and nature of wave “tails.” It has been shown by other authors^{2,7} that, in general, this field is not stationary. It is shown in this paper that in second order the most interesting part of this field is an *advanced* 2ℓ -pole wave imploding from the first-order outgoing wave to the Schwarzschild mass. This wave is associated with a multipole moment which is $O(v^{-\ell})$ in the advanced time v ,⁸ and the wave profile is the $(2 + \ell)$ th derivative of the moment. The case $\ell = 2$ is considered in detail, and it is shown that the remainder of the wave “tail” for this case consists of a set of nonradiative motions,¹ which are monopole, quadrupole, and 16-pole in their angular dependence. The imploding wave, whose associated moment involves an integral over the entire outgoing wave, arises from the mass \times radiation interaction, and we interpret it as a back-scattering, or reflection, of the outgoing wave by the curvature of the Schwarzschild space. The nonradiative motions are due to the self-interaction of the quadrupole radiation and seem interpretable as changes in the source as a consequence of the radiation emitted. For example, the monopole part is the well-known Bondi mass loss effect,¹ while the failure of a dipole part to appear testifies to the nonexistence of a source recoil in the $\ell = 2$ case.² Our work represents an advance over previous work in that, for the first time (we believe), an exact second-order wave tail is calculated and can be given a reasonable physical interpretation in terms of linear concepts.

This paper is organized in the following way. Section 2 gives some necessary background on the null-tetrad formalism. Section 3 presents a systematic method of calculating second-order corrections, and discusses the second-order initial data problem. In Sec. 4 the mass \times radiation interaction is analyzed in detail and the radiation \times radiation interaction is considered. The main results of the paper are summarized in a short concluding section.

2. NULL-TETRAD FORMALISM

We use throughout this paper the spin-coefficient formalism of NP,³ and give a brief review of it in this section.

Consider a tetrad of basis vectors $\ell_\mu, n_\mu, m_\mu,$ and \bar{m}_μ

⁷ E. Newman and R. Penrose, Phys. Rev. Letters 15, 231 (1965).
⁸ We define $f(x) = O[g(x)]$ to mean that there exist A and a independent of x such that $f(x) \leq Ag(x)$ for $a < x$.

in the four-dimensional Riemannian space satisfying⁹

$$\ell^\mu n_\mu = -m^\mu \bar{m}_\mu = 1, \quad \ell^\mu \ell_\mu = n^\mu n_\mu = m^\mu m_\mu = \ell^\mu m_\mu = n^\mu m_\mu = 0. \quad (2.1)$$

The vectors m_μ and \bar{m}_μ are defined by $m_\mu = (1/\sqrt{2})(a_\mu - ib_\mu)$ where a_μ and b_μ are real, unit, mutually orthogonal spacelike vectors, and the bar denotes complex conjugation. It follows from Eq. (2.1) that the metric is given in terms of the tetrad by

$$g_{\mu\nu} = \ell_\mu n_\nu + \ell_\nu n_\mu - m_\mu \bar{m}_\nu - m_\nu \bar{m}_\mu. \quad (2.2)$$

After certain definitions are made, a set of partial differential equations are derived in NP which are equivalent to the Einstein field equations. The needed definitions are: intrinsic derivatives,

$$D\Phi \equiv \Phi_{;\mu}\ell^\mu, \quad \Delta\Phi \equiv \Phi_{;\mu}n^\mu, \quad \delta\Phi \equiv \Phi_{;\mu}m^\mu, \quad \bar{\delta}\Phi \equiv \Phi_{;\mu}\bar{m}^\mu, \quad (2.3)$$

for an arbitrary scalar Φ ; combinations of Ricci rotation coefficients (called spin-coefficients),

$$\kappa \equiv \ell_{\mu;\nu}m^\mu\ell^\nu, \quad (2.4a)$$

$$\pi \equiv -n_{\mu;\nu}\bar{m}^\mu\ell^\nu, \quad (2.4b)$$

$$\epsilon \equiv \frac{1}{2}(\ell_{\mu;\nu}n^\mu\ell^\nu - m_{\mu;\nu}\bar{m}^\mu\ell^\nu), \quad (2.4c)$$

$$\rho \equiv \ell_{\mu;\nu}m^\mu\bar{m}^\nu, \quad (2.4d)$$

$$\lambda \equiv -n_{\mu;\nu}\bar{m}^\mu\bar{m}^\nu, \quad (2.4e)$$

$$\alpha \equiv \frac{1}{2}(\ell_{\mu;\nu}n^\mu\bar{m}^\nu - m_{\mu;\nu}\bar{m}^\mu\bar{m}^\nu), \quad (2.4f)$$

$$\sigma \equiv \ell_{\mu;\nu}m^\mu m^\nu, \quad (2.4g)$$

$$\mu \equiv -n_{\mu;\nu}\bar{m}^\mu m^\nu, \quad (2.4h)$$

$$\beta \equiv \frac{1}{2}(\ell_{\mu;\nu}n^\mu m^\nu - m_{\mu;\nu}\bar{m}^\mu m^\nu), \quad (2.4i)$$

$$\nu \equiv -n_{\mu;\nu}\bar{m}^\mu n^\nu, \quad (2.4j)$$

$$\gamma \equiv \frac{1}{2}(\ell_{\mu;\nu}n^\mu n^\nu - m_{\mu;\nu}\bar{m}^\mu n^\nu), \quad (2.4k)$$

$$\tau \equiv \ell_{\mu;\nu}m^\mu n^\nu; \quad (2.4l)$$

and independent components of the Weyl tensor,

$$\Psi_0 \equiv -C_{\mu\nu\rho\sigma}\ell^\mu m^\nu \ell^\rho m^\sigma, \quad (2.5a)$$

$$\Psi_1 \equiv -C_{\mu\nu\rho\sigma}\ell^\mu n^\nu \ell^\rho m^\sigma, \quad (2.5b)$$

$$\Psi_2 \equiv -C_{\mu\nu\rho\sigma}\bar{m}^\mu n^\nu \ell^\rho m^\sigma, \quad (2.5c)$$

$$\Psi_3 \equiv -C_{\mu\nu\rho\sigma}\bar{m}^\mu n^\nu \ell^\rho n^\sigma, \quad (2.5d)$$

$$\Psi_4 \equiv -C_{\mu\nu\rho\sigma}\bar{m}^\mu n^\nu \bar{m}^\rho n^\sigma. \quad (2.5e)$$

Before writing down the field equations we impose (without loss of generality) several simplifying coordinate and tetrad conditions which are both

⁹ Tensor indices denoted by Greek letters range and sum from 0 to 3. Tensor indices denoted by lower case Latin letters range and sum from 2 to 3. The letter A used as a subscript will be understood to range from 0 to 4 unless otherwise noted. Ordinary partial differentiation is denoted in the usual way or by a comma, and covariant differentiation is denoted by a semicolon. The metric has signature -2 .

convenient and natural to adopt in dealing with radiation problems. A parameter u which labels null hypersurfaces of the hyperbolic Riemannian space by $u = \text{const}$ is taken as the timelike coordinate, i.e., $x^0 = u$. The null-tetrad vector ℓ_μ is chosen as $\ell_\mu = u_{,\mu} = \delta_\mu^0$ so that it is normal to the null hypersurfaces. It is then also geodesic. Let $x^1 = r$ denote the affine parameter along the null geodesics lying in the null surfaces, and let the x^i label the different geodesics on each hypersurface. Then we have $\ell^\mu = dx^\mu/dr = \delta_1^\mu$. The relation $g^{\mu\nu}\ell_\nu = \ell^\mu$ implies $g^{01} = 1$ and $g^{00} = g^{0i} = 0$. This form of ℓ^μ produces the following results on the spin-coefficients:

$$\begin{aligned}\kappa &= 0 \quad (\ell^\mu \text{ is geodesic}), \\ \epsilon + \bar{\epsilon} &= 0 \quad (\ell^\mu \text{ in terms of an affine parameter}), \\ \rho &= \bar{\rho} \quad (\ell^\mu \text{ is hypersurface orthogonal}), \\ \tau &= \bar{\alpha} + \beta \quad (\ell^\mu \text{ is equal to a gradient}).\end{aligned}$$

We also parallelly propagate n_μ along ℓ_μ , which gives $\pi = 0$, and m_μ and \bar{m}_μ along ℓ_μ , which gives $\epsilon - \bar{\epsilon} = 0$.

In compliance with Eqs. (2.1) the vectors n^μ , m^μ , and \bar{m}^μ take the form

$$n^\mu = \delta_0^\mu + U\delta_1^\mu + X^i\delta_i^\mu, \quad (2.6a)$$

$$m^\mu = \omega\delta_1^\mu + \xi^i\delta_i^\mu, \quad (2.6b)$$

with

$$g^{11} = 2(U - \omega\bar{\omega}), \quad (2.7a)$$

$$g^{1i} = X^i - (\xi^i\bar{\omega} + \bar{\xi}^i\omega), \quad (2.7b)$$

$$g^{ij} = -(\xi^i\bar{\xi}^j + \bar{\xi}^i\xi^j). \quad (2.7c)$$

In terms of the quantities introduced in Eqs. (2.6), Eqs. (2.3) become

$$D = \frac{\partial}{\partial r}, \quad \delta = \omega \frac{\partial}{\partial r} + \xi^i \frac{\partial}{\partial x^i}, \quad (2.8)$$

$$\Delta = U \frac{\partial}{\partial r} + \frac{\partial}{\partial u} + X^i \frac{\partial}{\partial x^i}.$$

We call the set of quantities defined by Eqs. (2.4), (2.5), and (2.6) tetrad-formalism (TF) variables.

The empty-space field equations may now be written as³

$$D\xi^i = \rho\xi^i + \sigma\bar{\xi}^i, \quad (2.9a)$$

$$D\omega = \rho\omega + \sigma\bar{\omega} - \tau, \quad (2.9b)$$

$$DX^i = \bar{\tau}\xi^i + \tau\bar{\xi}^i, \quad (2.9c)$$

$$DU = \bar{\tau}\omega + \tau\bar{\omega} - (\gamma + \bar{\gamma}), \quad (2.9d)$$

$$D\rho = \rho^2 + \sigma\bar{\sigma}, \quad (2.10a)$$

$$D\sigma = 2\rho\sigma + \Psi_0, \quad (2.10b)$$

$$D\tau = \rho\tau + \sigma\bar{\tau} + \Psi_1, \quad (2.10c)$$

$$D\alpha = \rho\alpha + \bar{\sigma}\beta, \quad (2.10d)$$

$$D\beta = \rho\beta + \sigma\alpha + \Psi_1, \quad (2.10e)$$

$$D\gamma = \tau\alpha + \bar{\tau}\beta + \Psi_2, \quad (2.10f)$$

$$D\lambda = \rho\lambda + \bar{\sigma}\mu, \quad (2.10g)$$

$$D\mu = \rho\mu + \sigma\lambda + \Psi_2, \quad (2.10h)$$

$$D\nu = \tau\lambda + \bar{\tau}\mu + \Psi_3, \quad (2.10i)$$

$$\delta X^i - \Delta\xi^i = (\mu + \bar{\gamma} - \gamma)\xi^i + \bar{\lambda}\bar{\xi}^i, \quad (2.11a)$$

$$\delta\bar{\xi}^i - \bar{\delta}\xi^i = (\bar{\beta} - \alpha)\xi^i - (\beta - \bar{\alpha})\bar{\xi}^i, \quad (2.11b)$$

$$\delta\bar{\omega} - \bar{\delta}\omega = (\bar{\beta} - \alpha)\omega + (\bar{\alpha} - \beta)\bar{\omega} + \mu - \bar{\mu}, \quad (2.11c)$$

$$\delta U - \Delta\omega = (\mu + \bar{\gamma} - \gamma)\omega + \bar{\lambda}\bar{\omega} - \bar{\nu}, \quad (2.11d)$$

$$\Delta\lambda - \bar{\delta}\nu = 2\alpha\nu + (\bar{\gamma} - 3\gamma - \mu - \bar{\mu})\lambda - \Psi_4, \quad (2.12a)$$

$$\delta\rho - \bar{\delta}\sigma = \tau\rho + (\bar{\beta} - 3\alpha)\sigma - \Psi_1, \quad (2.12b)$$

$$\delta\alpha - \bar{\delta}\beta = \mu\rho - \lambda\sigma - 2\alpha\beta + \alpha\bar{\alpha} + \beta\bar{\beta} - \Psi_2, \quad (2.12c)$$

$$\delta\lambda - \bar{\delta}\mu = \bar{\tau}\mu + (\bar{\alpha} - 3\beta)\lambda - \Psi_3, \quad (2.12d)$$

$$\delta\nu - \Delta\mu = \gamma\mu - 2\beta\nu + \bar{\gamma}\mu + \mu^2 + \lambda\bar{\lambda}, \quad (2.12e)$$

$$\delta\gamma - \Delta\beta = \tau\mu - \sigma\nu + (\mu - \gamma + \bar{\gamma})\beta + \alpha\bar{\lambda}, \quad (2.12f)$$

$$\delta\tau - \Delta\sigma = 2\tau\beta + (\bar{\gamma} - 3\gamma + \mu)\sigma + \bar{\lambda}\rho, \quad (2.12g)$$

$$\Delta\rho - \bar{\delta}\tau = (\gamma + \bar{\gamma} - \bar{\mu})\rho - 2\alpha\tau - \lambda\sigma - \Psi_2, \quad (2.12h)$$

$$\Delta\alpha - \bar{\delta}\gamma = \rho\nu - \tau\lambda - \beta\lambda + (\bar{\gamma} - \gamma - \bar{\mu})\alpha - \Psi_3. \quad (2.12i)$$

Equations (2.9) and (2.10) are called radial equations because of the absence of a u derivative; Eqs. (2.11) and (2.12) are called nonradial equations. The Bianchi identities in this tetrad formalism are

$$\begin{aligned}D\Psi_A - \bar{\delta}\Psi_{A-1} \\ = (5 - A)\rho\Psi_A - 2(3 - A)\alpha\Psi_{A-1} + (1 - A)\lambda\Psi_{A-2},\end{aligned} \quad (2.13)$$

$$\begin{aligned}\Delta\Psi_{A-1} - \delta\Psi_A \\ = (A - 1)\nu\Psi_{A-2} - 2(A - 3)\gamma\Psi_{A-1} + (A - 5)\tau\Psi_A \\ - A\mu\Psi_{A-1} + 2(A - 2)\beta\Psi_A - (A - 4)\sigma\Psi_{A+1},\end{aligned} \quad (2.14)$$

where $A = 1, 2, 3, 4$.

3. PERTURBATION THEORY

In this section the structure of the field equations and their solutions is examined in perturbation theory. We assume that each TF variable is expandable in a small parameter, and we write, for example, $\rho = {}_{(0)}\rho + {}_{(1)}\rho + {}_{(2)}\rho \cdots$, where the subscript "zero" denotes flat space. For ease of writing, in the remainder of this section we omit the subscript from those TF

variables which are the order indicated by the title of the subsection in which they occur.

A. Zeroth-Order Field

The zeroth-order field is simply flat space-time, and in null-spherical polar coordinates it is given by

$$\begin{aligned} \Psi_A &= 0, \\ r\mu = U &= -\frac{1}{2}, \quad \rho = -\frac{1}{r}, \\ \alpha = -\beta &= -\frac{\sqrt{2}}{4r} \cot \theta, \\ \xi^i &= \frac{\sqrt{2}}{2r} \left(1, \frac{i}{\sin \theta} \right), \end{aligned} \tag{3.1}$$

with all other TF variables equal to zero. (We have restricted our attention in this paper to spaces with asymptotically Euclidean topology.)

B. First-Order Field

The first-order field satisfies the linearized field equations of general relativity, and we use here a form of the linear-theory analysis given in JN and TJ. The Riemann tensor (given by the Ψ_A) is considered to be a field in flat space which satisfies the linearized form of Eqs. (2.13) and (2.14). Thus we have

$$\Psi_0 - \frac{1}{2} D\Psi_0 - \frac{1}{2r} \Psi_0 + \frac{\sqrt{2}}{2r} \delta\Psi_1 = 0, \tag{3.2}$$

$$D\Psi_1 + \frac{4}{r} \Psi_1 + \frac{\sqrt{2}}{2r} \delta\Psi_0 = 0, \tag{3.3a}$$

$$D\Psi_2 + \frac{3}{r} \Psi_2 + \frac{\sqrt{2}}{2r} \delta\Psi_1 = 0, \tag{3.3b}$$

$$D\Psi_3 + \frac{2}{r} \Psi_3 + \frac{\sqrt{2}}{2r} \delta\Psi_2 = 0, \tag{3.3c}$$

$$D\Psi_4 + \frac{1}{r} \Psi_4 + \frac{\sqrt{2}}{2r} \delta\Psi_3 = 0, \tag{3.3d}$$

$$\Psi_1 - \frac{1}{2} D\Psi_1 - \frac{1}{r} \Psi_1 + \frac{\sqrt{2}}{2r} \delta\Psi_2 = 0, \tag{3.4a}$$

$$\Psi_2 - \frac{1}{2} D\Psi_2 - \frac{3}{2r} \Psi_2 + \frac{\sqrt{2}}{2r} \delta\Psi_3 = 0, \tag{3.4b}$$

$$\Psi_3 - \frac{1}{2} D\Psi_3 - \frac{2}{r} \Psi_3 + \frac{\sqrt{2}}{2r} \delta\Psi_4 = 0, \tag{3.4c}$$

where the dot denotes $\partial/\partial u$, and δ and $\bar{\delta}$ are angular differential operators defined in Appendix A.

We make the assumption that $\Psi_0 = O(r^{-5})$ as $r \rightarrow \infty$ (and the same assumption will be made for all perturbative orders); this insures us that the space is

asymptotically flat at the corresponding null infinity³ and permits the following analysis of Eqs. (3.2), (3.3), and (3.4). Equations (3.3) are satisfied by

$$\Psi_1 = \frac{\Psi_1^0}{r^4} - \frac{\sqrt{2}}{2r^4} \int_{\infty}^r r'^3 \bar{\delta}\Psi_0 dr', \tag{3.5a}$$

$$\Psi_2 = \frac{\Psi_2^0}{r^3} - \frac{\sqrt{2}}{2r^3} \int_{\infty}^r r'^2 \bar{\delta}\Psi_1 dr', \tag{3.5b}$$

$$\Psi_3 = \frac{\Psi_3^0}{r^2} - \frac{\sqrt{2}}{2r^2} \int_{\infty}^r r' \bar{\delta}\Psi_2 dr', \tag{3.5c}$$

$$\Psi_4 = \frac{\Psi_4^0}{r} - \frac{\sqrt{2}}{2r} \int_{\infty}^r \bar{\delta}\Psi_3 dr', \tag{3.5d}$$

where Ψ_A^0 ($A = 1, 2, 3, 4$) are arbitrary functions independent of r . When Eqs. (3.5) are substituted into Eqs. (3.4) and (3.2), it is found that, respectively, they are equivalent to

$$\Psi_3^0 = -\frac{1}{2}\sqrt{2} \delta\Psi_4^0, \tag{3.6a}$$

$$\Psi_2^0 = -\frac{1}{2}\sqrt{2} \delta\Psi_3^0, \tag{3.6b}$$

$$\Psi_1^0 = -\frac{1}{2}\sqrt{2} \delta\Psi_2^0, \tag{3.6c}$$

and

$$\begin{aligned} \Psi_0 - \frac{1}{2} D\Psi_0 - \frac{1}{2r} \Psi_0 - \frac{1}{2r^5} \int_{\infty}^r r'^3 \delta\bar{\delta}\Psi_0 dr' \\ + \frac{\sqrt{2}}{2} \frac{\delta\Psi_1^0}{r^5} = 0. \end{aligned} \tag{3.7}$$

The Ψ_A can now be determined. The Ψ_A^0 ($A = 1, 2, 3, 4$) are known through Eqs. (3.6) when $\Psi_4^0(u, \theta, \phi)$ is given and Ψ_3^0, Ψ_2^0 , and Ψ_1^0 are given at one particular value of u . Next Ψ_0 is found by solving Eq. (3.7). The remaining Ψ_A are then found from Eqs. (3.5) by r integration. Thus the major equation which must be solved is Eq. (3.7)—for all the other equations are solved by straightforward r or u integrations.

In order to complete the first-order calculations it is necessary to find the rest of the TF variables. They are obtained by linearizing Eqs. (2.9) and (2.10) and integrating them. The results are

$$\rho = 0, \tag{3.8a}$$

$$\sigma = \frac{\sigma^0}{r^2} + \frac{1}{r^2} \int_{\infty}^r r'^2 \Psi_0 dr', \tag{3.8b}$$

$$\alpha = -{}_{(0)}\alpha \int_{\infty}^r \bar{\sigma} dr', \tag{3.8c}$$

$$\beta = -\bar{\alpha} + \frac{1}{r} \int_{\infty}^r r' \Psi_1 dr', \tag{3.8d}$$

$$\tau = \bar{\alpha} + \beta, \tag{3.8e}$$

$$\gamma = -r {}_{(0)}\alpha \int_{\infty}^r \frac{1}{r'} (\bar{\tau} - \tau) dr' + \int_{\infty}^r \Psi_2 dr', \tag{3.8f}$$

$$\mu = \frac{1}{r} \int_{\infty}^r r' \Psi_2^{\circ} dr', \quad (3.8g)$$

$$\lambda = -\frac{1}{2r} \int_{\infty}^r \bar{\sigma} dr' + \frac{\lambda^{\circ}}{r}, \quad (3.8h)$$

$$\nu = -\int_{\infty}^r \frac{\bar{\tau}}{2r'} dr' + \int_{\infty}^r \Psi_3^{\circ} dr', \quad (3.8i)$$

$$U = -\int_{\infty}^r (\gamma + \bar{\gamma}) dr', \quad (3.8j)$$

$$\omega = -\frac{1}{r} \int_{\infty}^r r' \tau dr' + \frac{\omega^{\circ}}{r}, \quad (3.8k)$$

$$X^i = r_{(0)} \bar{\xi}^i \int_{\infty}^r \frac{\tau}{r'} dr' + r_{(0)} \xi^i \int_{\infty}^r \frac{\bar{\tau}}{r'} dr', \quad (3.8l)$$

$$\xi^i = {}_{(0)} \bar{\xi}^i \int_{\infty}^r \sigma dr', \quad (3.8m)$$

where σ° , λ° , and ω° are functions independent of r arising from the integration of Eqs. (2.9) and (2.10). The remaining field equations, Eqs. (2.11) and (2.12), are then satisfied if

$$\lambda^{\circ} = \bar{\sigma}^{\circ}, \quad (3.9a)$$

$$\omega^{\circ} = -\frac{\sqrt{2}}{2} \bar{\delta} \sigma^{\circ}, \quad (3.9b)$$

$$\Psi_2^{\circ} - \bar{\Psi}_2^{\circ} = \frac{1}{2} (\bar{\delta}^2 \sigma^{\circ} - \delta^2 \bar{\sigma}^{\circ}), \quad (3.9c)$$

$$\Psi_3^{\circ} = \frac{\sqrt{2}}{2} \delta \bar{\sigma}^{\circ}, \quad (3.9d)$$

$$\Psi_4^{\circ} = -\bar{\sigma}^{\circ}. \quad (3.9e)$$

A statement of the first-order initial data problem, consistent with that of the full theory,⁴ can now, with the aid of Eqs. (3.9), be given. If one gives $\Psi_2^{\circ}(\theta, \phi) + \bar{\Psi}_2^{\circ}(\theta, \phi)$, $\Psi_1^{\circ}(\theta, \phi)$, and $\Psi_0^{\circ}(\theta, \phi, r)$ at a particular value of u and $\sigma^{\circ}(u, \theta, \phi)$, one can obtain all of the first-order TF variables. Information about the retarded and advanced radiation is contained in $\sigma^{\circ}(u, \theta, \phi)$ and $\Psi_0^{\circ}(r, \theta, \phi)$,⁶ respectively, while these quantities combine with the rest of the initial data to fix the nonradiative (including the stationary) part of the solution. [It should perhaps be pointed out that there is no inconsistency between Eq. (3.6b) and Eq. (3.9c), only redundancy.] To avoid angular singularities we must take $\sigma^{\circ} = \sum_{\ell=2}^{\infty} \sum_{m=-\ell}^{\ell} \sigma_{\ell m}^{\circ} {}_2 Y_{\ell m}$ (the ${}_2 Y_{\ell m}$ are defined in Appendix A); then Eq. (3.9c) prevents imaginary ("magnetic") monopoles from appearing.⁵

It is interesting to note that this exclusion of the magnetic monopoles is the only modification of the first-order Ψ_A° , as obtained from the Bianchi identities, imposed by the field equations. This is always true in

the highest perturbative order treated, i.e., the ${}_{(n)} \Psi_A^{\circ}$ may be obtained from the n th-order Bianchi identities knowing the other TF variables to order $n-1$, then the n th-order field equations will only restrict the ${}_{(n)} \Psi_A^{\circ}$ by requiring that

$${}_{(n)} \sigma^{\circ} = \sum_{\ell, m} {}_{(n)} \sigma_{\ell m}^{\circ} {}_2 Y_{\ell m}.$$

Solutions of Eqs. (3.2), (3.3), and (3.4) are given here, representing axially symmetric, first-order fields. The retarded and advanced 2^{ℓ} -pole fields are, respectively, given by

$$\Psi_A^{\circ} = 2^{(2-A)/2} K_{2-A}(\ell) {}_{2-A} Y_{\ell 0} r^{\ell-2} \times d^{\ell-2+A} [a_{\ell}(u)/r^{(\ell+3-A)}], \quad (3.10)$$

and

$$\Psi_A^{\circ} = 2^{(2-A)/2} K_{A-2}(\ell) {}_{2-A} Y_{\ell 0} r^{\ell-2} \times D^{\ell+2-A} [b_{\ell}(u+2r)/r^{(\ell-1+A)}], \quad (3.11)$$

where

$$d \equiv -2 \frac{\partial}{\partial u} + \frac{\partial}{\partial r},$$

$$K_p(\ell) \equiv \left[\frac{(\ell+p)!}{(\ell-p)!} \right]^{\frac{1}{2}},$$

$${}_s Y_{\ell 0} = 0, \quad \text{for } |s| > \ell.$$

Here $a_{\ell}(u)$ and $b_{\ell}(u+2r)$ are arbitrary functions of their arguments (for $\ell \geq 2$) and are proportional to the retarded and advanced 2^{ℓ} th-multipole moments, respectively. For $\ell=1$, a_1 and b_1 are real and constant; and for $\ell=0$, a_0 and b_0 are real and constant. That Eqs. (3.10) and (3.11) are solutions to Eqs. (3.2), (3.3), and (3.4) may be verified by direct substitution, use of Eq. (B3), and use of Appendix A (where the spin- s spherical harmonics ${}_s Y_{\ell m}$ are defined and their pertinent properties given).

Several comments should be made about the solutions. Their designation as retarded and advanced solutions, respectively, follows from the designation of u as the label of the future null cones. Nothing in the formalism forces this choice on us. If u is assumed to increase into the future, then $v \equiv u+2r$ will label past null cones, and it also will increase into the future. The second point is that, although Eqs. (3.10) and (3.11) were presented as first-order solutions, we call any Ψ_A° of this form "formally linear," no matter what its order happens to be. This terminology will prove useful in what follows.

The designation of $a_{\ell}(u)$ and $b_{\ell}(u+2r)$ as proportional to the 2^{ℓ} th-multipole moments of the retarded and advanced fields, respectively, is consistent with the moment definitions suggested in JN and with

their generalization to include advanced solutions.⁶ In addition we say that a linear retarded solution is a nonradiative motion¹ if $\dot{\sigma}^\circ = 0$, i.e., if

$$\frac{d^{\ell+1}}{du^{\ell+1}} a_\ell(u) = 0,$$

with a corresponding definition for advanced solutions. It should be noted that a given nonradiative motion can be looked upon as a special case of either a retarded or an advanced solution.

C. Second-Order Field

The perturbation theory considered here consists of finding the higher-order corrections to a given linear combination of first-order multipole fields. With all first-order quantities known, the equations governing the second-order Ψ_A are Eqs. (2.13) and (2.14), keeping only second-order terms. These equations are the same as Eqs. (3.2), (3.3), and (3.4), except that the zero right-hand sides are replaced by terms (called driving terms) involving quadratic products of known first-order quantities. The terms driving Eqs. (3.3) are, respectively, R_1 , R_2 , R_3 , and R_4 , where

$$R_1 = (-4 {}_{(1)}\alpha + {}_{(1)}\delta) {}_{(1)}\Psi_0, \quad (3.12a)$$

$$R_2 = (-2 {}_{(1)}\alpha + {}_{(1)}\delta) {}_{(1)}\Psi_1 - {}_{(1)}\lambda {}_{(1)}\Psi_0, \quad (3.12b)$$

$$R_3 = {}_{(1)}\delta {}_{(1)}\Psi_2 - 2 {}_{(1)}\lambda {}_{(1)}\Psi_1, \quad (3.12c)$$

$$R_4 = (2 {}_{(1)}\alpha + {}_{(1)}\delta) {}_{(1)}\Psi_3 - 3 {}_{(1)}\lambda {}_{(1)}\Psi_2. \quad (3.12d)$$

We do the same integrations and substitutions with the driven equations as we did for Eqs. (3.2), (3.3), and (3.4) before, and we find that the second-order counterpart of Eq. (3.7) is

$$\Psi_0 - \frac{1}{2} D \Psi_0 - \frac{1}{2r} \Psi_0 - \frac{1}{2r^5} \int_{-\infty}^r r' {}_{(1)}\delta \bar{\delta} \Psi_0 dr' + \frac{\sqrt{2}}{2r^5} \frac{\partial \Psi_1^\circ}{\partial r} = D_0, \quad (3.13)$$

where

$$D_0 \equiv \left(-{}_{(1)}U \frac{\partial}{\partial r} - {}_{(1)}X^i \frac{\partial}{\partial x^i} + 4 {}_{(1)}\gamma - {}_{(1)}\mu \right) {}_{(1)}\Psi_0 + (-4 {}_{(1)}\tau - 2 {}_{(1)}\beta + {}_{(1)}\delta) {}_{(1)}\Psi_1 + 3 {}_{(1)}\sigma {}_{(1)}\Psi_2 - \frac{\sqrt{2}}{2r^5} \int_{-\infty}^r r' {}_{(1)}\delta \bar{\delta} R_1 dr', \quad (3.14)$$

and the second-order counterparts to Eqs. (3.6) and (3.9) are⁴

$$\Psi_4^\circ = -\ddot{\sigma}^\circ, \quad (3.15a)$$

$$\Psi_3^\circ = \frac{\sqrt{2}}{2} \delta \dot{\sigma}^\circ, \quad (3.15b)$$

$$\Psi_2^\circ - \bar{\Psi}_2^\circ = \frac{1}{2} (\bar{\delta}^2 \sigma^\circ - \delta^2 \bar{\sigma}^\circ) + {}_{(1)}\bar{\sigma}^\circ {}_{(1)}\dot{\sigma}^\circ - {}_{(1)}\sigma^\circ {}_{(1)}\dot{\bar{\sigma}}^\circ, \quad (3.15c)$$

$$\Psi_2^\circ + \bar{\Psi}_2^\circ = \frac{1}{2} (\bar{\delta}^2 \dot{\sigma}^\circ + \delta^2 \dot{\bar{\sigma}}^\circ) - {}_{(1)}\sigma^\circ {}_{(1)}\ddot{\sigma}^\circ - {}_{(1)}\bar{\sigma}^\circ {}_{(1)}\ddot{\bar{\sigma}}^\circ, \quad (3.15d)$$

$$\Psi_1^\circ = -\frac{\sqrt{2}}{2} \delta \Psi_2^\circ + 2 {}_{(1)}\sigma^\circ {}_{(1)}\Psi_3^\circ, \quad (3.15e)$$

$$\omega^\circ = -\frac{\sqrt{2}}{2} \bar{\delta} \sigma^\circ, \quad (3.15f)$$

$$\lambda^\circ = \dot{\bar{\sigma}}^\circ. \quad (3.15g)$$

The second-order problem is solved, except for straightforward r integrations, when Ψ_0 is found as a solution to Eq. (3.13). We present here a general method of solving this equation. The method is applied in Sec. 4 to two special cases, and what we believe to be physically interesting results are obtained.

Consider the linear field to be a superposition of fields associated with multipoles, both retarded and advanced. Then D_0 will be a sum of quadratic terms in which one factor comes from a 2^ℓ multipole and the other from a $2^{\ell'}$ multipole. We can symbolize this statement by

$$D_0 = \sum [D_0(R^\ell \times R^{\ell'}) + D_0(R^\ell \times A^{\ell'}) + D_0(A^\ell \times A^{\ell'})],$$

where R indicates that the multipole is retarded and A indicates that the multipole is advanced. Because Eq. (3.13) is linear it is sufficient to be able to solve for that part of Ψ_0 [denoted by $\Psi_0(\ell \times \ell')$] which arises from a single prototype driving term, which we denote by $D_0(\ell \times \ell')$, leaving the retarded or advanced nature of the multipoles unspecified for the time being.

The angular dependence of $\Psi_0(\ell \times \ell')$ is easily determined. When the first-order quantities are put into Eq. (3.14), it is seen that the angular dependence of $D_0(\ell \times \ell')$ arises from the products

$${}_0Y_{\ell 0} {}_2Y_{\ell' 0}, \quad {}_2Y_{\ell 0} {}_0Y_{\ell' 0}, \quad {}_1Y_{\ell 0} {}_1Y_{\ell' 0}, \quad {}_{-1}Y_{\ell 0} {}_3Y_{\ell' 0}, \quad {}_3Y_{\ell 0} {}_{-1}Y_{\ell' 0}, \\ {}_{-2}Y_{\ell 0} {}_4Y_{\ell' 0}, \quad \text{and} \quad {}_4Y_{\ell 0} {}_{-2}Y_{\ell' 0}.$$

Each of these products has a finite expansion of the form $C_{2\ 2} Y_{20} + C_{4\ 2} Y_{40} + \dots + C_{\ell+\ell'\ 2} Y_{\ell+\ell' 0}$ if $\ell + \ell'$ is even or $C_{3\ 2} Y_{30} + C_{5\ 2} Y_{50} + \dots + C_{\ell+\ell'\ 2} Y_{\ell+\ell' 0}$ if $\ell + \ell'$ is odd. Therefore $D_0(\ell \times \ell')$ is also a finite series in ${}_2Y_{L0}$ which we write as¹⁰

$$D_0(\ell \times \ell') = \sum_{L=2,3}^{\ell+\ell'} D_{0L}(u, r) {}_2Y_{L0}. \quad (3.16)$$

¹⁰ Throughout the remainder of this section L is to be understood to range and sum over even integers 2 to $\ell + \ell'$ if $\ell + \ell'$ is even, and over odd integers 3 to $\ell + \ell'$ if $\ell + \ell'$ is odd.

Since ${}_2Y_{\ell 0}$ is an eigenfunction of $\delta\bar{\delta}$ (see Appendix A) we may eliminate the angular dependence from Eq. (3.13) by taking $\Psi_0(\ell \times \ell') = \sum_{L=2,3}^{\ell+\ell'} f_{0L} {}_2Y_{L0}$, and, since $\delta\Psi_1^\circ$ has spin weight 2 (see Appendix A), expanding $\delta\Psi_1^\circ$ as $\delta\Psi_1^\circ = \sum_{L=2,3}^{\ell+\ell'} (\delta\Psi_1^\circ)_L {}_2Y_{L0}$. We obtain

$$f_{0L} - \frac{1}{2} D f_{0L} - \frac{1}{2r} f_{0L} + \frac{L(L+1) - 2}{2r^5} \times \int_{\infty}^r r'^3 f_{0L} dr' + \frac{\sqrt{2}(\delta\Psi_1^\circ)_L}{2r^5} = D_{0L}(\ell \times \ell'). \tag{3.17}$$

We now confine our efforts to interactions driven by $D_0(R\ell \times R\ell')$. From the expressions given previously for the retarded 2ℓ -pole fields, we find the form of the retarded \times retarded driving terms to be

$$D_{0L}(R\ell \times R\ell') = \sum_{n=0}^{\ell+\ell'} \frac{h_{Ln}(u)}{r^{n+5}}, \tag{3.18}$$

where $h_{Ln}(u)$ depends on $\ell, \ell', ({}_1)a_\ell(u)$, and $({}_1)a_{\ell'}(u)$. The solution of Eq. (3.17) for a prototype $D_{0L} = h_{Ln}(u)/r^{n+5}$ is discussed in Appendix B, and the reader is referred there for the details. The term $\sqrt{2}(\delta\Psi_1^\circ)_L/2r^5$ has a first-order counterpart in the linear theory, and an examination of Eq. (3.10) for $A = 0$ enables us to obtain the piece of second-order Ψ_0 corresponding to this term. Using this information and the results of Appendix B, we define¹¹

$$A_L \equiv [K_2(L)]^2 r^{L-2} d^{L-2} [G_L(u)/r^{L+3}], \tag{3.19a}$$

and¹²

$$B_L \equiv \frac{2(-1)^{L+2}}{K_{-2}(L)} \int_{-\infty}^u \sum_{n=0}^{\ell+\ell'} \frac{2^{n-L} n!}{(n+L+2)!} \times \frac{h_{Ln}(u') du'}{(u+2r-u')^{n-L+2}} + H_L, \tag{3.20}$$

where H_L is an arbitrary function of $u+2r$, the $K_p(L)$ were defined after Eq. (3.11), and

$$\frac{d^{L-1}}{du^{L-1}} G_L(u) = (-1)^{L-1} \frac{\sqrt{2}(\delta\Psi_1^\circ)_L}{2^{L+1}}. \tag{3.19b}$$

The complete solution for the second-order Ψ_A for the interaction of retarded fields can now be given:

$$\Psi_0(\ell \times \ell') = \sum_{L=2,3}^{\ell+\ell'} 2K_{-2} Y_{L0} \left[A_L + r^{L-2} D^{L+2} \left(\frac{B_L}{r^{L-1}} \right) \right], \tag{3.21a}$$

¹¹ An integral similar to that in Eq. (3.20) is used by Bonnor in Ref. 2.

¹² The lower limit on this integral in B_L could be replaced by an arbitrary function of $u+2r$. This freedom however can be absorbed by H_L .

$$\Psi_1(\ell \times \ell') = \frac{\Psi_1^\circ}{r^4} + \sum_{L=2,3}^{\ell+\ell'} \sqrt{2} K_{-1} Y_{L0} \times \left[\frac{1}{r^4} \int_{\infty}^r r'^3 A_L dr' + r^{L-2} D^{L+1} \left(\frac{B_L}{r^L} \right) \right] + \frac{1}{r^4} \int_{\infty}^r r'^4 R_1 dr', \tag{3.21b}$$

$$\Psi_2(\ell \times \ell') = \frac{\Psi_2^\circ}{r^3} + \sum_{L=2,3}^{\ell+\ell'} {}_0Y_{L0} \left[\frac{1}{r^3} \int_{\infty}^r \frac{1}{r'^2} \int_{\infty}^{r''} r'^3 A_L dr' dr'' + r^{L-2} D^L \left(\frac{B_L}{r^{L+1}} \right) \right] + \frac{1}{r^3} \int_{\infty}^r r'^3 R_2 dr', \tag{3.21c}$$

$$\Psi_3(\ell \times \ell') = \frac{\Psi_3^\circ}{r^2} + \sum_{L=2,3}^{\ell+\ell'} \frac{1}{\sqrt{2}} K_{-1} Y_{L0} \left[\frac{1}{r^2} \int_{\infty}^r \frac{1}{r'^2} \int_{\infty}^{r''} \frac{1}{r''^2} \times \int_{\infty}^{r'''} r'^3 A_L dr' dr'' dr''' + r^{L-2} D^{L-1} \left(\frac{B_L}{r^{L+2}} \right) \right] + \frac{1}{r^2} \int_{\infty}^r r'^2 R_3 dr', \tag{3.21d}$$

$$\Psi_4(\ell \times \ell') = \frac{\Psi_4^\circ}{r} + \sum_{L=2,3}^{\ell+\ell'} \frac{1}{2} K_{-2} Y_{L0} \left[\frac{1}{r} \int_{\infty}^r \frac{1}{r'^2} \int_{\infty}^{r''} \frac{1}{r''^2} \times \int_{\infty}^{r'''} \frac{1}{r''^2} \int_{\infty}^{r''''} r'^3 A_L dr' dr'' dr''' dr'''' + r^{L-2} D^{L-2} \left(\frac{B_L}{r^{L+3}} \right) \right] + \frac{1}{r} \int_{\infty}^r r' R_4 dr', \tag{3.21e}$$

with $\Psi_4^\circ, \Psi_3^\circ, \Psi_2^\circ$, and Ψ_1° given by Eqs. (3.15). Note that the integrations indicated in Eqs. (3.21) are trivial to perform.

We now discuss some important features of the solution. Equations (3.21) contain all of the second-order initial data explicitly. The set of functions $H_L(u+2r)$ determines the incoming radiation and gives part of the nonradiative information. Clearly solutions corresponding to different choices of H_L differ by "formally linear" solutions. The news function $\sigma^\circ(u, \theta)$ may be freely given and, along with Ψ_1° and $(\Psi_2^\circ + \bar{\Psi}_2^\circ)$ at a particular time, determines [through Eqs. (3.15)] $\Psi_4^\circ, \Psi_3^\circ, \Psi_2^\circ$, and Ψ_1° . Thus the outgoing radiation and the rest of the nonradiative information is now fixed. The $L-1$ constants of integration needed in evaluating $G_L(u)$ from Eq. (3.19b) can be conveniently fixed by demanding that

$$\frac{d^k}{du^k} G_L(u) \Big|_{-\infty} = 0, \quad k = 0, 1, \dots, L-2. \tag{3.22}$$

The freedom corresponding to different choices of these constants can be absorbed by H_L .

It is of interest to note that in some cases the solution may be quite simple. For arbitrary $h_{nL}(u)$ the integrand appearing in B_L is, of course, not generally a total derivative with respect to u' ; however, for some ℓ and ℓ' and arbitrary ${}_{(1)}a_\ell(u)$ and ${}_{(1)}a_{\ell'}(u)$ it may be, and in such a case the parts of the Ψ_A not arising from H_L are finite series in $1/r$. As we see later, for $D_{0L}(R0 \times R2)$ this integrand is not a total derivative but for $D_{0L}(R2 \times R2)$ it is.

Finally, it should be noted that the methods presented in Appendix B allow us to solve the second-order problem for any type interaction term $D_{0L}(\ell \times \ell')$, not only $D_{0L}(R\ell \times R\ell')$; also these methods are applicable at higher orders because the n th-order field corresponding to f_{0L} will satisfy Eq. (3.17) with different driving terms.

IV. SECOND-ORDER CORRECTIONS

A. Mass \times Radiation Interaction

The method of calculating second-order corrections to first-order retarded solutions, presented in Sec. 3, is applied in this section. The problem being considered is that of a Schwarzschild mass of small magnitude at the focus of a retarded quadrupole wave of gravitational radiation with small amplitude. The second-order correction divides naturally into two parts: the interaction between the mass and the quadrupole radiation, denoted by $(R0 \times R2)$, and the interaction of the quadrupole radiation with itself, denoted by $(R2 \times R2)$. First the $(R0 \times R2)$ interaction is solved exactly, then the $(R2 \times R2)$ interaction is considered; in addition, the main features of the $(R0 \times R\ell)$ interaction are discussed. In this section all quantities are second order unless otherwise specified or obviously something else from their context.

The first-order solution which we wish to correct is

$$\begin{aligned}
 {}_{(1)}\Psi_0 &= \frac{q}{r^5} \left[2K_2(2) \left(\frac{\pi}{5} \right) {}_2Y_{20} \right], \\
 {}_{(1)}\Psi_1 &= \left(\frac{\dot{q}}{r^4} + \frac{2q}{r^5} \right) \left[-\frac{4}{\sqrt{2}} K_1(2) \left(\frac{\pi}{5} \right) {}_1Y_{20} \right], \\
 {}_{(1)}\Psi_2 &= \frac{m}{r^3} 2\pi {}_0Y_{00} + \left(\frac{\ddot{q}}{r^3} + \frac{3\dot{q}}{r^4} + \frac{3q}{r^5} \right) \left[4 \left(\frac{\pi}{5} \right) {}_0Y_{20} \right], \\
 {}_{(1)}\Psi_3 &= \left(\frac{\ddot{\bar{q}}}{r^2} + \frac{3\dot{\bar{q}}}{r^3} + \frac{9\bar{q}}{2r^4} + \frac{3q}{r^5} \right) \\
 &\quad \times \left[-4\sqrt{2} K_{-1}(2) \left(\frac{\pi}{5} \right) {}_{-1}Y_{20} \right], \\
 {}_{(1)}\Psi_4 &= \left(\frac{\ddot{\bar{q}}}{r} + \frac{2\dot{\bar{q}}}{r^2} + \frac{3\ddot{q}}{r^3} + \frac{3\dot{q}}{r^4} + \frac{3q}{2r^5} \right) \\
 &\quad \times \left[8K_{-2}(2) \left(\frac{\pi}{5} \right) {}_{-2}Y_{20} \right],
 \end{aligned} \tag{4.1}$$

where m is equal to the Schwarzschild mass and $q = q(u)$ is the first-order quadrupole moment associated with the radiation field. Using Eqs. (3.8), (3.9), and these Ψ_A , it is a simple matter to obtain all the other TF variables to first order. The results are given in Appendix C. If we evaluate the driving terms defined by Eqs. (3.12) and (3.14), we obtain

$$D_0(R0 \times R2) = m \left(\frac{\ddot{\bar{q}}}{r^5} + \frac{15q}{r^7} \right) \left[\frac{1}{2} K_2(2) {}_2Y_{20} \right], \tag{4.2a}$$

$$R_1(R0 \times R2) = R_2(R0 \times R2) = 0, \tag{4.2b}$$

$$R_3(R0 \times R2) = m \left(\frac{\ddot{q}}{r^5} - \frac{3\dot{\bar{q}}}{2r^6} - \frac{\bar{q}}{r^7} \right) \left[3\sqrt{2} K_{-1}(2) {}_{-1}Y_{20} \right], \tag{4.2c}$$

$$R_4(R0 \times R2) = m \left(\frac{\ddot{\bar{q}}}{r^4} + \frac{\dot{\bar{q}}}{2r^5} + \frac{\bar{q}}{2r^7} \right) \left[-6K_{-2}(2) {}_{-2}Y_{20} \right]. \tag{4.2d}$$

We now fix some of the second-order initial data. We pick $\sigma^\circ = 0$ and set $\Psi_2^\circ + \bar{\Psi}_2^\circ = \Psi_1^\circ = 0$ at $u = -\infty$, implying, through Eqs. (3.15), that

$$\sigma^\circ = \Psi_4^\circ = \Psi_3^\circ = \Psi_2^\circ = \Psi_1^\circ = 0. \tag{4.3}$$

Eq. (3.19b) gives us $[dG_2(u)/du] = 0$, so, in accordance with Eq. (3.22), we get

$$G_2 = 0. \tag{4.4a}$$

A comparison of Eq. (4.2a) and Eq. (3.18) gives us

$$\begin{aligned}
 h_{20} &= \ddot{\bar{q}}(u) \left(\frac{mK_2(2)}{4} \right), \\
 h_{21} &= 0, \\
 h_{22} &= q(u) \left(\frac{15mK_2(2)}{4} \right).
 \end{aligned} \tag{4.4b}$$

Substituting Eqs. (4.4) into Eqs. (3.19) and (3.20), we can evaluate A_2 and B_2 . If we substitute A_2 and B_2 into Eqs. (3.21) we obtain

$$\begin{aligned}
 \Psi_0(R0 \times R2) &= \left[\frac{6\dot{\bar{q}}}{r^5} + D^4 \left(\frac{I_2}{r} \right) \right] \left(\frac{m}{2} \right) \left[2K_{-2}(2) {}_2Y_{20} \right], \\
 \Psi_1(R0 \times R2) &= \left[-\frac{6}{r^5} + D^3 \left(\frac{I_2}{r^2} \right) \right] \\
 &\quad \times \left(\frac{m}{2} \right) \left[\sqrt{2} K_{-1}(2) {}_1Y_{20} \right], \\
 \Psi_2(R0 \times R2) &= \left[\frac{3}{r^5} + D^2 \left(\frac{I_2}{r^3} \right) \right] \left(\frac{m}{2} \right) ({}_0Y_{20}), \\
 \Psi_3(R0 \times R2) &= \left[-\frac{\dot{\bar{q}}}{r^4} + \frac{\bar{q}}{2r^6} + D \left(\frac{I_2}{r^4} \right) \right] \\
 &\quad \times \left(\frac{m}{2} \right) \left[2^{-\frac{1}{2}} K_1(2) {}_{-1}Y_{20} \right], \\
 \Psi_4(R0 \times R2) &= \left(\frac{\ddot{\bar{q}}}{2r^3} + \frac{\dot{\bar{q}}}{2r^4} + \frac{I_2}{r^5} \right) \left(\frac{m}{2} \right) \left[\frac{1}{2} K_2(2) {}_{-2}Y_{20} \right],
 \end{aligned} \tag{4.5}$$

with

$${}_{(1)}I_2 \equiv \int_{-\infty}^u \frac{q(u')}{(v-u')^2} du' + H_2, \quad v \equiv u + 2r. \quad (4.6)$$

Since we will be dealing with sandwich waves, we have assumed that ${}_{(1)}\dot{q}(-\infty) = 0$. The initial data has now been entirely fixed, except for H_2 , which is an arbitrary function of $u + 2r$. Two particular choices of H_2 are made below in our discussion of the interpretation of the solution.

To aid us in discussing the second-order effects accompanying first-order radiation we use the concept of "wave tails."^{1,2} If a linear radiative solution that vanishes in a certain region of space-time gives rise to a second-order correction that is nonvanishing in that same region, the second-order correction in that region will be called a *tail* of the linear radiative solution. On the other hand, terms in the second-order correction which are nonvanishing only where the first-order radiative solution is nonvanishing will be called *transient* terms. Although some transient effects are calculated here, we do not yet know how to interpret them. Our main interest is in wave tails because we shall find it possible to interpret those appearing in this paper by referring to the linear theory.

We now assume the wave to be a sandwich wave of thickness $2u_0$ and to be centered about $u = 0$, and we want to restrict our attention to the regions $u < -u_0$ and $u > +u_0$. In these regions the solution given by Eqs. (4.5) reduces to

$$\begin{aligned} \Psi_0 &= \frac{1}{2}mD^4 \left(\frac{I_2}{r}\right) [2K_{-2}(2) {}_2Y_{20}], \\ \Psi_1 &= \frac{1}{2}mD^3 \left(\frac{I_2}{r^2}\right) [\sqrt{2} K_{-1}(2) {}_1Y_{20}], \\ \Psi_2 &= \frac{1}{2}mD^2 \left(\frac{I_2}{r^3}\right) ({}_0Y_{20}), \\ \Psi_3 &= \frac{1}{2}mD \left(\frac{I_2}{r^4}\right) \left[\frac{1}{\sqrt{2}} K_1(2) {}_{-1}Y_{20}\right], \\ \Psi_4 &= \frac{1}{2}m \left(\frac{I_2}{r^5}\right) \left[\frac{1}{2}K_2(2) {}_{-2}Y_{20}\right], \end{aligned} \quad (4.7)$$

where I_2 was defined by Eq. (4.6). Thus the sandwich wave in a Schwarzschild space has a tail on either side of the sandwich, and it is given by Eq. (4.7).

It is helpful to substitute particular choices of ${}_{(1)}q(u)$ and $H_L(u + 2r)$ into Eqs. (4.7). For clarity in describing the results we shall use Penrose's pictures to help in the presentation. In Fig. 1 an "empty" Penrose picture is shown as an example. Future and past temporal infinity and spatial infinity are labeled I^+ , I^- , and I^0 , respectively. Future and past null infinity are labeled by \mathcal{J}^+ and \mathcal{J}^- , respectively.

Constant retarded and advanced time surfaces are also shown.

As an example let us choose ${}_{(1)}q(u) = Q\delta(u)$. Thus the first-order solution is simply a shock wave of magnitude Q exploding from a mass equal to m . Let us consider two possible $H_2(u + 2r)$'s: (a) $H_2(u + 2r) = 0$; (b) $H_2(u + 2r) = -Q/v^2$. The integral in $I_2(u, v)$ is easily evaluated, and we find that the solutions are

$$\begin{aligned} \Psi_A &= 0, \quad u < 0, \\ \Psi_A &= D^{4-A} \left(\frac{M_2(v)}{r^{1+A}}\right) 2^{(2-A)/2} K_{A-2}(2) {}_{2-A}Y_{20}, \quad u > 0, \end{aligned} \quad (4.8a)$$

$$\begin{aligned} \Psi_A &= -D^{4-A} \left(\frac{M_2(v)}{r^{1+A}}\right) 2^{(2-A)/2} K_{A-2}(2) {}_{2-A}Y_{20}, \quad u < 0, \\ \Psi_A &= 0, \quad u > 0, \end{aligned} \quad (4.8b)$$

where $M_2 = (Q/v^2)(m/2)$. By comparison with Eq. (3.11) we see that the nonvanishing part of the solutions, Eqs. (4.8), is simply an incoming quadrupole radiation solution with moment $\pm M_2$. The two solutions are represented pictorially in Fig. 2. The case $H_2 = 0$ is the physically realistic one. It shows that the outgoing first-order shock wave is partially reflected by the curvature of the Schwarzschild space and is back-scattered to the point from which it came. The two cases differ by a "formally linear" solution. Thus, if for $v > 0$ just the right second-order quadrupole wave was coming in from infinity, it would cancel out the back-scattered wave, and the first-order wave would be preceded by a tail instead of followed by one. The case $\frac{1}{2}mH_2 = -M_2$ is an example of such an incoming wave.

It should be emphasized that the solutions given by Eqs. (4.8) are not "formally linear" solutions. In a *part* of space-time they are indistinguishable from "formally linear" solutions, which permits an easy interpretation.

Obviously the above results are in no way dependent on the use of the δ function. For example, a

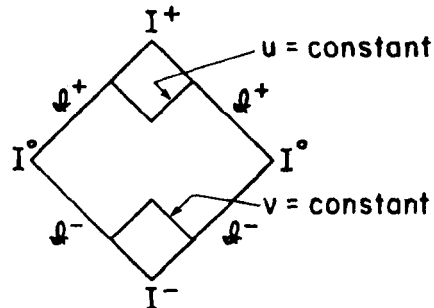


FIG. 1. Penrose Picture.

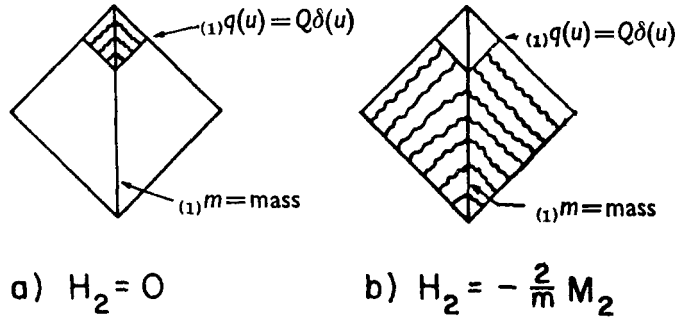


FIG. 2. Wave tails for ${}_{(1)}q(u) = Q\delta(u)$.

square wave of height Q and thickness $2u_0$ yields $M_2 = [2u_0Q/(v^2 - u_0^2)](-m/2)$.

It is possible to calculate a rather general result for mass \times radiation interactions. Using the 2^ℓ -pole retarded solutions of Eqs. (3.10) and the definition of D_0 , one can easily show that the corresponding second-order correction always gives rise to a wave tail equivalent to an advanced radiation solution, if and when the first-order radiation is shut off. The moment of the incoming wave is proportional to

$$\left[\frac{(\ell + 3)(\ell - 1)}{\ell - 2} + \frac{(\ell + 2)(\ell + 1)(\ell - 2)}{\ell(\ell - 1)} \right] \times \int_{-\infty}^{u_0} \frac{a_\ell(u')}{(v - u')^\ell} du', \quad (4.9)$$

which vanishes for no ℓ . For a wide class of $a_\ell(u')$'s the moment is clearly proportional to $1/v^\ell$ for large v , so the tail will be $O(v^{-\ell})$.

B. Radiation \times Radiation Interaction

The complete second-order correction to the first-order solution given by Eqs. (4.1) must include the $(R2 \times R2)$ interaction. The driving term $D_0(R2 \times R2)$ was evaluated in Ref. 6, and Eq. (3.15) was solved there for $\Psi_0(R2 \times R2)$. (The other Ψ_A were not calculated.) In this paper we are primarily interested in wave tails, and, having limited ourselves to a sandwich wave, we now obtain the entire set of Ψ_A for those regions of space-time outside of the sandwich wave.

In part A of this section we set

$$\sigma^\circ = (\Psi_2^\circ + \bar{\Psi}_2^\circ)|_{u=-\infty} = \Psi_1^\circ|_{u=-\infty} = 0.$$

We also assumed, for simplicity, that the sandwich wave was centered on $u = 0$ and of thickness $2u_0$. Using Eqs. (3.15a) and (3.15b), we retain from the above assumptions $\Psi_4^\circ = \Psi_3^\circ = 0$. On the other hand, Ψ_2° and Ψ_1° are now nonvanishing due to the $(R2 \times R2)$ interaction. We could evaluate $\Psi_2^\circ(u)$ and $\Psi_1^\circ(u)$ with

the aid of Eqs. (3.15c), (3.15d), and (3.15e). We could then, in principle, solve Eq. (3.13) for $\Psi_0(R2 \times R2)$, using Eq. (3.12a), but it is considerably easier to obtain $\Psi_0(R2 \times R2)$ from Ref. 6. With Ψ_4° , Ψ_3° , Ψ_2° , Ψ_1° , and Ψ_0° known, we can then use the second-order counterpart of Eqs. (3.5) to obtain all the Ψ_A . The details of this tedious but straightforward calculation are omitted and only the result given here. (It should be pointed out, however, that since we are only concerned with the regions $u < -u_0$ and $u > +u_0$, the quantities R_0 , R_1 , R_2 , and R_3 play no part in the calculation, which is a big simplification.) The result, in the notation of Eqs. (3.10), is

$$\Psi_A(R2 \times R2) = 0, \quad u < -u_0, \quad (4.10a)$$

and

$$\begin{aligned} \Psi_A(R2 \times R2) &= r^{-2}d^{A-2} \left(\frac{M_0}{r^{3-A}} \right) 2^{(2-A)/2} K_{2-A}(0) {}_{2-A}Y_{00} \\ &+ d^A \left(\frac{M_2}{r^{5-A}} \right) 2^{(2-A)/2} K_{2-A}(2) {}_{2-A}Y_{20} \\ &+ r^2d^{A+2} \left(\frac{M_4}{r^{7-A}} \right) 2^{(2-A)/2} K_{2-A}(4) {}_{2-A}Y_{40}, \quad u > u_0, \end{aligned} \quad (4.10b)$$

where

$$\begin{aligned} M_0 &= -\frac{1}{6\pi^{\frac{1}{2}}} F_2(u_0), \\ M_2 &= \frac{1}{84} \left(\frac{5}{\pi} \right)^{\frac{1}{2}} \left[\frac{1}{2}(u - u_0)^2 F_2(u_0) \right. \\ &\quad \left. + (u - u_0) F_1(u_0) + F_0(u_0) \right], \\ M_4 &= -\frac{1}{336(\pi)^{\frac{1}{2}}} \left[\frac{1}{2^{\frac{1}{4}}}(u - u_0)^4 F_2(u_0) \right. \\ &\quad \left. + \frac{1}{2}(u - u_0)^3 F_1(u_0) + \frac{1}{2}(u - u_0)^2 F_0(u_0) \right. \\ &\quad \left. + (u - u_0) F_{-1}(u_0) + F_{-2}(u_0) \right], \end{aligned} \quad (4.11)$$

with

$$\begin{aligned}
 F_2(u) &= \frac{1}{4} \int_{-u_0}^u (\ddot{q}\ddot{q} + \ddot{q}\ddot{q}) du', \\
 F_1(u) &= \int_{-u_0}^u [F_2(u') + \frac{1}{4}(5\ddot{q}\ddot{q} - \ddot{q}\ddot{q})] du', \\
 F_0(u) &= \int_{-u_0}^u [F_1(u') + \frac{3}{2}\ddot{q}\ddot{q}] du', \\
 F_{-1}(u) &= \int_{-u_0}^u [F_0(u') + 2\ddot{q}\ddot{q}] du', \\
 F_{-2}(u) &= \int_{-u_0}^u [F_{-1}(u') + \frac{5}{2}(\frac{3}{2}q\ddot{q} + q\ddot{q})] du'.
 \end{aligned} \tag{4.12}$$

Thus the ($R2 \times R2$) interaction also contributes a tail. If the tail is made to vanish before the sandwich wave (as is the case with the choice of initial data here), it takes the form given by Eq. (4.10b) after the sandwich wave. In this later region it is once again a second-order “formally linear” solution (but not a global “formally linear” solution). If we use the interpretive language of the linear theory $\Psi_0(R2 \times R2)$ can be described as a nonradiative motion characterized by nonvanishing monopole, quadrupole, and 16-pole moments given by Eqs. (4.11). The monopole moment is precisely the Bondi mass loss of the source due to the emission of the first-order wave. The vanishing dipole moment testifies to the lack of recoil when a mass emits quadrupole radiation. The quadrupole and 16-pole moments may describe the redistribution of the matter of the source as a consequence of the emission of radiation; their time dependence suggests that the radiation followed an “explosion” of the source and that the source material is spreading in time causing the moments to diverge in u as $u \rightarrow \infty$, but with no further radiation emitted.

It may be pointed out that in the third-order interaction ($R2 \times R2 \times R2$) the terms arising from ${}_{(1)}q(u) {}_{(2)}M_0$ will give rise to a radiative tail in the same manner as the ($R0 \times R2$) interaction of part A of this section. Hence, barring unexpected cancellations among terms arising from ${}_{(1)}q(u) {}_{(2)}M_0$ with other third-order terms, the third-order interaction of the quadrupole radiation with itself will show a back-scattering of the radiation.

5. CONCLUSION

We summarize the main physical results of the paper. We have shown that if a first-order mass emits a first-order gravitational, quadrupole sandwich wave, the second-order correction has the following properties.

(1) One can give initial data such that the second-order correction vanishes in that region of

space-time preceding the passage of the sandwich wave. If one does so, then

(2) the second-order correction is, in general, nonvanishing in that region of space-time following the passage of the sandwich wave.

This result is not new²; however, we have now succeeded in calculating this tail exactly in second order, and in making a sensible interpretation.

The solution mentioned in (2) is, in fact, a “formally linear” (albeit second-order) solution, and we interpret it by methods appropriate to the linear theory. We find that it consists of three parts:

(a) An arbitrary second-order retarded solution which can be, and was, eliminated from the problem by the choice of initial data.

(b) A second-order advanced radiation solution focused on the source of the original wave. This effect is proportional to the first-order mass, and thus represents a continuing partial reflection of the first-order outgoing wave by the first-order Schwarzschild curvature of the space.

(c) A second-order nonradiative motion characterized by nonvanishing monopole, quadrupole, and 16-pole moments. The monopole moment is a constant and is the mass loss of the source due to the radiation emitted. The quadrupole and 16-pole moments diverge as $u \rightarrow \infty$ and are consistent with the idea that the emission of radiation was accompanied by an explosion of the source, so that the source material is now spreading out from its original location.

An interesting and surprising property of the second-order solution is that the radiation \times radiation interaction is responsible only for the tail described in (c) and makes no contribution to that described in (b). This suggests that one set up a first-order solution representing an imploding-exploding quadrupole sandwich wave without sources, and investigate the second-order correction. Our result suggests that there may be no tail either before the incoming wave or after the outgoing wave.

ACKNOWLEDGMENTS

The authors wish to thank Professor P. G. Bergmann and Professor J. N. Goldberg for many helpful conversations during the completion of this work. Manfred Leiser and Dr. Andrzej Staruszkiewicz also participated in many useful discussions.

APPENDIX A

We give the definition of the angular differential operator δ , denoted by δ , and the spin-weighted spherical harmonics ${}_s Y_{\ell m}$, which were introduced in

Sec. 3. For more details we refer the reader to the papers where these quantities were first introduced.^{13,14}

A quantity η is said to have spin-weight s if a transformation $m'^\mu = e^{i\psi}m^\mu$ of the flat-space vector $m^\mu = (\sqrt{2}/2r)(0, 0, 1, i/\sin \theta)$ induces the transformation $\eta' = e^{is\psi}\eta$. The operator δ is then defined in terms of the spin weight of the quantity on which it operates by

$$\delta\eta = -(\sin \theta)^s \left(\frac{\partial}{\partial \theta} + \frac{i}{\sin \theta} \frac{\partial}{\partial \phi} \right) [(\sin \theta)^{-s}\eta], \quad (A1)$$

where s is the spin weight of η . We also define $\bar{\delta}$ by

$$\bar{\delta}\eta = -(\sin \theta)^{-s} \left(\frac{\partial}{\partial \theta} - \frac{i}{\sin \theta} \frac{\partial}{\partial \phi} \right) [(\sin \theta)^s\eta]. \quad (A2)$$

It follows from the definitions, Eqs. (2.5) and (2.4g), that the $(1)\Psi_A$ have spin weights $2 - A$ and that $(1)\sigma^\circ$ has spin weight 2. For the second-order quantities it can be shown that the $(2)\Psi_A$ have spin weights $2 - A$, that $(2)\sigma^\circ$ has spin weight 2, and that R_1 has spin weight 1. With this information all occurrences of δ and $\bar{\delta}$ in this paper are well defined.

The spin-weighted spherical harmonics ${}_sY_{\ell m}$ are defined by

$${}_sY_{\ell m} = \begin{cases} K_{-s}(\ell)\delta^s Y_{\ell m}, & 0 \leq s \leq \ell, \\ (-1)^s K_s(\ell)\bar{\delta}^{-s} Y_{\ell m}, & -\ell \leq s \leq 0, \end{cases} \quad (A3)$$

where the ${}_0Y_{\ell m} = Y_{\ell m}$ are the ordinary spherical harmonics and the ${}_sY_{\ell m}$ are of spin weight s . The operators δ and $\bar{\delta}$ are raising and lowering operators, respectively, for spin-weighted quantities. In particular, we have

$$\begin{aligned} \delta {}_s Y_{\ell m} &= [(\ell - s)(\ell + s + 1)]^{\frac{1}{2}} {}_{s+1} Y_{\ell m}, \\ \bar{\delta} {}_s Y_{\ell m} &= -[(\ell + s)(\ell - s + 1)]^{\frac{1}{2}} {}_{s-1} Y_{\ell m}, \end{aligned} \quad (A4)$$

from which it follows that

$$\delta \bar{\delta} {}_s Y_{\ell m} = -(\ell + s)(\ell - s + 1) {}_s Y_{\ell m}. \quad (A5)$$

The functions ${}_sY_{\ell 0}$ are familiar; they are proportional to the usual associated Legendre polynomials P_ℓ^s :

$$\begin{aligned} {}_s Y_{\ell 0} &= \frac{1}{(2\pi)^{\frac{1}{2}}} \left(\frac{2\ell + 1}{2} \right)^{\frac{1}{2}} K_{-s}(\ell) P_\ell^s, & s \geq 0, \\ {}_s Y_{\ell 0} &= \frac{(-1)^s}{(2\pi)^{\frac{1}{2}}} \left(\frac{2\ell + 1}{2} \right)^{\frac{1}{2}} K_s(\ell) P_\ell^{-s}, & s < 0. \end{aligned} \quad (A6)$$

APPENDIX B

Let the linear operator occurring in Eq. (3.17) be defined by

$$\mathfrak{L}\Phi \equiv \dot{\Phi} - \frac{1}{2}D\Phi - \frac{1}{2r}\Phi + \frac{L(L+1)-2}{2r^5} \int_\infty^r r'^3 \Phi dr'. \quad (B1)$$

If in Eq. (B1) we make the substitution

$$\Phi = r^{L-2} D^{L+2} \left(\frac{B}{r^{L-1}} \right), \quad (B2)$$

and use the identity

$$D \left[r^{L+1} D^L \left(\frac{F}{r} \right) \right] \equiv r^L D^{L+1} F, \quad (B3)$$

which holds for all integer $L \geq 0$ and for any arbitrary function F , we find that

$$\mathfrak{L}\Phi = r^{L-2} D^{L+2} \left[\frac{1}{r^{L-1}} (\dot{B} - \frac{1}{2}DB) \right], \quad (B4)$$

where the assumption is made that

$$r^{L+2} D^{L+1} \left(\frac{B}{r^L} \right) \Big|_\infty = 0. \quad (B5)$$

For simplicity we require Eq. (B5), although all that is needed is that $r^{L+2} D^{L+1} (Br^{-L})|_{r=\infty}$ be defined. For all solutions used in this paper Eq. (B5) is satisfied.

It is now clear that to solve

$$\mathfrak{L}\Phi_{Ln} = \frac{h_{Ln}(u)}{r^{n+5}}$$

it is sufficient to solve

$$\dot{B}_{Ln} - \frac{1}{2}DB_{Ln} = \frac{(-1)^{L+2}n!}{(n+L+2)!} \frac{h_{Ln}(u)}{r^{n-L+2}}. \quad (B6)$$

We note that if the independent variables are changed from u and r to u and $u + 2r = v$, then $\partial/\partial u - \frac{1}{2}(\partial/\partial r)$ becomes just $\partial/\partial u$, and the above equation can be easily integrated to yield

$$\begin{aligned} B_{Ln} &= \frac{(-1)^{L+2}2^{n-L+2}n!}{(n+L+2)!} \\ &\times \int_{-\infty}^u \frac{h_{Ln}(u') du'}{(u+2r-u')^{n-L+2}} + H_L, \end{aligned} \quad (B7)$$

where H_L is an arbitrary function of $u + 2r$. It should be noticed that B_{Ln} given by Eq. (B7) is still a solution of Eq. (B6) when the $-\infty$'s are replaced by arbitrary functions of $u + 2r$; however, we consider this freedom to be absorbed into the arbitrary function H_L . [It is conceivable that for some $H_{Ln}(u')$ the integral in Eq. (B7) may not be defined, but that the integral expression for the solution would still be possible with some lower limit on the integral other than $-\infty$.] The solution arising from the whole

¹³ E. T. Newman and R. Penrose, *J. Math. Phys.* **7**, 863 (1966).
¹⁴ J. N. Goldberg, A. J. Macfarlane, E. T. Newman, F. Rohrlich, and E. C. G. Sudarshan, *J. Math. Phys.* **8**, 2155 (1967).

driving term is then

$$2K_{-2}(L)B_L = \sum_{n=0}^{\ell+\ell'} B_{Ln}.$$

We have thus shown how to solve Eq. (3.17) for that part of the solution arising from $D_{0L}(R\ell \times R\ell')$ and have derived Eq. (3.21a).

Let us now consider the equation

$$\Omega\Phi = R, \quad (\text{B8})$$

where R is a given function of u and r . The investigation of the cases of advanced \times retarded or advanced \times advanced interactions (or even higher-order perturbation) problems can be reduced to the solving of Eq. (B8) (this was indicated in Sec. 3). The analysis which led to Eq. (B6) will lead from Eq. (B8) to

$$\dot{B}_L - \frac{1}{2}DB_L = g_L(u, r), \quad (\text{B9})$$

where g_L is a known function constructed in a definite manner from the given function R . The solution to Eq. (B9) is found in a manner similar to that of Eq. (B6) to be

$$B_L = \int_{-\infty}^u g_L(u', \frac{1}{2}(u + 2r - u')) du' + H_L, \quad (\text{B10})$$

where H_L is an arbitrary function of $u + 2r$ and the previous remarks on the lower limit $-\infty$ apply here also.

APPENDIX C

The first-order corrections to the other TF variables for a given $(1)\Psi_A$ are calculated by means of Eqs. (3.8) and (3.9). These corrections for the Ψ_A of Eqs. (4.1) are used in calculations in this paper, and are given here for completeness¹⁵:

$$\begin{aligned} (1)\rho &= 0, \\ (1)\sigma &= \left(\frac{\ddot{q}}{r^2} + \frac{3\dot{q}}{r^4}\right) \left[-\frac{1}{3}\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_2(2)_2 Y_{20}\right], \\ (1)\alpha &= \left(\frac{\dot{q}}{r^2} + \frac{\bar{q}}{r^4}\right) \left[-\left(\frac{\pi}{5}\right)^{\frac{1}{2}} \frac{1}{\sqrt{2}} K_{-1}(2)_{-1} Y_{20}\right], \end{aligned}$$

¹⁵ It should be emphasized that in some of these expressions the ${}_s Y_{\ell m}$ are serving only as complete sets, and do not necessarily represent the spin weights correctly. The quantity γ , for example, does not have a well-defined spin weight.

$$(1)\tau = \left(\frac{\dot{q}}{r^3} + \frac{4q}{3r^4}\right) \left[\left(\frac{\pi}{5}\right)^{\frac{1}{2}} \sqrt{2} K_1(2)_1 Y_{20}\right],$$

$$(1)\beta = (1)\tau - (1)\bar{\alpha},$$

$$(1)\mu = \left(\frac{\ddot{q}}{r^2} + \frac{3\dot{q}}{2r^3} + \frac{q}{r^4}\right) \left[-4\left(\frac{\pi}{5}\right)^{\frac{1}{2}} {}_0 Y_{20}\right] + \frac{m}{r} \pi^{\frac{1}{2}} {}_0 Y_{00},$$

$$\begin{aligned} (1)\gamma &= \left(\frac{\ddot{q}}{r^2} + \frac{\dot{q} + 11\dot{q}}{6r^3} + \frac{\bar{q} + 8q}{6r^4}\right) \left[-2\left(\frac{\pi}{5}\right)^{\frac{1}{2}} {}_0 Y_{20}\right] \\ &+ \left(\frac{\dot{q} - \dot{q}}{r^3} + \frac{\bar{q} - q}{r^4}\right) \left(-\frac{\pi^{\frac{1}{2}}}{6} {}_0 Y_{00}\right) \\ &+ \frac{m}{r^2} \pi^{\frac{1}{2}} {}_0 Y_{00}, \end{aligned}$$

$$(1)\lambda = \left(\frac{\ddot{q}}{r} + \frac{\dot{q}}{2r^2} + \frac{\bar{q}}{2r^4}\right) \left[-8\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_{-2}(2)_{-2} Y_{20}\right],$$

$$\begin{aligned} (1)\nu &= \left(\frac{\ddot{q}}{r} + \frac{3\dot{q}}{2r^2} + \frac{6\dot{q} - \dot{q}}{4r^3} + \frac{3q - \bar{q}}{4r^4}\right) \\ &\times \left[4\left(\frac{\pi}{5}\right)^{\frac{1}{2}} \sqrt{2} K_{-1}(2)_{-1} Y_{20}\right], \end{aligned}$$

$$\begin{aligned} (1)U &= \left(\frac{\dot{q} + \ddot{q}}{r} + \frac{\dot{q} + \dot{q}}{r^2} + \frac{q + \bar{q}}{2r^3}\right) \left[-2\left(\frac{\pi}{5}\right)^{\frac{1}{2}} {}_0 Y_{20}\right] \\ &+ \frac{m}{r} 2\pi^{\frac{1}{2}} {}_0 Y_{00}, \end{aligned}$$

$$(1)\omega = \left(\frac{\ddot{q}}{r} - \frac{3\dot{q}}{2r^2} - \frac{q}{r^3}\right) \left[-\frac{2\sqrt{2}}{3}\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_1(2)_1 Y_{20}\right],$$

$$\begin{aligned} (1)X^2 &= \left(\frac{\dot{q}}{r^3} + \frac{q}{r^4}\right) \left[-\frac{1}{3}\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_1(2)_1 Y_{20}\right] \\ &+ \left(\frac{\dot{q}}{r^3} + \frac{\bar{q}}{r^4}\right) \left[2\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_{-1}(2)_{-1} Y_{20}\right], \end{aligned}$$

$$(1)X^3 = \left(\frac{\dot{q} - \dot{q}}{r^3} + \frac{q - \bar{q}}{r^4}\right) i \left(\frac{\pi}{3}\right)^{\frac{1}{2}} {}_0 Y_{10},$$

$$(1)\xi^2 = \left(\frac{\ddot{q}}{r^2} + \frac{q}{r^4}\right) \left[\frac{1}{3\sqrt{2}}\left(\frac{\pi}{5}\right)^{\frac{1}{2}} K_2(2)_2 Y_{20}\right],$$

$$(1)\xi^3 = \left(\frac{\ddot{q}}{r^2} + \frac{q}{r^4}\right) \left[-i\left(\frac{\pi}{6}\right)^{\frac{1}{2}} K_1(1)_1 Y_{10}\right].$$