

redox choices and make it a unique organic molecule in cellular redox metabolism. Recent advances in studies on both synthetic and natural flavin analogues have begun to dissect out features which control sites of hydrogen transfer, control of two-electron vs. one-electron pathways, and modes of O₂ reductive activation as well as to define a new biochemical niche in natural gas production. With the role of flavins in bacterial bioluminescence and plant photoreceptor systems as precedent, we can anticipate that flavins may continue

to provide their own spotlight on their roles in redox biochemistry.

I thank the National Institutes of Health, the Alfred P. Sloan Foundation, and the Camille and Henry Dreyfus Foundation for support of research from my laboratory described in this article. It is also a great pleasure to thank my co-workers, both at M.I.T. and Merck, whose names appear on specific references, for their essential contributions in the design and execution of that work. Special thanks go to Fred Jacobson for the artwork in this article.

Deciphering the Protein-DNA Recognition Code

MARVIN H. CARUTHERS

Department of Chemistry, University of Colorado, Boulder, Colorado 80309

Received September 12, 1979

How do proteins interact with specific DNA sequences and as a consequence regulate the expression of genes? Our research is directed toward deciphering this recognition code. The overall objective is to determine why certain DNA sequences are the recognition sites for control proteins. In this way we may be able to learn how to regulate gene expression, to turn genes on, and to turn genes off.

This Account outlines our progress in understanding the recognition or binding of an *E. coli* protein, the *lac* repressor, for its complementary DNA, the *E. coli lac* operator. The binding of *lac* repressor to its operator site on DNA is the biological function of the protein. When the repressor is bound to the DNA, the genes involved in lactose metabolism, the *lac* operon, are not expressed. Conversely when the repressor is not bound to the operator, these genes are expressed.¹ The equilibrium association constant for the binding of repressor protein and operator DNA has been measured and is extremely large ($1 \times 10^{13} \text{ M}^{-1}$ at $I = 0.05 \text{ M}$, pH 7.4, 24 °C).² In contrast, for *E. coli* DNA deleted for the *lac* operator, the equilibrium association constant for the binding of *lac* repressor has been measured³ as $(1-3) \times 10^6 \text{ M}^{-1}$. Thus the binding of *lac* operator for its complementary DNA is highly specific. The objective of our research is to understand how *lac* repressor recognizes *lac* operator, binds specifically to this DNA, and thereby controls the expression of the *lac* operon genes. The approach has involved a series of chemical and enzymatic investigations. Initially *lac* operator DNA was synthesized by a combination of chemical and enzymatic procedures. This DNA was then modified in a sequence-specific manner and tested for altered stability of the *lac* repressor-*lac* operator (RO) complex. In this Account I will review the methodology used for synthesizing various *lac* operators, the results obtained

from measuring the stability of these modified *lac* operator-*lac* repressor complexes, and the model for the RO recognition process that was derived from these results.

The basic assumption underlying this research is that the recognition process involves a series of specific hydrophobic and hydrogen bonds between repressor and operator. We have focused our attention on deciphering those functional groups on *lac* operator that are involved in this interaction. The DNA sites are summarized in Figure 1. On a thymine-adenine base pair and in the major groove, the adenine 6-amino group is potentially a hydrogen-bond donor whereas the thymine 4-carbonyl and adenine N7 are potentially hydrogen-bond acceptors. Furthermore the thymine 5-methyl group could interact by hydrophobic bonding. The guanine-cytosine base pair also has hydrogen-bond acceptor groups (guanine N7 and 6-carbonyl) and a hydrogen-bond donor group (cytosine 4-amino) located in the major groove. In the minor groove, both base pairs have hydrogen-bond acceptor groups on the 2-carbonyl (thymine and cytosine) or the 3-nitrogen (adenine and guanine). However guanine is the only base which contains a hydrogen-bond donor located in the minor groove.

Our efforts have been directed toward probing the significance of these various functional groups as potential protein recognition sites. We have altered base pairs by inserting various analogues at specific sites in the *lac* operator sequence and then measured how these alterations affect the repressor-operator interaction. The adenine-thymine base pair has been changed by insertion of bromine or hydrogen for the thymine 5-methyl group. These analogues therefore probe a potential hydrophobic recognition site. The guanine-cytosine base pair has been modified in both major and minor grooves. The cytosine 5-hydrogen has been re-

Marvin H. Caruthers was born in Des Moines, Iowa, in 1940. He received his B.S. from Iowa State University and his Ph.D. from Northwestern University with Robert L. Letsinger. He then was a research associate with H. G. Khorana, first at the University of Wisconsin and then at the Massachusetts Institute of Technology. Since 1973 he has been on the faculty of the University of Colorado where he is now Professor of Chemistry. His research interests are in nucleic acid chemistry and biochemistry.

(1) For a comprehensive review of this system, the reader should consult "The Operon", J. H. Miller and W. S. Reznikoff, Eds., Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 1978.

(2) A. D. Riggs, H. Suzuki, and S. Bourgeois, *J. Mol. Biol.*, 48, 67 (1970).

(3) S.-Y. Lin and A. D. Riggs, *J. Mol. Biol.*, 72, 671 (1972).

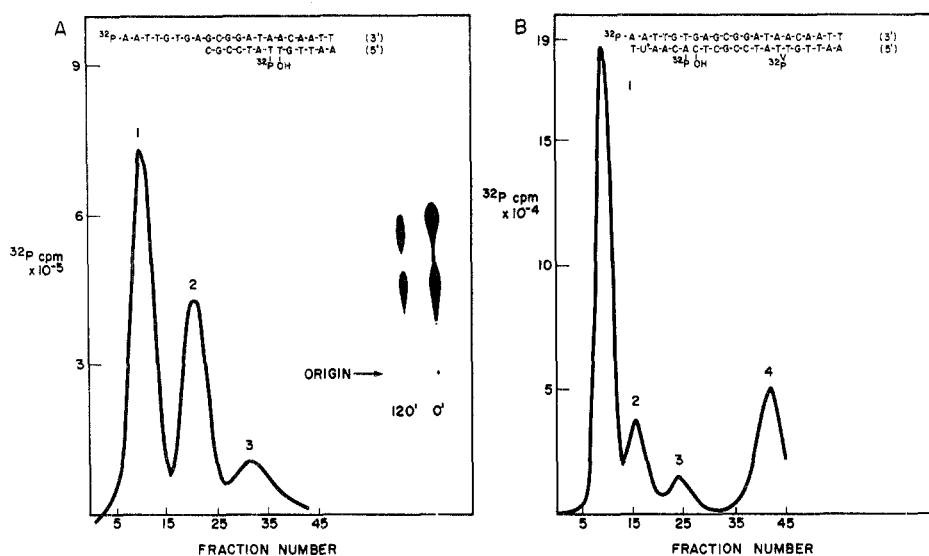


Figure 4. Synthesis of a *lac* operator containing 5-bromouracil at a specific position. The symbol U^+ designates 5-bromouracil. (Reprinted with permission from ref 21. Copyright 1978, Academic Press.)

5'-nucleotide. The initial condensation involved a reaction of $\text{d}(\text{MeOTr})\text{T}$ with $\text{d}[\text{pbzA-T}(\text{Ac})]$. The product in 62% yield was isolated by a partition procedure.¹⁶ The trinucleotide was condensed with $\text{d}[\text{panC-anC}(\text{Ac})]$ to form the pentanucleotide (40%). The pentanucleotide was condensed with $\text{d}[\text{pibG-anC}(\text{ib})]$ to give the heptanucleotide in 25% yield. Finally the undecanucleotide was formed (12%) by condensation of the heptanucleotide and $\text{d}[\text{pT-anC-bzA-anC}(\text{Ac})]$. By using various segment 6 intermediates, this sequence has been altered chemically to contain deoxycytosine and deoxyadenosine at positions 10 and 13 and deoxyguanosine at position 10.^{17,18} (See Figure 2 for the location of specific positions in segment 6.) By enzymatic manipulation of segment 6 intermediates, base analogues have been inserted into this sequence.¹⁹⁻²³ These include hypoxanthine (position 15), uracil (positions 13, 18, and 20), 5-bromouracil (positions 13, 18, and 20), 5-methylcytosine (position 13), and 5-bromocytosine (positions 10, 12, and 13). In a similar way, the remaining segments have been modified. Therefore this plan has proven to be very flexible. Presently we have synthesized in excess of 40 *lac* operator analogues using various intermediates and without initiating the re-synthesis of any segments.

The incorporation of nucleotide analogues involved considerable enzymatic manipulation primarily with T4 kinase, T4 DNA ligase, *E. coli* DNA polymerase I, and deoxynucleotidyl terminal transferase. An example summarized in Figure 4 illustrates how 5-bromouracil was inserted at position 7. Part A shows the first steps in the synthesis. Initially $[5'\text{-}^{32}\text{P}]\text{d}(\text{pT-A-T-C-C-G-C})$ was joined to $\text{d}(\text{A-A-T-T-G-T})$ by using T4 ligase. This

ligation step was completed on a template of joined segments 2 and 3. The elution profile from Sephadex G-75 is shown. The template was contained within peak 1 and the ligation product was found in peak 2. Peak 3 contained excess $[5'\text{-}^{32}\text{P}]\text{d}(\text{pT-A-T-C-C-G-C})$. The ligation product was next extended two nucleotides by using *E. coli* DNA polymerase, dTTP and dCTP. The template was joined segments 2 and 3. The inset shows the analysis of this reaction by thin-layer chromatography. The ligation product has greater mobility than the template. After 120-min reaction, the mobility of the template has not changed. However, the ligation product now has reduced mobility which indicates the addition of two nucleotides.

The final synthesis step is diagrammed in part B of Figure 4. A chemically synthesized deoxyoligonucleotide containing 5-bromouracil, $[5'\text{-}^{32}\text{P}]\text{d}(\text{pA-C-A-A-BrU-T})$, was joined to the deoxyoligonucleotide whose synthesis was outlined in part A, $\text{d}(\text{A-A-T-T-G-T-T-A-T-C-C-G-C-T-C})$. Once again T4 ligase was used to catalyze this ligation step, and joined segments 2 and 3 served as a template. The elution profile of this reaction mixture is shown in part B. The deoxyoligonucleotides corresponding to the product duplex were isolated from peak 1. Intermediates from various synthesis steps were found in the remaining peaks. Therefore the synthesis of a duplex containing 5-bromouracil at position 7 required three different enzymatic reactions and several chemically synthesized deoxyoligonucleotides. Similarly the same modified hexanucleotide, $\text{d}(\text{A-C-A-A-BrU-T})$, was used to insert 5-bromouracil at position 25.²¹ In this way one chemically synthesized intermediate containing 5-bromouracil when coupled with appropriate enzymatic steps could be used for the preparation of two 5-bromouracil-modified *lac* operator duplexes. This basic methodology involving the most efficient combination of chemical and enzymatic procedures was used for the synthesis of all modified *lac* operators.

The relative stabilities of *lac* repressor with these modified *lac* operators were determined by using the membrane filter assay.^{2,24} This assay is based on the

(24) A. D. Riggs, S. Bourgeois, and M. Cohn, *J. Mol. Biol.*, 53, 401 (1970).

(16) M. H. Caruthers and H. G. Khorana, *J. Mol. Biol.*, 72, 407 (1972).

(17) H. S. Sista, R. T. Loder, and M. H. Caruthers, *Nucleic Acids Res.*, 6, 2583 (1979).

(18) E. Brand and M. H. Caruthers, unpublished experiments.

(19) D. G. Yansura, D. V. Goeddel, D. L. Cribbs, and M. H. Caruthers, *Nucleic Acids Res.*, 4, 723 (1977).

(20) D. V. Goeddel, D. G. Yansura, and M. H. Caruthers, *Nucleic Acids Res.*, 4, 3039 (1977).

(21) D. V. Goeddel, D. G. Yansura, C. Winston, and M. H. Caruthers, *J. Mol. Biol.*, 123, 661 (1978).

(22) E. F. Fisher and M. H. Caruthers, *Nucleic Acids Res.*, 7, 401 (1979).

(23) D. G. Yansura, D. V. Goeddel, A. Kundu, and M. H. Caruthers, *J. Mol. Biol.*, 133, 117 (1979).

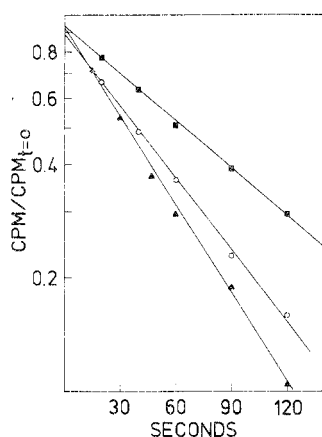


Figure 5. Dissociation kinetics of *lac* operators from *lac* repressor. Unmodified *lac* operator was the synthetic 26-base-paired operator shown in Figure 2. 5-Bromouracil-containing operators were also 26 base pairs in length. (O) Unmodified operator; (■) operator containing 5-bromouracil at position 19; (▲) operator containing 5-bromouracil at position 13. (Reprinted with permission from ref 21. Copyright 1978, Academic Press.)

observation that protein-DNA complexes are retained on nitrocellulose membranes whereas duplex DNA passes through these filters. *Lac* repressor and 5'-³²P-labeled, modified *lac* operator were mixed in approximately equimolar amounts. At time zero, unlabeled, naturally occurring *lac* operator in 50-fold molar excess was added. The decay of radioactivity bound to membranes was observed to be first order and was a measure of the stability of a repressor-modified operator complex. Different modified operators were observed to display different decay rates. We interpreted destabilization of these complexes by a specific modification as an indication of an important *lac* repressor recognition site. Examples of the results obtained by measuring dissociation rates are presented in Figure 5. In this particular case, decay kinetics for RO complexes involving base pair changes at positions 13 and 19 were examined. The insertion of adenine-5-bromouracil at position 13 destabilized the complex, whereas insertion of 5-bromouracil-adenine at position 19 (the symmetrically related position) stabilized the complex. Occasionally a modified operator forms an RO complex which is completely dissociated within 10-15 s, a rate too rapid for measurement by the membrane filtration technique. For these cases, quantitative data were obtained by using the competition equilibrium method.³ Thus when these analyses procedures were used, quantitative data for a large number of modified operators were obtained. The results are summarized in Figure 6. The numbers below the sites are the corresponding affinities of repressor for modified operators, expressed as a percent of the affinity for unmodified operator. This was possible since the kinetics of the unmodified, synthetic operator-repressor interaction have been examined in considerable detail.²⁵

The data as summarized in Figure 6 reveal considerable information on the mechanism of the *lac* operator-*lac* repressor interaction. Clearly the *lac* repressor recognizes the 5-methyl group at position 13. The evidence is as follows. Removal of the methyl group (an adenine-uracil base pair) reduced the affinity of repressor for operator to 8% of the unmodified affinity.

(25) D. V. Goeddel, D. G. Yansura, and M. H. Caruthers, *Proc. Natl. Acad. Sci. U.S.A.*, **74**, 3292 (1977).

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26								
T	G	T	G	G	A	A	T	T	G	T	G	A	G	C	G	G	A	T	A	A	C	A	A	T	T								
A	C	A	C	C	T	T	A	A	C	A	C	T	C	G	C	C	T	A	T	T	G	T	T	A	A								
													A	T	G	T	T	A	C	T													
													T	A	C	A	A	T	G	A													
													8	14	5	19	4	3	37	15													
													H	H	H	C	H	H	C	T													
													C	C	C	H	C	C	H	A													
													105	165	19	1	1	240	92														
													A	U	A U A																		
													8	70	86	90																	
A A U*U*			U* A			A U* A A			A A U*U*																								
U*U* A A			A U*			U* A U*U*			U*U* A A																								
98 88 85 170			190			60			90 120 90 95			90 90 76 85																					
				G G						C*																							
				C* C*						G																							
				80						95																							
													G																				
													C*																				
													105																				
													T																				
													A																				
													12																				

Figure 6. A summary of single-site alterations in the *lac* operator sequence. The *lac* operator sequence, including heavy lines which delineate twofold symmetric positions, is shown at the top of the figure. The dyad axis is indicated by the arrow. Various sequence changes are listed below the appropriate *lac* operator base pairs. A number located directly below each modification is the corresponding affinity of repressor for the operator, expressed as a percent of the affinity for the unmodified operator. Symbols for base analogues are as follows: H, hypoxanthine; U, uracil; U*, 5-bromouracil; C*, 5-bromocytosine; C⁰, 5-methylcytosine.

A similar loss was observed for the naturally occurring O^c mutation (a guanine-cytosine base pair). However insertion of 5-methylcytosine for cytosine stabilized the repressor-operator interaction to the same extent as observed for the naturally occurring adenine-thymine base pair. The methyl group on cytosine sterically occupies the same position as does the methyl group on thymine. Therefore the 5-methyl attached to either pyrimidine is the important functional group. The remainder of either base pair is of no consequence. We have additional data which support this conclusion. The transversion from adenine-thymine to thymine-adenine destabilized the RO interaction as much as does conversion to guanine-cytosine. This transversion shifts the methyl group approximately 13 Å and to the opposite side of the major groove. Therefore the repressor interaction with the thymine 5-methyl is positionally specific. 5-Bromouracil and 5-bromocytosine have also been inserted into the operator at position 13. Both reduced the stability of the RO complex by about the same percent, and the effect is much less dramatic than insertion of hydrogen for the methyl group. This was expected if the methyl group is interacting with repressor through hydrophobic contacts. A bromine atom, although highly polarizable, is still much less hydrophobic than a methyl group.^{26,27} Thus all these substitutions are consistent with the interpretation that the *lac* repressor recognizes the methyl group attached to the 5 position of a pyrimidine at position 13. However additional uracil and 5-bromouracil substitutions have failed to uncover any major interaction sites in-

(26) S. J. Gill and I. Wädsö, *Proc. Natl. Acad. Sci. U.S.A.*, **73**, 2955 (1976).

(27) E. Wilhelm, R. Battino, and R. Wilcox, *Chem. Rev.*, **77**, 219 (1977).

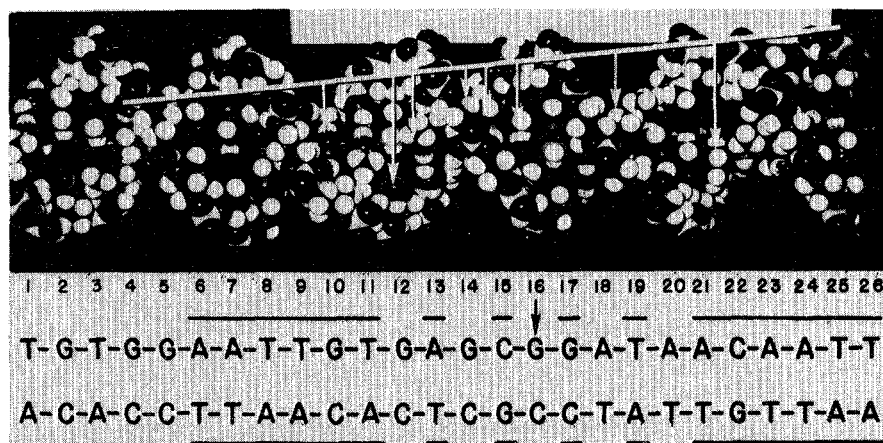


Figure 7. Postulated location of the repressor relative to the operator and location of major DNA interaction sites. The general, very schematic location of the repressor is shown by the solid, white line. White arrows extending from this line point to the DNA functional groups that interact strongly with repressor. Most of these functional groups are visible in this view (sites 10, 12, 13, 15, 16, 22). However, two sites (14, 19) are slightly hidden but readily visible when viewed from above. The black arrow marks the dyad axis. (Reprinted with permission from ref 28.)

volving the 5-methyl group.

We have also determined that the 2-amino of guanine is an important recognition site at positions 14, 15, and 16. At each of these positions, introduction of hypoxanthine caused a considerable reduction in affinity of repressor for *lac* operator. In contrast the same analogue did not dramatically alter the repressor-operator affinity when substituted at positions 10, 12, 17, and 22. These data therefore explain the biochemical significance of the O^c mutations at positions 14, 15, and 16. Each O^c mutation contains adenine which lacks the important 2-amino group. These results also reveal that the sequence symmetry is not relevant to the repressor-operator interaction. The 2-amino at position 15 is an important recognition site whereas the same group at position 17 is not recognized by repressor. Furthermore the 5-methyl of thymine at position 19 does not appear to be a recognition site whereas the symmetrically positioned methyl group at position 13 is recognized by repressor. Therefore site-specific substitution experiments with 5-bromouracil, uracil, and hypoxanthine all suggest that repressor does not recognize base pairs 13 and 19 or 15 and 17 in the same way.

Four additional repressor recognition sites have been identified by using various nucleotide analogues. However the evidence is indirect. These sites are at positions 10, 12, 19, and 22. We now know that recognition at positions 10 and 22 is through either the 6-carbonyl on guanine and/or the 4-amino on cytosine. Similarly at position 19, recognition must be through either the thymine 4-carbonyl and/or adenine 6-amino groups. At position 12 the protein appears to recognize the guanine N7. The evidence regarding these sites has been presented in detail elsewhere.^{23,28}

A model for the repressor-operator interaction follows from these data.²⁸ Strong contacts are illustrated in Figure 7. These are the 5-methyl of thymine (position 13), the 2-amino of guanine (positions 14, 15, and 16), the N7 of guanine (position 12), and the central major groove functional groups at positions 10, 19, and 22. Additionally minor but specific interactions involving thymine 5-methyl groups have been detected at posi-

tions 6, 7, 25, and 26. As can be seen by inspection of Figure 7, all these functional groups are located on the same side of the DNA. Other data from W. Gilbert's laboratory involving chemical modification of the repressor-operator complex support and expand this model. The data are summarized in Figure 2.

Ogata and Gilbert have shown that the binding of repressor to operator specifically protects five guanines and five adenines against methylation with dimethyl sulfate.^{5,6} Dimethyl sulfate methylates double-stranded DNA at the N7 of guanine and the N3 of adenine exposed, respectively, in the major and minor grooves. Therefore this experiment tells us that the repressor covers the operator and protects specific sites against alkylation. They have also completed cross-linking experiments which identify thymine residues that are covered by repressor in the major groove.⁸ All the protected sites deciphered by these investigators align on the same side of the operator as our contact sites. We therefore propose that the repressor covers the operator in the major and minor grooves at positions 6-9, the major groove at 10-13, and the minor groove at 14-16. Positions 16-22 are covered in the major groove and 23-26 in both grooves. The outer boundaries of the RO interaction are not defined by these experiments. This spacing would allow repressor to contact the operator from one side and interact with both grooves through the sites outlined above. The equilibrium constant for the interaction of *lac* repressor with *lac* operator at 24 °C and pH 7.4 is $1 \times 10^{13} \text{ M}^{-1}$. We must therefore account for -18 kcal/mol change in free energy upon formation of the RO complex. Approximately one-half this amount can be attributed to non-specific ionic interactions between the phosphodeoxyribose backbone of the operator and repressor.²⁹ The remainder, approximately -9 kcal/mol, is the nonelectrostatic component of the binding energy. We can account for this favorable binding free energy by summing the individual contributions from the sites outlined in Figure 7. Thus we feel that this model satisfactorily accounts for most of the specific interactions between repressor and operator. Quite possibly no more

(28) D. V. Goeddel, D. G. Yansura, and M. H. Caruthers, *Proc. Natl. Acad. Sci. U.S.A.*, **75**, 3578 (1978).

(29) M. D. Barkley and S. Bourgeois in "The Operon", J. H. Miller and W. S. Reznikoff, Eds., Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 1978, p 203.

than one or two additional major recognition sites remain to be detected.

Deciphering the location of specific *lac* repressor recognition sites on *lac* operator enables us to consider several other investigations. I will comment on two possibilities. We would like to know how this operator sequence relates to its function which is control of gene expression in the *lac* operon. Our present hypothesis is that the *lac* sequence is maximized for control rather than binding of *lac* repressor. If we are correct, then certain naturally occurring base-pair changes should not alter the RO interaction. For example, if the guanine N7 nitrogen is recognized by repressor, then insertion of adenine should not lead to a change in stability of the complex. The location of the N7 nitrogen is identical for both purines. As a test of this hypothesis we

are synthesizing a series of operators containing sequence changes which should not alter the stability of the interaction.

Another question we want to answer relates to the stability of the RO complex. Our hypothesis is that we can design a *lac* operator sequence that will bind repressor even more tightly than the naturally occurring sequence. Two candidate sequences are presently being synthesized. If we are successful in this venture, then we will be well on our way toward understanding how to alter gene regulatory regions in a predesigned fashion.

The work summarized herein is that of several excellent and imaginative graduate students. Their names appear in the references cited. The financial support of the National Institutes of Health, The National Science Foundation, Research Corporation, and the University of Colorado is gratefully acknowledged.