# TOWARDS GENERALIZATION AND OPTIMIZATION OF IMPLICIT METHODS

AVI LIN

*Computer Science Department, Technion-Israel Institute of Technology, Haifa, Israel 32000*

## SUMMARY

A general implicit (GI) method for solving iteratively the algebraic system arising from a finite difference approximation of an elliptic partial differential equation is formulated. Under certain assumptions this method can be reduced to the already known implicit techniques. It is shown that the GI method has a very special physical meaning when solving fluid flow problems. It is shown also how this method can be optimized to achieve the maximum rate of convergence. Finally it is shown how this new strategy is applied by solving some classical numerical fluid dynamics problems.

KEY WORDS   Iterative Methods   Implicit Procedures

## INTRODUCTION

When solving numerically elliptic or parabolic problems in several spatial dimensions, in most cases some iteration procedure has to be used in order to solve the set of linear finite difference (or finite element) algebraic equations.[1] Sometimes the inversion of the resultant matrix coefficients can be done directly,[2] but usually they are limited to maximum number of points in the field and/or to certain kinds of PDEs (for example to the Poisson like equation). Generally the iteration procedure is constructed from two basic elements: a stationary iteration method and/or accelerating method. The stationary iterative methods are commonly divided into two groups[1,3] (a) the point iterative or explicit methods and (b) the line iteration or implicit methods. The only condition for an algorithm to serve as an iteration method is that it will converge to the solution of the system. Thus, it is enough to use the iteration method as a procedure by which the algebraic system is solved. In order to minimize the CPU computer time of the iteration procedure some accelerating techniques are also used. The gradient methods are usually used to speed the rate of convergence of the numerical solution. It is well known[4] that using the gradient method with a system which is preconditioned by one of the iteration methods may produce a technique which is faster then any of the original methods.[5] Some examples of combining the conjugate gradient (CG) technique with some explicit iteration techniques (such as the Gauss–Seidel or the Jacobi methods) are given in Reference 1, and with some implicit iteration techniques (such as ADI or line relaxation methods) are given in Reference 5. Also it was shown[4] that when the CG method was preconditioned by the strongly implicit (SI) procedure[6] it produced an extremely fast rate of convergence.

It turns out that the gradient method is not always easily implemented in the iteration procedure and not always so efficient (as the system matrix $[A]$ begins to differ from symmetry, the rate of convergence of the CG method decreases). Moreover, when combining the gradient method with

stationary procedures for preconditioning, it is always better to have a stationary procedure which is comparatively fast by itself to get fast convergence from the combined procedure. That is, as the rate of convergence of the stationary technique increases, the preconditioning of the system's matrix is 'less incomplete', and a higher rate of convergence is expected. Therefore, the present study is solely concerned with the improvement of basic stationary iterative techniques and specifically of the implicit methods, since they appear to be much more stable and faster than the explicit methods.[3]

Most of the implicit methods are line methods, such as the line relaxation (LR) or the ADI techniques.[3] The SI method as well as the modified strongly implicit methods (MSI)[7] are also implicit techniques, since before every iteration the algorithm's coefficients are changed 'elliptically'. When examining the various existing implicit techniques, one might ask if it is possible that all of them fall into a more general family of techniques. The goal of the present study is to define and find such a family which will be classified under the category of the 'general implicit' (GI) methods for solving the algebraic equations arising from an elliptic system. The generalization is such that under certain assumptions this family can be reduced to any of the known implicit techniques. The basis for such a strategy is discussed in the first part of this paper. As it turns out, this method is partially dependent on the nodal algebraic equation of the field's grid points and it might be limited because of the way the iteration is executed. Consequently, in the second part, the so-called optimal general implicit technique is presented. Some stability and convergence features of these schemes will be also studied here. Most of the developments and applications will be given from the computational fluid dynamic area. Numerical results show that the GI method is better than the other methods in solving high Reynolds number flows.

## THE GENERAL IMPLICIT (GI) METHOD

### Some basic existing implicit methods

Generally, implicit methods converge faster and exhibit much more stability, than explicit methods. Usually, the variables in the implicit methods are solved implicitly along lines: rows and columns in the ADI techniques[3] and diagonals in the MSI techniques.[7] Let us name this line the 'implicit line'. Since the variables in the implicit techniques are solved in a coupled manner along the implicit line, their stability conditions (which are roughly conditions on the short wave disturbances in the field) become less severe than in the explicit methods.[3] Another feature of the implicit techniques is that the effect of the boundary conditions is felt immediately in the field, since the implicit line always runs between two boundary points, and all the equations along this line are coupled.[7] Therefore, the long wave mode of disturbances is also decreased more rapidly by the implicit techniques than by the explicit techniques.[8] The SI method does not differ from the regular implicit methods, since the algorithm's coefficients depend on the coefficients of the neighbouring points. The MSI methods, which are modifications of the SI method,[7] are also line implicit techniques. All the above description can be explained mathematically by presenting the incomplete LU decomposition for the various implicit formulations, as in Reference 8. The main idea of the present study is to try, in the spirit of the MSI technique, to generalize the line implicit techniques and to establish an optimal implicit technique which can be reduced under particular conditions to the well known existing implicit procedures. It turns out that there is no *ad hoc* proof of the observations that the present method is better than other iterative techniques which were examined. Thus the near-optimal factorization of equation (1) will not be studied in detail. The GI method will be the first method to be presented here as a preliminary study of the techniques to be described in this paper.

*The general implicit (GI) method*

For simplicity, the basic principle of the GI method will be demonstrated by considering the following one variable elliptic partial differential equation:

$$L(\Phi) = 0 \tag{1}$$

where $\Phi$ is the variable, $L$ is a second order elliptic differential operator in the spatial co-ordinates and, possibly, parabolic in the time-like direction. Equation (1) presents a model for many of the convection–diffusion phenomena in nature.[9] By discretizing equation (1) on a finite grid, which is spread over the finite spatial domain of the solution $\Omega$, the value of $\Phi$ at every grid point $p \in \Omega$ is connected algebraicaly to the $\Phi$ values at the surrounding points e, w, n and s, as shown in Figure 1. We consider only a linear operator $L$; however, the solution of the non-linear algebraic system can frequently be found by quasi-linearizing the system and solving a sequence of linear problems.[3] The linear difference equation can be written generally at the point p as follows:

$$P\Phi_p + E\Phi_e + W\Phi_w + S\Phi_s + N\Phi_n + D_p = 0 \tag{2}$$

where $P$, $E$, $W$, $N$ and $S$ are the coefficients of $\Phi_p$, $\Phi_e$, $\Phi_w$, $\Phi_s$ and $\Phi_n$, respectively centred at the point p and $D_p$ is the source term. For two dimensional fields with Dirichlet boundary conditions for $\Phi$, which has $m \times n$ discrete points, there are

$$K = (m - 2) \times (n - 2) \tag{3}$$

algebraic relations of the form of equation (2). This system of equations is generally very large and very sparse. Assuming that the finite difference algorithm is convergent, which means that there are no possible growing solutions, the goal is to solve the algebraic system given by equation (2) iteratively in a stable and rapid manner. This, as mentioned in the introduction, is the requirement that the technique be as implicit and simultaneous as possible.

*The general implicit line*

Basically, implicit methods can be formulated by defining a spatial ordering among the grid points in $\Omega$ so that the system of algebraic equations has the possibility of being solved with maximum implicitness. This ordering creates the general implicit line (GIL) which can be defined as follows:

*Definition* 1. The GIL is a continuous line in $\Omega$ through all the discrete points with the following properties:
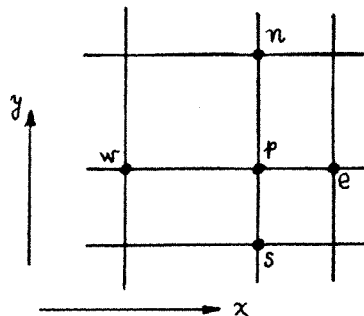


Figure 1. The elementary finite difference cell

(1) The number of intersections with itself should be minimized,

(2) It is preferable that the GIL begins and ends at different grid points or on the field boundaries. The second property is not essential from the theoretical point of view, but it was deduced from the convergence rates of some numerical experiments. Examples of two possible GILs are given in Figure 2(a) and 2(b). Generally, the GIL is constructed in such a way that if it is drawn to some inner point p, the next point to which this line is to be continued should be one of the e, w, n or s neighbouring points to p, as is demonstrated in Figure 2(c). Only three out of these four possible points need be considered, since one of them has already been used to bring this line to the point p. Therefore, the complexity in creating a GIL is less than $O(3^K)$. Once it has been decided to which point the GIL is to be continued, the other two points will be defined as the 'side points' of the central point p (or of the $p$th equation) for solving a given GIL. The other three points which are on the GIL will be defined as the 'main points'. After constructing the desired GIL, all the points of the field are numbered in the same order as they appear on the GIL. Equation (2) is then written for every point on the GIL, carrying the number of this point. The final system of equations is:

$$[A]\{\Phi\} = \{D\} \qquad (4)$$

where the system matrix coefficient $[A]$ can be written

$$[A] = [M] - [N] \qquad (5)$$

where $[M]$ is a tri-diagonal matrix built from the coefficients of the main points. Since the GIL is continuous, $[M]$ has the possibility of having filled diagonals. The $[N]$ matrix describes the contribution of the 'side points' to the balance of equation (2). Usually, the non-zero entries of $[N]$ appear in symmetric positions.



(a)

(b)

(c)

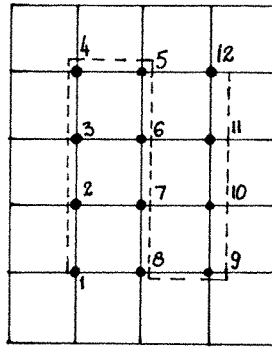Figure 2. (a), (b) examples for GIL; (c) generation of a GIL

Figure 3. Finite difference domain for example 1

*A possible iterative procedure*

The second part of the GI method is associated with defining an iterative technique for solving the difference equation (4) along a GIL. Since there are efficient block tri-diagonal systems solvers, one possibility is to use such a technique by writing

$$[M]\{\Phi\}^{n+1} = [N]\{\Phi\}^n + \{b\}^n \tag{6}$$

where the upper index $n$ is the iteration number counter. Here, the contributions of the side points are treated explicitly, whereas all the main variables (along the GIL) are treated implicitly (and are solved, for example, by the generalized Thomas algorithm). This possibility resembles in some sense the ADI and the LR techniques as will be mentioned later. Other stationary techniques, such as the SI or the MSI procedures for solving equation (4) will also be considered shortly. Let us demonstrate such a solution by a simple example.

*Example 1*

Figure 3 describes a field of $6 \times 5$ discrete points, governed by an elliptic discretized equation for $\Phi$ such as equation (2), with Dirichlet boundary conditions. The chosen GIL is shown also on this Figure. The source term vector $\{b\}$ and the matrix $[M]$ are

$$\{b\} = \begin{bmatrix} D_1 + S_1\Phi_{1_s} + W_1\Phi_{1_w} + E_1\Phi_8 \\ D_2 + W_2\Phi_{2_w} + E_2\Phi_7 \\ D_3 + W_3\Phi_{3_w} + E_3\Phi_6 \\ D_4 + N_4\Phi_4 + W_4\Phi_{4_w} \\ D_5 + N_5\Phi_{5_n} + E_5\Phi_{12} \\ D_6 + W_6\Phi_3 + E_6\Phi_{11} \\ D_7 + W_7\Phi_2 + E_7\Phi_{10} \\ D_8 + S_8 + W_8\Phi_1 \\ D_9 + E_9\Phi_{9_e} + S_9\Phi_{9_s} \\ D_{10} + W_{10}\Phi_7 + E_{10}\Phi_{10_e} \\ D_{11} + W_{11}\Phi_6 + E_{11}\Phi_{11_e} \\ D_{12} + W_{12}\Phi_5 + E_{12}\Phi_{12_e} + N_{12}\Phi_{12_n} \end{bmatrix} \tag{7a}$$
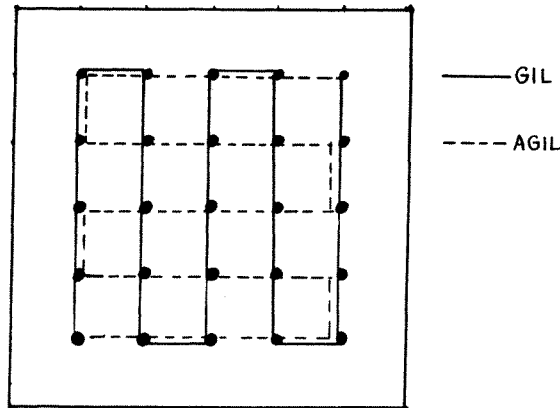
and

$$[M] = \begin{bmatrix}
P_1 & N_1 \\
S_2 & P_2 & N_2 \\
 & S_3 & P_3 & N_3 \\
 & & S_4 & P_4 & E_4 \\
 & & & W_5 & P_5 & S_5 \\
 & & & & N_6 & P_6 & S_6 \\
 & & & & & N_7 & P_7 & S_7 \\
 & & & & & & N_8 & P_8 & E_8 \\
 & & & & & & & W_9 & P_9 & N_9 \\
 & & & & & & & & S_{10} & P_{10} & N_{10} \\
 & & & & & & & & & S_{11} & P_{11} & N_{11} \\
 & & & & & & & & & & S_{12} & P_{12}
\end{bmatrix} \quad (7b)$$

The underlined terms in $\{b\}$ are the 'side' points' contributions that are treated as known (explicitly) from iteration to iteration. It is clear that the technique is more like the Jacobi block iterative method than the GS block iterative method, since the old variables are not updated during the course of the iteration over all the field points. It is apparent that many other GILs can be constructed but the procedure is the same.

*Generalization of the method*

The GI method has two degrees of freedom:
(1) the choice of the GIL according to *Definition* 1
(2) the choice of the means of splitting of the coefficients matrix $[A]$ according to the chosen iterative procedure to solve equation (4).
It can easily be verified that with a choice of certain GILs and certain splittings of $[A]$, one can obtain all the implicit techniques. For example, the ADI, the LR and the block Jacobi procedures may be obtained by letting the GIL (such as in Figure 2, say) include the boundary points, as can be seen in Figure 4(a). Here the algebraic system of $K$ equations may be separated, for example, into $m - 2$ subsystems of $n - 2$ equations. Every subsystem is solved implicitly for $\Phi$, along a column (or a row) of the field. The above procedures will be different only in the choice of the splitting procedure [as in equation (5)]. A similar technique to the MSI method can be obtained by choosing the GIL to have the 'stairs' mode along diagonals as shown in Figure 4(b). When the SI procedure[6] is compared qualitatively to the GI procedure, one can see that the special way of treating the two dimensional nature of the problem in the SI procedure is inherent in the GI procedure. Moreover, with the SI procedure a group of half of the points near the boundaries (which, in some sense, carry the short wave disturbances) do not satisfy the difference equation, whereas in the GI procedure all of the points in this group are very close to satisfying or exactly satisfy the governing difference equation.

The way of treating the boundary conditions in the GI method is not the main reason for the success of this method (which is much faster in convergence than other implicit methods), but as will be discussed later, it is one of the reasons to the smooth rate of convergence as well.

Figure 4. Reduction of the GI method to the (a) LR and (b) MSI procedures

## An alternating direction general implicit (AGI) method

For every given GIL it is possible to find at least one other GIL that fulfils the following conditions:

(1) It will be normal to the original line at all points except at the boundary points or at the near boundary points.

(2) On the boundary points or at the near boundary points this line may be parallel or normal to the original GIL as the continuity of this line requires.

In other words: the 'side' points of the GIL will be converted into 'main' points of the AGIL, and some of the 'main' points of the original GIL will be converted into 'side' points on the AGIL. An example of an AGIL is given in Figure 5. One iteration of this procedure is obtained by two sweeps; the first along GIL and the second along the AGIL. This technique will be called the alternating general implicit (AGI) procedure. It is reasonable to assume that the AGI method should converge faster then the GI method, just as the alternating line relaxation (ALR) techniques converge faster than the regular one-direction implicit procedures. In the present paper the AGI method, which is a topic by itself, will not be discussed in any depth. Next, we shall derive the basic features of the GI method for any elliptic operator $L$ of equation (1) in conjunction with a convection–diffusion-like problem.

Figure 5. Example of an AGIL

## The difference scheme of the convection–diffusion equation

The present study will be concentrated on the linear convection–diffusion equation, which is a reasonable model for many of the physical transport phenomena, such as the Navier–Stokes equations and the energy equation, among others. We shall deal with a time dependent self adjoint operator $L$ which represents the following convection–diffusion problem for $\Phi$:

$$\Phi_t + u(x, y)\Phi_x + v(x, y)\Phi_y = \frac{1}{R}(\Phi_{xx} + \Phi_{yy}) + S(x, y) \tag{8}$$

where $t$ is the time or the time-like (iteration) co-ordinate, $x$ and $y$ are the physical co-ordinates of the two dimensional domain, $u$ and $v$ are prescribed velocity-like convection coefficient functions which are partially continued over $\Omega$, and $S(x, y)$ is the source term. We shall assume also that $\max(|u|, |v|)$ is of the order of 1. For the Navier–Stokes equations, $R$ is the Reynolds number of the flow. For simplicity without losing generality, let us spread over $\Omega$ an equally spaced grid in both directions $x$ and $y$ (see Figure 1). The diffusion term in the $x$ direction, $\Phi_{xx}$, is discretized by a second order central difference at the point p as follows:

$$(\Phi_{xx})_p = \frac{\Phi_e - 2\Phi_p + \Phi_w}{\Delta x^2} \tag{9}$$

A similar expression for the $y$ diffusion term, $\Phi_{yy}$, may also be written. The convection terms are modelled in general by including the second-order correction.[10] The convection in the $x$ direction, $\Phi_x$, for example, is differenced as follows:

$$\Phi_x^{n+1} = \begin{cases} \dfrac{\Phi_p^{n+1} - \Phi^{n+1}}{\Delta x} + \dfrac{\Phi_e^l - 2\Phi_p^l + \Phi_w^l}{2\Delta x}, & \text{for } u_p \leqslant 0 \\[4mm] \dfrac{\Phi_e^{n+1} - \Phi_p^{n+1}}{\Delta x} - \dfrac{\Phi_e^l - 2\Phi_p^l + \Phi_w^l}{2\Delta x}, & \text{for } u_p > 0 \end{cases} \tag{10}$$

where the upper index $n$ indicates the time step or the iteration number. Generally, for $l = n + 1$, the two equations (10) are the same, being the second-order accurate finite difference model for the first derivative. For $l = n$ equation (10) results in the upwind differencing scheme which recovers its second order accuracy only in the steady state. For one dimensional cases it has been shown[10] that the $l = n$ case is unconditionally stable and consistent. The $l = 0$ case is the classical upwind

difference for $\Phi_x$. Here both $l = n + 1$ and $l = n$ will be considered. The time derivative, $\Phi_t$, is represented by the one sided differencing:

$$(\Phi_t)_p = \frac{\Phi_p^{n+1} - \Phi_p^n}{\Delta t} \tag{11}$$

Substituting equations (9), (10) and (11) into equation (8), an equation of the form of equation (2) is obtained for every point p in the field. It can be shown, by incorporating similar finite difference approximation to the various derivatives of equation (8), that equation (2) can always be derived with the following features:

(1) $P > 0$                                $P > 0$                      (12a)

and only for $l = n$:

(2)                              $E, W, N, S \leqslant 0$                    (12b)

(3)                        $|P| > \begin{cases} |E| + |W| \\ \text{and} \\ |N| + |S| \end{cases}$            (12c)

It should be noted that constructing the coefficient matrix $[A]$ with any GIL, results in the $P$ coefficients being on the main diagonal of $[A]$. Equations (12c) means that for $l = n$ in equation (10) the matrix $[A]$ is strictly diagonal dominant. Next we shall examine the features of the matrix coefficients of the GI algorithm.

*Some features of the GI algebraic system*

In the general case of a convection–diffusion equation, $[A]$ is not necessarily symmetric and therefore it is not necessarily positive definite unless there is no convection in the field, in which case a Poisson-like equation governs the transport phenomenon.

All the iteration methods seek to solve the non-singular $K \times K$ system of equations (2), (4) by splitting $[A]$ into the $[M]$ and $[N]$ matrices as was given by equation (5). The iterative method then evaluates $\Phi$ as follows:

$$[M]\{\Phi\}^{(n+1)} = [N]\{\Phi\}^{(n)} + \{D\} \tag{13}$$

or

$$\{\Phi\}^{(n+1)} = [T]\{\Phi\}^{(n)} + \{F\} \tag{14}$$

where

$$[T] = [M]^{-1}[N] \tag{15}$$

is the iteration matrix, and

$$\{F\} = [M]^{-1}\{D\} \tag{16}$$

The first theorem can be proved directly from the definition of $[A]$ of the GI method:

*Theorem 1.* The coefficient matrix is not singular and $[A]^{-1} > 0$.

*Proof.* Since there are non-positive entries in $[A]$, except those on the main diagonal which are always positive, this theorem is proved by Theorem (3.10) of Reference 8.

Now consider iterative solutions that resemble the LR or the ADI methods in some sense; in this case the matrices $[M]$ and $\{D\} + [N]\{\Phi\}$ are very similar to those in equations (7). This means that $[M]$ consists of the three main diagonals of $[A]$, and $[N]$ consists of the rest of the entries in $[A]$. Now we can define the features of $[M]$ and $[N]$.

*Theorem* 2. For every GIL, $[M]^{-1} > 0$.

*Proof.* All the entries in $[M]$ are zero except the main three diagonals. Let us denote the main diagonal by $\{b_j\}$, the upper diagonal by $\{c_j\}$ and the lower diagonal by $\{a_j\}$. According to the GI method's features [equation (12), Theorem 1], $b_j \geqslant 0$, $c_j < 0$, $a_j < 0$ and $a_j + b_j + c_j > 0$. Therefore $[M]^{-1} \geqslant 0$ according to Theorem (3.10) of Reference 8.

Also we have:

$$\text{For every GIL } [N] \geqslant 0$$

*Proof.* Since it is apparent from equation (12c) that the entries of $-[N]$ are the same as those of the off-diagonal entries of $[A]$, which are not positive, therefore those of $[N]$ are not negative.

From these Theorems it can be concluded that the splitting defined by equation (13) is a regular splitting, which might be defined as follows:

*Definition* 2. For $K \times K$ real matrices, $[A]$, $[M]$ and $[N]$, equation (13) is a regular splitting if $[A]$ and $[M]$ are non-singular with $[M]^{-1} \geqslant 0$ and $[N] \geqslant 0$. Thus, equations (13) and (14) present a regular splitting.

By defining the matrix $[G]$ as

$$[G] = [A]^{-1}[N] \tag{17}$$

it can be shown that

$$[T] - [M]^{-1}[N] = ([I] + [G]^{-1})[G] \tag{18}$$

and

$$\mu = \frac{\lambda}{1 + \lambda} \tag{19}$$

where $[I]$ is the unitary matrix, $\lambda$ is an eigenvalue of $[G]$ and $\mu$ is the respective eigenvalue of $[T]$. This leads to the following spectral radius relation:

*Theorem* 3. if $[A] = [M] - [N]$ is a regular splitting with $[A]^{-1} > 0$, then

$$\rho(T) = \frac{\rho(G)}{1 + \rho(G)} < 1 \tag{20}$$

where $\rho(H)$ is the spectral radius of $[H]$. Thus, $[T]$ is convergent and the iterative procedure, equation (14), converges for any initial vector $\{\Phi\}$. Usually, this convergence theorem gives little, if any, information regarding the rate of convergence. Even if a method converges, it may converge more slowly than other existing methods, and therefore such a method has no practical use. On the other hand this theorem is useful to find an appropriate matrix $[N]$ to optimize the rate of convergence.

It can be seen from equation (20) that as $\rho(G)$ becomes smaller than 1, the method will converge faster. This means, that as $[M]$ resembles $[A]$, and $[N] \to 0$ the method is improved. Thus the superiority of the method can be examined by its choice for $[N]$. In the GI method we have some

control on the structure of $[N]$ by the right definition of the GIL, as is shown in the following section.

*The optimal GI method*

In the above method the matrix $[A]$ contains three diagonals with non-zero entries, and another $2N$ non-zero elements spreading in a symmetric fashion. All the rest of the entries are zero. An example is given later in equation (29). Here, the splitting of the matrix $[A]$ is done in a similar way to equation (15) by letting its three main diagonals form the matrix $[M]$ and all the other $2N$ non-zero elements are used to construct the $[N]$ matrix. For the $l = n$ case [see equation (10)], all the entries in $[N]$ are non-negative. The optimal GI method will be obtained from solving the following problem:

Find such a continuous ordering (and such a splitting of $[A]$) that $\rho(G) \rightarrow \min$.

It is obvious that the GIL connects all the field's grid points except those that have a Dirichlet boundary condition. The following theorem will lead us to a possible optimal GIL procedure.

*Theorem 4.* The matrix $[N]$ is singular at least of order 4.

*Proof.* For a GIL that begins in one field's corner grid point ($d_1$ say, see Figure 6), and ends in the opposite corner point ($d_2$, say), the equations for $\Phi$ in the two corner points ($d_3$ and $d_4$, say) do not include any unknowns, except those along the GIL as is seen in Figure 6. Therefore the entries of $d_3$ and $d_4$ are zero, as well as in $(d_3)_{+1}$ and $(d_4)_{+1}$, and the theorem is proved. Moreover, it can be seen that the number of non-zero eigenvalues of $[G]$ is reduced by the number of grid points on two opposite boundaries.

Let us consider the effect of the $[N] = (n_{ij})$ values on the spectral radius. Using theorem (1.5) of Reference 8 one may conclude that for any $n \times n$ matrix $[A] = (\alpha_{ij})$:

$$\rho(A) \leqslant \max_{i \leqslant j \leqslant n} \sum_{i=1}^{i=N} |a_{ij}| = \beta \tag{21}$$

Thus the maximum of the row or column sums of the absolute values of the elements of $[A]$ gives an upper bound to $\rho(A)$. This may imply that in order to reduce the spectral radius, the value of $\beta$ should be reduced. In view of this conclusion, $\rho$ will be as small as the entries of $[N]$ are small compared to the entries of $[M]$. Thus the optimal GIL can be defined as follows:



Figure 6. The effects of the boundaries on $[N]$

Figure 7. Building a GIL

*Definition 3*. An optimal GIL is the line on which any three successive main points have larger absolute values of the difference between the absolute values of the coefficients than those of the side points' coefficients.

In other words: if the optimal GIL moves to the point p from the point s (see Figure 7) and if none of the points e, w and n are already on the line, then the three of them are legitimate candidates to continue the current GIL. According to this definition, the line will be continued to the point with the smallest absolute value coefficient. The following example will demonstrate an application of the GI procedure as formulated above and the main difficulties associated with it.

*Example 2*

In the second example let us consider the following Dirichlet problem

$$\Delta\Phi = \frac{\partial^2\Phi}{\partial x^2} + \frac{\partial^2\Phi}{\partial y^2} = -Q \tag{22}$$

with the boundary condition

$$\Phi(x, y) = \Phi_\omega(x, y), \quad (x, y)\in\partial\Omega$$

where $Q = Q(x, y)$ is the source term.

Let us assume for simplicity that $\Omega = (x, y):0 \leqslant x, y \leqslant 1$. The above equation is discretized in the



Figure 8. Finite difference notation for the Poisson equation solved over a unit square

usual way; a grid is spread over $\Omega$ with lines parallel to the $x$ and $y$ co-ordinates. It is also assumed that the $m$ discrete points along the $x$ axis and the $n$ points along the $y$ axis are spaced evenly, so that the spacings (see Figure 8)

$$\Delta x = \frac{1}{m-1}; \quad \Delta y = \frac{1}{n-1} \tag{23}$$

are constant. Also, without losing generality, we will assume that

$$\alpha = \frac{\Delta x}{\Delta y} \leqslant 1, \quad \text{i.e. } m \geqslant n \tag{24}$$

The finite difference approximation to equation (22) is

$$E = W = 1 \tag{25a}$$

$$N = S = \alpha^2 \tag{25b}$$

$$P = -2(1 + \alpha^2) \tag{25c}$$

where $E$, $W$, $N$, $S$ and $P$ are defined in equation (2). Say that in this case we begin to draw the GIL from the lower left point of $\Omega$, as in Figure 8.

Since the $N$ coefficient is smaller than the $E$ coefficient, the GIL will move up to the point 2. The same argument holds along all the way up to the point $n - 2$. Here, since the $N$ coefficient is zero, the only possibility to continue the GIL is by turning it to the right towards the point $n - 1$. From this point we have again two possibilities: to continue the line to the right [to the point $(i, j) = (4, n - 1)$] or downwards [to the point $(i, j) = (3, N - 2)$]. Since the $S$ coefficient is smaller than the $E$ coefficient the GIL will move down to the point $n$, etc.

Results for this example were obtained for $(N - 1) \times (M - 1) = 900$ with $Q = 0$ and $\Phi_\omega(x, y) = 0$. As initial conditions we have chosen

$$\Phi(x, y)|_{x, y \in \Omega} = 1$$

Figure 9 describes the variation of the rate of convergence for the GI method as a function of $\alpha$ as well as the variation of some other similar methods such as the JLR method[11] and the ADI



Figure 9. Rate of convergence of the GI and other similar methods

method.[5] The rate of convergence presented here is normalized with respect to the work units used by these methods, where the GS method is defined to use one working unit.

It can be seen that the Jacobi implicit line relaxation (JLR) is slower than the GI method, but both of them are close. The Gauss–Siedel implicit line relaxation (GSLR) is better than the GI method. The reason is that in the GI method the variables are updated only after the whole field is calculated, whereas in the regular (successive) line relaxation technique the variables are updated on every line once they have been calculated. The main disadvantage of the GI procedure is that the variables cannot be updated during the course of one sweep over the field. For the Laplace equation it can be shown that the GSLR converges twice as fast as the JLR technique for $\alpha = 1$; the rate of convergence of the JLR is similar to that of the GI procedure, and, as $\alpha$ goes up, the rate of convergence of the GI procedure goes up by about 20 per cent relative to that of the JLR technique. The main conclusion that might be derived is that the power of the spatial ordering of the grid points as is expressed in the creation of the GIL should be considered in another frame of iterative procedures. In this frame we have to emphasize the possibility of using an optimal GIL more carefully, as is suggested in *Definition* 3. A possible way to accomplish this is by examining the different factorizations of $[A]$.

*The two-dimensional GI method*

Let $t$ be the index marking the elements along the GIL, then a possible factorization of the GI method can be expressed as follows:

$$\Phi_t = E_t \Phi_{t+1} + F_t \tag{26}$$

Since the GI method with the factorization (26) [which is the same as in Theorem 3], was found not to be beneficial, a fully multi-dimensional iteration procedure was considered. Here, just as in the SI technique, the solution depends on the variables located along all the directions of the domain co-ordinates. The procedure may be written for one variable $\Phi$ in the two-dimensional case as

$$\Phi_{ij} = A_{ij} \Phi_{i+1\,j} + B_{ij} \Phi_{ij+1} + C_{ij} \tag{27}$$

where $i$ and $j$ are the two direction's indices of the domain.

With not much additional work it can be shown that the above factorization for the GI method, which is like the LR method is, in some sense, a special case of the SI technique. In the spirit of the SI method, we can define the direction of the GIL as $t$ and that of the AGIL as $n$, and formulate an SI-like procedure for $\Phi_t$:

$$\Phi_t = A_t \Phi_{t+1} + B_t \Phi_{np1} + C_t \tag{28}$$

where the point np1 is a side point which is located on the GIL, and should be one of the four neighbour points to the point $t$ that appears on GIL adjacent to the point $t$. Figure 10 summarizes several possible cases, and shows how the point np1 should be chosen. In cases where two possibilities for choosing the point np1 are available, numerical experiments have shown that as $|np1 - t|$ is reduced, the rate of convergence is increased. The two dimensional GI (TDGI) technique is illustrated in the following example which is very similar to Example 2.

*Example 3*

Let us consider the domain $\Omega$ and the GIL as are defined in Figure 11, and let us not be specific about the elliptic operator that has to be resolved here.
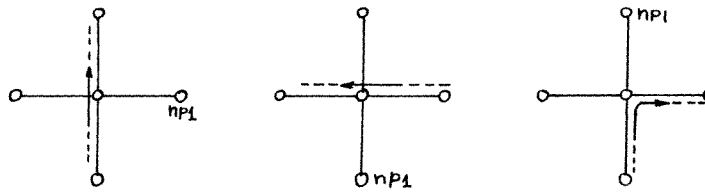
Figure 10. Possible positions for the point np1 on the OGIL

The system matrix $[A]$ is

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|
| 1 | $P_1$ | $N_1$ | | $E_1$ | | | | | | | | | | | | |
| 2 | $S_2$ | $P_2$ | $E_2$ | | | | | | | | | | | $N_2$ | | |
| 3 | | $W_3$ | $P_3$ | $S_3$ | | | | $E_3$ | | | | | | $N_3$ | | |
| 4 | $W_4$ | | $N_4$ | $P_4$ | $E_4$ | | | | | | | | | | | |
| 5 | | | | $W_5$ | $P_5$ | $E_5$ | | $N_5$ | | | | | | | | |
| 6 | | | | | $W_6$ | $P_6$ | $N_6$ | | | | | | | | | |
| 7 | | | | | | $S_7$ | $P_7$ | $W_7$ | | $N_7$ | | | | | | |
| 8 | | | $W_8$ | | $S_8$ | | $E_8$ | $P_8$ | $N_8$ | | | | | | | |
| 9 | | | | | | | | $S_9$ | $P_9$ | $E_9$ | | $N_9$ | | $W_9$ | | |
| 10 | | | | | | | $S_{10}$ | | $W_{10}$ | $P_{10}$ | $N_{10}$ | | | | | |
| 11 | | | | | | | | | | $S_{11}$ | $P_{11}$ | $W_{11}$ | | | | |
| 12 | | | | | | | | | $S_{12}$ | | $W_{12}$ | $P_{12}$ | $E_{12}$ | | | |
| 13 | | | | | | | | | | | | $E_{13}$ | $P_{13}$ | $S_{13}$ | | $W_{13}$ |
| 14 | | $S_{14}$ | | | | | | | $E_{14}$ | | | $N_{14}$ | $P_{14}$ | $W_{14}$ | | |
| 15 | $S_{15}$ | | | | | | | | | | | | $E_{15}$ | $P_{15}$ | $N_{15}$ | |
| 16 | | | | | | | | | | | | $E_{16}$ | | $S_{16}$ | $P_{16}$ | |

A possible algorithm point by point is

$$\Phi_1 = A_1\Phi_2 + B_1\Phi_4 + C_1 \tag{30}$$

$$\Phi_2 = A_2\Phi_3 + B_2\Phi_{15} + C_2$$

$$\Phi_3 = A_3\Phi_4 + B_3\Phi_8 + C_3$$

$$\Phi_4 = A_4\Phi_5 + C_4$$

$$\Phi_5 = A_5\Phi_6 + B_5\Phi_8 + C_5$$

$$\Phi_6 = A_6\Phi_7 + C_6$$

$$\Phi_7 = A_7\Phi_8 + B_7\Phi_{10} + C_7$$

$$\Phi_8 = A_8\Phi_9 + C_8$$



Figure 11. GIL for example 3

$$\Phi_9 \ = A_9\Phi_{10} + B_9\Phi_{12} + C_9$$

$$\Phi_{10} = A_{10}\Phi_{11} + C_{10}$$

$$\Phi_{11} = A_{11}\Phi_{12} + C_{11}$$

$$\Phi_{12} = A_{12}\Phi_{13} + C_{12}$$

$$\Phi_{13} = A_{13}\Phi_{14} + B_{13}\Phi_{16} + C_{13}$$

$$\Phi_{14} = A_{14}\Phi_{15} + C_{14}$$

$$\Phi_{15} = A_{15}\Phi_{16} + C_{15}$$

$$\Phi_{16} = A_{16}\Phi_{16_b} + C_{16}$$

where $A_i$, $B_i$ and $C_i$ are the algorithm coefficients.

The algorithm coefficients can be found by comparing the governing differenced equation at the point $t$ along the GIL from $t = 1$ to $t = 16$:

$$P_t\Phi_t + N_t\Phi_{t+1} + S_t\Phi_{t-1} \rightarrow main\ points$$
$$+ E_t\Phi_u + W_t\Phi_d \quad \rightarrow side\ points$$
$$+ D_t = 0 \tag{31}$$

For $t = 1$ the only side point is No. 4 so that $A_1$, $B_1$ and $C_1$ in equation (30) are known (in terms of $P_1$, $N_1$, $S_1$, $E_1$, $W_1$ and $D_1$). For $t = 2$ the only side point is No. 15, so that the governing equation (31) can be written as:

$$P_2\Phi_2 + N_2\Phi_3 + S_2\Phi_1 + E_2\Phi_{15} + D_2 = 0 \tag{32}$$

substitution of $\Phi_1$ from equation (30) gives the following relations:

$$A_2 = \frac{N_2}{K}$$

$$B_2 = \frac{E_2}{K}$$

$$C_2 = \frac{D_2 + S_2(B_1\Phi_4 + C_1)}{K}$$

where

$$K = -\frac{1}{P_2 + A_1 S_2}$$

The same procedure is repeated for $t = 3, 4, \ldots, 16$. During the evaluation of the coefficients, just as in the standard SI procedure, some off-diagonal values such as $\Phi_4$ in the equation at $t = 2$, are assumed to be known from the previous iteration. After evaluating the coefficients, new $\Phi$ values are calculated using equation (28).

## OPTIMIZATION OF THE GI METHOD

Although using one GIL results in a low computer time, for long GILs or for complex geometries it may be helpful to reduce the programing effort to split the GIL into several parts. It is reasonable to choose the cut-off points as the grid points where the convection is much smaller than the diffusion. That is since, as was discussed previously for the pure diffusion equation, the rates of convergence of the GI method and the GSLR method are comparable; what is gained with the GI method by having fewer and smaller eigenvalues, is lost because the iteration technique using the

Figure 12. A boundary-layer-like flow

new value of the variable is not immediately used, as is the case with the GSLR method. Naturally, in most of the cases the grid points near a solid boundary are good candidates for marking the devisions of the GIL. Inner points can also serve as the dividing points of a GIL. For example, in solving the fluid flow equations, the points at which the flow velocities are very low, such as at the centre of a vortex (or where the stream function reaches its extremum values) are a reasonable choice for dividing points. We will define an optimal GIL as a GIL which is divided in an optimal way and the method as the optimal GI (OGI) method. The optimality will be defined after some special cases have been considered

*Example 4. The boundary layer like flow*[11]

A boundary layer like flow is a unidirectional flow field in which the diffusion effects are important only in direction normal to the outer flow direction, say $x$ in Figure 12. We will consider here the following reduced Navier–Stokes equations for such a flow:[11]

$$uu_x + vv_y = \varepsilon(u_{xx} + u_{yy}) \tag{33a}$$

$$u_x + v_y = 0 \tag{33b}$$

where $\varepsilon$ is the inverse of the flow Reynolds number. Usually equations (33) are solved by marching techniques along $x$ in which the $u$ and $v$ velocities are solved implicitly along the normal co-ordinate $y$. Let us choose the computational domain as in Figure 13, with the boundary conditions as depicted on the Figure.

For a square grid $\Delta x = \Delta y$ for standard central differences one will find that the $|E|$ and $|W|$ coefficients in equation (33a) are bigger than the $|N|$ and $|S|$ coefficients; this is due to the nature of the boundary layer in which $|u| > |v|$. It follows that the GIL will be mostly normal to $x$. This means that the GI procedure resembles the marching technique for a boundary layer. Thus we may conclude that the break points along the GIL will be the upper $(y = y_n)$ and the lower $(y = 0)$



Figure 13. Computational domain for the boundary layer fow

Figure 14. A boundary layer with injection

boundaries. Thus the OGI technique is reduced to the standard marching procedure. This conclusion depends very much on the way the convection terms are differenced. However, if the convection terms in the $y$ direction were differenced in the upwind fashion then the GIL would not have to trace the boundary-layer-like marching steps. Further, if $\Delta x \neq \Delta y$, there will exist some points where the GIL will be parallel to the plate's $x$ direction rather than normal to it. Since the GIL direction depends on the order of the finite difference scheme, we can extend the definition of a GIL to include terms like 'the order of the GIL'. Thus it can be said that to the second order of accuracy the GIL is similar to the boundary layer natural co-ordinates.

*Example 5. flow with small recirculation*[11]

Let us consider the previous example of a uniform flow over a flat plate where fluid is also injected into the flow at the wall and normal to it, and the same amount of flow is sucked back just downstream of the injection region as in Figure 14.

The steady-state incompressible flow governing equations (Navier–Stokes) are

$$uu_x + vu_y = -p_x + \varepsilon(u_{xx} + u_{yy}) \tag{34a}$$

$$uv_x + vv_y = -p_y + \varepsilon(v_{xx} + v_{yy}) \tag{34b}$$

$$u_x + v_y = 0 \tag{34c}$$

An equal spaced staggered mesh[12] $\Delta x = \Delta y$ is spread over the computational domain for solving this system numerically. All the derivatives are modelled by the standard central differences (except



Figure 15. GILs for boundary layer with injection

near the right and upper boundaries and for the continuity equation (34c) at one inner point).

If the injection velocity $\bar{V}_s = u_s/u_\infty$ is small enough (say $\bar{v}_s < 0\cdot1$) so that $u \geqslant 0$ in all the flow field, the construction of the GIL near the recirculation zone for $\Delta x = \Delta y$ can be constructed as follows:

Beginning from the point 'a' on the wall close to the injection region, as in Figure 15, the GIL will move up and to the left since for high Reynolds numbers $|v| > u$ in this zone whose dimensions are dependent on the Reynolds number. That is, because of the continuity equation, $v$ is more or less the injection velocity and $u$ varies linearly with $y$. The line will remain parallel to the $x$ direction until a point at which $|v| < u$ is reached. Here the GIL will move in the direction normal to the wall. This direction will be kept until a point is reached at which $|v| > u$ again and the GIL will change its direction again to the left, etc. This pass will eventually lead the GIL out from the recirculation zone, in a direction parallel to $y$. Thus the velocity along this pass is $O(\sqrt{\varepsilon})$ and will stay more or less with the same value (which may have small changes because of the changes in the thickness of the boundary layer), and the boundary layer behaviour is recovered again. The GIL which begins at the point 'b' near the end of the sucking zone will behave similarly to that which begin at the point 'a', except that it will move to the right instead of moving left. Thus, the inner recirculation centre 'c' acts like a point of symmetry; here the GIL which passes through 'c' is more or less normal to the wall. Examining this line, it can be seen that it roughly follows the normals to the streamlines $\psi$. Such lines resemble the potential-lines $\varphi$. This observation is not a surprise since the normal velocity to $\psi$ along the streamline is zero (by definition), and all the convection is parallel to it. Thus just as the potential lines carry the maximum change of the (streamwise) velocities in the physical plane, the GIL points towards the maximum change of the variables on the finite difference grid. This conclusion results also from the fact that theoretically the GIL follows the normal direction to the velocity vector and the GIL pass will follow as close as possible to the normals to $\psi$. All these conclusions may result from the following proposition:

### Proposition

If the derivatives of the quantities appearing in the elliptic PDE describe the motion of a fluid flow and these equations are solved numerically by standard central differences, then the GILs are the potential-like lines.

We prefer to put this proposition in this way although the potential lines are not defined for viscous flow. However, if we define $\varphi$ to be normal to $\psi$ (and not by the continuity equation) then it is well posed. The proof for this proposition is rather simple in the light of the above explanation. However, it is not so simple to generalize this proposition to cases where the convection terms are treated by upwind difference methods[13] $[l = 0$ in equation (10)]. Sometimes it may be found that in these cases the best GIL is parallel to $\psi$ lines, but now there is no *ad hoc* proof that this is the case in general. In addition to these conclusions, the last proposition has much more basic applications than the previous GI techniques.

In recent years many studies in the various areas in science in which elliptic equations govern phenomena emphasize the mapping of the domain over which the problem has to be solved as a basic tool of the solution.[14] The philosophy is that in many problems it is preferable to map a complicated region and to solve a more complicated equation, rather than to solve the original equations on a complicated region. Many mapping techniques have been suggested.[15] For problems in fluid mechanics, the field boundaries are usually simple functions of the stream function (such as a constant stream function) the OGI method suggests that the potential-like lives be used as one of the independent field co-ordinates not only because the new region becomes much more simple, but because certain numerical schemes (such as the one presented here) become

optimal along such a co-ordinate. These methods will be discussed elsewhere.[16] In the next section we will demonstrate how to apply the schemes discussed above (mainly the OGI method) to some simple fluid dynamic problems.

## APPLICATIONS

This section presents some applications of the OGI method to some fluid dynamics problems. The purpose here is not to solve new problems, nor even to study extensively the problem's results, but the main point is to compare the convergence properties of the OGI method with similar iterative methods (comparisons with methods which involve some special acceleration techniques such as the CG procedures to solve equation (1) will not be considered). It should be emphasized that the present method is only an iterative technique and does not overcome any difficulties which arise from the finite difference scheme, for example the oscillatory behaviour of the steady-state solution for high Reynolds number flows when central differences are used everywhere is not prevented. With the Proposition in the last section it is possible to make sure that the iterative procedure will converge the difference equations faster than the same procedure when it is used on a regular (horizontal–vertical) mesh, using the same iterative update technique. According to the OGI iterative procedure before solving a flow field, the patterns of the various parts of the GIL have to be guessed, more or less in the expected potential-like line directions. These patterns are corrected after several iterations, especially when a significant change in the convection coefficients (flow velocities) is encountered. If the iterative procedure is convergent then the GIL converges to its final OGIL pass. The procedure which generates the GIL pattern automatically is quite complicated. In addition to the request that this procedure will construct enough parts of the GIL so that all the grid points will share this line, it should be also able to deal with every part separately. In the present paper we are only concerned with the principles of the OGI method and not with the automatic generation of the GIL. Since for many problems it is not easy to formulate the right GIL, we have chosen in this study to present only fields with simple geometries. As a result we have chosen to present the OGI method for fluid dynamic problems for which the solutions are very well established.

### Flow over a step

The steady-state incompressible flow over a step as it is presented in Figure 16 is popular for examining numerical schemes.[17] This flow has one main recirculating eddy, and, as the Reynolds number increases, smaller eddies may appear close to the corner. The results are very sensitive to the finite difference technique. If the convection terms were modelled by the standard central difference scheme then the high Reynolds number iterative solutions might diverge. If the upwind scheme is adopted then false diffusion effects will change the conditions of the problem.
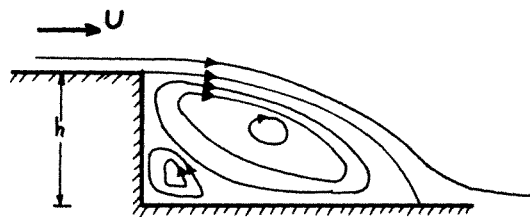


Figure 16. Flow over a step

We consider here two systems of equations. The first is the primitive variable system, equations (34). The second is the stream function $\psi$–vorticity $\omega$ formulation:

$$\Delta^2\psi + \omega = 0 \qquad (35a)$$

$$u\omega_x + v\omega_y = v\Delta^2\omega \qquad (35b)$$

where

$$\omega = v_x - u_y \qquad (36a)$$

$$\psi_x = -v \qquad (36b)$$

$$\psi_y = u \qquad (36c)$$

For all the derivatives central differences are employed. System (34) is solved on a staggered mesh as shown in Figure 17. Since the coefficients of the convection terms in this system are the same, both of these equations can use the same OGI method with the same GIL. Since the continuity equation is of the first order, it does not have a specific direction, and the diagonal dominance is always kept. Figure 18 presents some optimal GILs for Reynolds numbers $R_h = Uh/v = 1500$. The



Figure 17. The staggered mesh



Figure 18. GILs for the flow over a step

Figure 19. Comparisons of the OGI method with similar iterative methods for the flow over a step

GILs follow more or less the potential-like lines pattern. For the equal spaced mesh spreadings, Figure 19 presents a comparison between the OGI method and other iterative methods. The iteration co-ordinate in this Figure is normalized with respect to the number of work units expended in doing one iteration. One normalized iteration is defined as that of the LR method. The rate of convergence of the OGI method is about 3·5 times better than that of the MSI2 method[7] and about 10 times better than the ADI method. Figure 21 depicts the rate of convergence in the linear zone of $\rho$ as a function of $R_h$. Again it can be seen that the OGI method is superior to the other methods.



Figure 20. The driven cavity geometry and OGILs for Reynolds number = 100
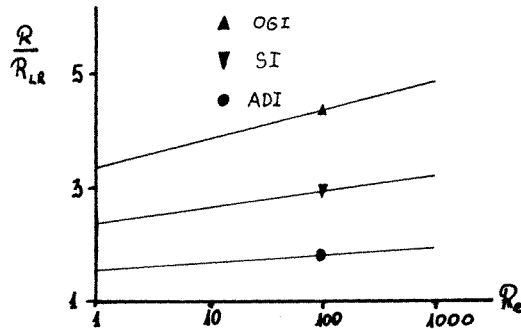
Figure 21. Rates of convergence for the driven cavity problem

Now consider the $\psi$–$\omega$ formulation. It is well known that the two equations should be solved in a coupled manner in order to avoid divergence even for relatively low Reynolds numbers. This is mainly due to the linear dependence of the vorticity on the stream-function gradients at the solid boundaries. Since $\psi$ changes mostly along the potential-like lines, it is easy to see that equation (35a) will have small explicit perturbations along a near optimal GIL. Thus, theoretically, if there were an infinite number of grid points, it would be sufficient to treat equation (35a) totally explicit except near solid boundaries. It was found that for the $\psi$–$\omega$ system the OGI method is a slightly better than the LR method, which is slightly worse than the ADI method. However, on checking the various procedures to generate the OGI it appears that there exists an OGIL along which the $\psi$–$\omega$ system will converge much faster than other implicit methods as will be shown below.

*The driven cavity*

The field geometry of this well known problem is shown in Figure 20.

The driven cavity case is a simple problem that has often been used to test and to compare numerical methods.[18] The governing equations for the incompressible steady state are once again equations (35) and (36). Usually this problem has been solved by the $\psi$–$\omega$ formation [equations (35)], but as has been discussed before, the optimal GIL which is very suitable for solving the $\omega$ equation, is a bad choice for solving the first equation [equation (35)] for $\psi$] implicitly. Figure 20 presents a typical OGIL for half of the cavity, for $Re = 100$, when the primitive variable system, equations (34). is used. In the region close to the main vortex, the implicit lines are more or less normal to the streamlines, whereas the implicit lines near the boundaries give only some indication[7] of the existence of small secondary recirculation vortices near the corners. Figure 21 presents a comparison of the convergence rate between the OGI methods and some other similar iterative methods, as function of the Reynolds number. It can be seen that the OGI method is faster than other implicit techniques and as the Reynolds number increases the rate of convergence of all the methods decreases.

## CONCLUSIONS

In this paper, an iterative strategy to solve the the algebraic system of equations arising from discretizing an elliptic equation by some finite difference technique has been presented. The GI and the OGI methods which result from this strategy reflect the effect of the current solution on the way the iterative process is executed. The lines along which the algebraic system is solved implicitly are

no longer columns or rows of the grid, but a general pass through the grid points in the field. It has also been shown that there is a so-called optimal line with which the maximum rate of convergence is achieved.

This family of methods is new and has therefore not been exhaustively discussed. Although it is difficult to apply this procedure, it has been shown that it has a pronounced superiority to other implicit methods (which are particular cases of the GI method).

The main application of this method is in numerical fluid dynamics problems, in which the rate of convergence is reduced dramatically as the Reynolds number increases. The OGI method depends very weakly on the Reynolds number but is very sensitive to the OGIL pass. Generally the OGIL follow the normals to the streamlines which in inviscid flow are the potential lines and which can be defined in this way also for viscous flow.

This work represents a preleminary study of this new method in which the feasibility of this strategy is shown. More study is required in order to evaluate this method more generally.

## REFERENCES

1. A. Jennings, *Matrix Computation for Engineers and Scientists*, Wiley-Interscience, 1977.
2. P. N. Schwatztrauber and R. A. Sweet, "The direct solution of the discrete Poisson equation on a disk', *SIAM J. Number. Anal.*, **5**, 950–970 (1977).
3. W. F. Ames, *Numerical Methods for Partial Differential Equations*, Academic Press, 1977.
4. P. K. Khosla and S. G. Rubin, 'A conjugate gradient iterative method', *Computers and Fluids*, **9**, 109–121 (1981).
5. D. Kershaw, 'The incomplete Chulesky-conjugate gradient method for the iterative solution of systems of linear equations', *J. Comp. Phys.*, **26**, 43–65 (1978).
6. H. L. Stone, 'Iterative solution of implicit approximations of multidimensional partial differential equations', *SIAM J. Numer. Anal.*, **5**, 530–558 (1968).
7. A. Lin, 'The modified strongly implicit method for solving elliptic equations', *J. Comp. Appl. Math.*, 1982.
8. R. S. Varga, *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs NJ, MR 28H1725, 1962.
9. A. Lin and S. G. Rubin, 'A conditionally stable symmetric numerical scheme for solving convection–diffusion equations', in R. Shaw *et. al.* (eds) *Innovative Numerical Analysis in Applied Engineering*, University of Virginia Press, 1980.
10. P. K. Khosla, and S. G. Rubin, 'Diagonally dominant-second order accurate implicit scheme', *Computers and Fluids*, **2**, 207–209 (1974).
11. S. G. Rubin and A. Lin, 'Marching with the parabolized Navier–Stokes equations', *Isr. J. Tech.*, **18**, 21–31 (1981).
12. T. D. Taylor and E. Ndefo, *Proc. 2nd Int. Conf. Numer. Meth. Fluid Dynamics*, 1970, pp. 356–364.
13. M. Wolfshtein, 'Numerical smearing in one-sided difference approximations to the equations of non-viscous flow', Imperial College, Department of Mechanical Engineering, *Rep. EF/TN/A/3*, 1968.
14. R. T. Davis, 'Numerical methods for coordinate generation based on Schwarz–Cristoffel transformation', *AIAA Comp. Fluid Dynamic Conf.*, Williamsburg, Va., 1979.
15. P. R. Eisenman, 'Geometric methods in computational fluid dynamics', NASA Langley Research Centre, *ICASE Rep. 80–11*, 1980.
16. A. Lin, 'Applications of the optimal general implicit method', in preparation, 1982.
17. D. J. Atkins, S. J. Maskell and M. A. Patrick, 'Numerical predictions of separated flows' *Int. j. numer. mech. eng.*, **15**, 129–144 (1980).
18. R. E. Smith and Amy Kidd, 'Comparative study of two numerical techniques for the solution of viscous flow in a driven cavity', *NASA SP-378*, 1975.