

# Numerical Solution of Helmholtz's Equation by Implicit Capacitance Matrix Methods

WLODZIMIERZ PROSKUROWSKI

Lawrence Berkeley Laboratory

The numerical solution of Helmholtz's equation in an arbitrary bounded plane region is considered. Variants of the capacitance matrix method that greatly reduce storage requirements are developed. This allows the use of a very fine mesh with several hundred mesh points in each direction or the use of a computer with a small core storage.

**Key Words and Phrases:** Helmholtz's equation, fast Poisson solver, capacitance matrix, irregular regions

**CR Categories:** 5.17

## 1. INTRODUCTION

In recent years special techniques known as capacitance matrix methods have been developed for the numerical solution of Helmholtz's equation in a general planar bounded region  $\Omega$ ,

$$-\Delta u + cu = f \quad \text{in } \Omega,$$

where  $c$  is a real constant, and either Dirichlet or Neumann conditions are specified on the boundary  $\partial\Omega$ . These methods make use of fast solvers in regions that allow for the separation of variables. Operation count  $\theta(n, m)$  for such fast solvers is proportional to  $mn \log_2 n$ , where  $m$  and  $n$  are the number of mesh points in the two directions. For a detailed discussion of such methods and a history of their development, see Proskurowski and Widlund [9] and Widlund [12].

In this paper we confine ourselves to two-dimensional bounded regions; for problems in three dimensions see O'Leary and Widlund [6].

The relatively large amount of needed core storage limits the applicability of our algorithms, as previously developed. Normally, one generates and stores a dense capacitance matrix  $C$  of the order of  $p$ , where  $p$  is the number of mesh points inside  $\Omega$  adjacent to the boundary  $\partial\Omega$ . In this paper we develop an implicit method in which we avoid generating and storing the matrix  $C$ ; see also O'Leary and Widlund [6]. Moreover, we exploit the fact that only the mesh points in the close vicinity of the boundary  $\partial\Omega$  are involved in the main part of the computation,

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

This work was done with support from the U.S. Energy Research and Development Administration. Author's present address: Department of Mathematics, University of Southern California, Los Angeles, CA 90007.

© 1979 ACM 0098-3500/79/0300-0036\$00.75

that is, the capacitance matrix iterations. Using a solver on a rectangle, which takes the sparsity of the problem into account, developed originally by Banegas [1], we design an algorithm that requires only  $32p$  storage locations for its main part. Only two steps, computing the space potential term and final solution, are limited by the storage requirements of a fast solver on a rectangle, i.e.,  $m \cdot n$  locations, where  $m$  and  $n$  are the number of mesh points in a rectangle in which the region  $\Omega$  is imbedded. To remove that obstacle we propose a solver that requires only  $2nm^{1/2}$  storage locations at the expense of some computational effort.

For three-dimensional problems, where storage requirements are even more critical, the algorithm presented here could be implemented advantageously.

Numerical results from extensive experiments on a CDC 7600 computer are reported and a comparison of different programs is given.

## 2. CAPACITANCE MATRIX METHODS AND POTENTIAL THEORY

In this section we give a brief review of the potential theoretical approach leading to the capacitance matrix methods described in Proskurowski and Widlund [9]; see also Widlund [12].

We consider a problem on an arbitrary bounded planar region  $\Omega$ . The region  $\Omega$  is first imbedded in a rectangle, and a uniform mesh is introduced with the same mesh size in the two coordinate directions. The boundary conditions on the rectangle can be of arbitrary type as long as they allow for use of a fast solver; see Widlund [11] and Proskurowski and Widlund [9]. The set of mesh points is decomposed into three disjoint sets:  $\Omega_h$ ,  $\partial\Omega_h$ , and  $(C\Omega)_h$ . The set  $\Omega_h$  is the set of interior mesh points, i.e., each of its members has all its immediate neighbors in the open set  $\Omega$ . The remaining mesh points in  $\Omega$  constitute  $\partial\Omega_h$ , the set of irregular mesh points, while the set  $(C\Omega)_h$  contains all the remaining mesh points, i.e., the exterior mesh points. The Laplacian is represented by the five-point formula for all points in  $\Omega_h \cup (C\Omega)_h$ . The data for the exterior points are extended in an arbitrary way; for the proof that the solution on  $\Omega_h \cup \partial\Omega_h$  is independent of the solution and the data on  $(C\Omega)_h$ , see Proskurowski and Widlund [9, sec. 3]. For the irregular points we must introduce a formula that also takes the boundary conditions on  $\partial\Omega$  into account. We therefore combine the discrete Laplacian with an interpolation formula. The important problem of scaling these auxiliary equations is treated in detail in Proskurowski and Widlund [9] and Shieh [10]. We denote by  $A$  the  $n \cdot m \times n \cdot m$  matrix corresponding to the difference problem enlarged to a rectangle handling the given boundary conditions on  $\partial\Omega$ . The regularly structured problem for which a fast solver can be used is given by the  $n \cdot m \times n \cdot m$  matrix  $B$  representing the discrete Laplacian. With a proper ordering of the equations,  $A$  and  $B$  differ only in the rows corresponding to the irregular mesh points. For the Neumann problem we write  $A = B - UV^T$ , and for the Dirichlet problem  $A = B + UZ^T$ , where  $U$ ,  $V$ , and  $Z$  have  $p$  columns, and  $p$  is the number of irregular mesh points. The matrix  $U$  represents an extension operator, which maps  $\partial\Omega_h$  values onto the whole rectangle. It retains the values on  $\partial\Omega_h$  and makes the remaining values equal to zero. Its transpose,  $U^T$ , is a trace operator. Matrices  $V^T$  and  $-Z^T$  are compact representations of  $B - A$ , from which the zero rows corresponding to the regular mesh points have been deleted.

In potential theory the solution of the Neumann problem is given as a sum of a space potential  $u_s$  and a single layer potential of charge distribution  $\rho$  at the boundary  $\partial\Omega$ :

$$u(x) = u_s(x) + \mathcal{V}(x). \quad (2.1)$$

A discrete analog to eq. (2.1) is

$$u = Gf + GU\rho \quad (2.2)$$

where each of the  $p$  columns of  $U$  represents a unit charge placed at an irregular point, where the discrete operator  $G$  plays the same role as the integral operator defined by the fundamental solution of the continuous problem (see Proskurowski and Widlund [9]), and  $\rho$  is determined by solving the capacitance matrix equation

$$C\rho = (I_p - V^T GU)\rho = V^T Gf = g \quad (2.3)$$

where the  $p \times p$  matrix  $C$  is the capacitance matrix and  $\rho$  is a vector of  $p$  components. A proper approach for the Dirichlet problem is a double-layer potential  $\mathcal{V}$  of dipole density  $\mu$  at the boundary  $\partial\Omega$ :

$$u(x) = u_s(x) + \mathcal{V}(x). \quad (2.4)$$

A discrete analog to eq. (2.4) is

$$u = Gf + GD\mu, \quad (2.5)$$

where  $D$  has  $p$  columns, each of them representing a unit discrete dipole placed at an irregular point, and  $\mu$  is the solution of

$$C\mu = (I_p + Z^T GD)\mu = -Z^T Gf = g, \quad (2.6)$$

where  $\mu$  is a vector of  $p$  components. Shieh [10] has shown that the capacitance matrix  $C$  is equal to the sum of  $K_h$  and a matrix with a small condition number, where  $K_h$  is an approximation to the correct compact operator of the corresponding Fredholm integral equation of the second kind. The conjugate gradient method converges superlinearly for Fredholm integral equations of the second kind, as shown by Hayes [5]. Therefore, the conjugate gradient method applied for solving eqs. (2.3) and (2.6) converges rapidly; in practice, it is almost independent of the size of the mesh; see also Proskurowski and Widlund [9]. In summary, the algorithm consists of the following steps:

1. Generate the capacitance matrix  $C$ .
2. Compute  $g$ .
3. Solve eq. (2.3) or (2.6) by the conjugate gradient method
4. Use the fast solver to obtain  $u = G(f + U\rho)$  or  $u = G(f + D\mu)$ .

Another option for Step 3 is to factor  $C$  and to solve eq. (2.3) or (2.6) by Gaussian elimination. For the details of the algorithms and ways of fast generation of  $C$ , see Proskurowski and Widlund [9].

The total operation count for that algorithm is proportional to  $nm \log_2 n$  and  $p^2 \log_2 p$  if the conjugate gradient option is used.

Some alternatives to this algorithm that make it possible to avoid the explicit generation of  $C$  are described in the next sections.

### 3. AN IMPLICIT CAPACITANCE MATRIX METHOD

Methods in which we explicitly generate, store, and possibly factor the capacitance matrix may become inefficient when the mesh is refined. The capacitance matrix is a dense,  $p \times p$  matrix, where  $p$  is the number of irregular mesh points, which grows linearly with the number of mesh points in a coordinate direction. For example, somewhere between the values of  $p$  equal to 150 and 200 the small core memory (SCM) for the CDC 7600 computer becomes saturated. The CDC 7600 has two types of core storage: (1) the SCM which contains 65,536 decimal 60-bit words, and (2) the large core memory (LCM) which contains 512,000 decimal 60-bit words. The use of the LCM would allow us to increase the maximum values of  $p$  by a factor of 2 or slightly more, while for even larger  $p$  one must use a secondary memory device with a much longer access time. Therefore, we now present a method in which the capacitance matrix is used only implicitly without generating and storing it, thus saving  $p^2$  storage locations at the expense of a small increase in computational effort; see also O'Leary and Widlund [6], Widlund [12], and an early paper of George [4].

We describe the method for the Dirichlet boundary condition in which the double-layer potential approach is used. The method is very similar for the Neumann boundary condition in which a single-layer approach is used.

Once more we write the capacitance matrix equation

$$C\mu = (I_p + Z^T G D)\mu = -Z^T G f = g. \quad (3.1)$$

The capacitance matrix  $C$  can also be factored into the form

$$C = U^T A G D, \quad (3.2)$$

which we will use subsequently. Since matrix  $C$  is nonsymmetric and we intend to use the conjugate gradient method for solving eq. (3.1), we reformulate it in terms of a least squares problem

$$C^T C \mu = C^T g. \quad (3.3)$$

Thus each step of the conjugate gradient method requires the computation of a matrix-vector product  $C^T(Cx)$  for any vector  $x$  of length  $p$  given on the set of irregular mesh points, i.e.,

$$D^T G (U^T A)^T (U^T A) G D x. \quad (3.4)$$

Let us rewrite expression (3.4) as a sequence of equations:

$$\begin{aligned} x_1 &= D x, \\ x_2 &= G x_1 \quad \text{or} \quad B x_2 = x_1, \\ x_3 &= (U^T A) x_2, \\ x_4 &= (U^T A)^T x_3, \\ x_5 &= G x_4 \quad \text{or} \quad B x_5 = x_4, \\ y &= D^T x_5. \end{aligned} \quad (3.5)$$

Consequently, we first set an  $m \times n$  array equal to zero, then generate the mesh function  $Dx$  by setting up the discrete dipoles. This step costs  $2p$  multiplicative operations. Then we obtain  $G(Dx)$  by using the fast solver at a cost proportional to  $mn \log_2 n$  operations.  $U^T$  maps a mesh function defined for all mesh points into its restriction to irregular mesh points. Therefore, it is enough to apply the operator  $A$  to  $GDx$  only on the set of closest neighbors of irregular mesh points. Acting in this way we compute  $U^T A(GDx)$  at the expense of  $4p$  multiplicative operations. The part corresponding to the transpose of  $C$  is performed in a similar fashion. Thus the vector  $C^T Cx$  is obtained at a cost of two calls of the fast solver, proportional to  $mn \log_2 n$  operations, plus a lower-order term, proportional to  $p$  operations. In our program some operations were repeated in order to save storage space.

Summing up, this method requires  $mn + 8p$  storage locations ( $mn + 10p$  for the Neumann problem) compared with  $p^2 + nm$  plus a lower-order term proportional to  $p$  for the explicit capacitance matrix methods, as implemented in Proskurowski and Widlund [9]. On the other hand, the operation count for the present method is  $(2k + 3) \cdot \theta(m, n)$ , where  $k$  is the number of the conjugate gradient iterations and is proportional to  $\log_2 p$  and  $\theta(n, m)$  is the cost of a fast Helmholtz solver proportional to  $mn \log_2 n$ . By comparison, the cost for the explicit capacitance method with the conjugate gradient option is equal to  $3.5 \cdot \theta(n, m) + (2k + c) \cdot p^2$  operations, where  $c$  is a constant arising from the generation of the capacitance matrix.

An experimental comparison of the computation times for  $n = 64, p = 132$  and a circular region is given in Section 7. It shows that whenever the capacitance matrix is small enough to fit into the SCM, the explicit capacitance matrix methods are slightly faster. On the other hand, the present method does not make use of the translation invariance of the solution, which is exploited in our variant of the explicit capacitance matrix method, and alternative fast solvers might be easily used. A further development of this method is described in Section 5.

#### 4. A HELMHOLTZ SOLVER THAT TAKES ADVANTAGE OF SPARSITY

Consider the Helmholtz equation  $(-\Delta + c)u = f$  on a rectangular region with an  $m \cdot n$  mesh. Let the mesh values of  $f$  and  $u$  be called sources and targets, respectively. Denote then by  $s$  the number of nonzero sources and by  $t$  the number of targets where the solution is required. In our application we have either  $s \ll m \cdot n$  or  $t \ll m \cdot n$ , or both; see Section 5. We now describe how to make use of the sparsity of the sources and the fact that the solution is needed only at relatively few mesh points. A related idea for the case of a point source and targets on the lines parallel to the axes is mentioned in Buzbee and Dorr [2, p. 758]. The present method was developed by Banegas [1] and in our experiments we have been using a considerably altered variant of her algorithm. This algorithm was designed to be compatible with the one using fast Fourier transform (FFT) as described in Proskurowski and Widlund [9]. Nevertheless, there is no difficulty in adapting it to alternative fast Helmholtz solvers, if necessary. We remark also that there is no restriction on the location and number of sources and targets and that the only restriction on  $n$  is that it is even.

We first outline the fast Helmholtz solver described in [9]. The solution is

obtained by applying the FFT in one coordinate direction, i.e., for  $m$  vectors of length  $n$ , then by solving  $n$  very special tridiagonal systems of equations of order  $m$  by a Toeplitz method, and finally by using an inverse FFT on  $m$  vectors of length  $n$ . Storage requirements are here equal  $m \cdot n$  and the total operation count for this solver is  $\theta(m, n) = \frac{9}{2} mn \log_2(n/2) + 15mn$ , where each multiplication and each floating point addition is taken as a unit operation. We remark that the operation count always shows only a part of the actual computational expenses and is not an exact measure of it.

Let us now write the Fourier coefficients as inner products of the data vector  $f^{(l)}$  with the eigenvectors  $\phi^{(j)}$  of a certain block of the block matrix  $B$  representing the discrete Laplacian  $(-\Delta + c)$ :

$$\hat{f}_k^{(j)} = \frac{1}{n} \sum_{l=1}^n \phi_l^{(j)} \cdot f_k^{(l)}, \quad (4.1)$$

where  $j = 1, \dots, n$  and  $k = 1, \dots, m$ . In addition, the inverse Fourier coefficients can be written as

$$u_k^{(t)} = \sum_{j=1}^n \phi_l^{(j)} \cdot \hat{u}_k^{(j)}, \quad (4.2)$$

where  $l = 1, 2, \dots, n$  and  $k = 1, 2, \dots, m$ . For sparse sources the number of nonzero entries  $f_k^{(l)}$  in eq. (4.1) is reduced to  $s$ . Similarly, if the solution is needed at  $t$  mesh points we evaluate only  $t$  of the sums (4.2). Consequently, the operation count for the summation formulas (4.1) and (4.2) is reduced considerably. Moreover, we reorder the computation, performing it separately for each frequency  $j$ . We first compute the Fourier coefficients for all nonzero sources  $f_i$ ,  $i = 1, \dots, s$ , and simultaneously sum those having the same index  $k$ . We then solve the intermediate tridiagonal matrix problem (with  $\lambda_j$  corresponding to each frequency  $j$ ) with the previously computed Fourier coefficients as a right-hand side. We finally compute the inverse Fourier coefficients for all targets we need for the solution  $u_i$ ,  $i = 1, \dots, t$ , and sum simultaneously those having the same index  $k$ . At this point locations used for the temporary values of  $\hat{f}^{(j)}$  can be released and used for  $\hat{f}^{(j+1)}$ .

Thus, this procedure requires  $3(s + t)$  locations for sources, targets, and their coordinates, plus  $2m$  locations to store certain values of the sine and cosine functions and  $2n$  locations to store temporarily some of the Fourier coefficients. In our program we also store temporarily some indexes in order to avoid recomputing them. In all, we use

$$4(s + t) + 2(m + n) \quad (4.3)$$

storage locations. If the location of the sources and targets coincide, we can perform the computations in place (as we do in the fast solver using FFT), thereby further reducing these requirements to  $4s + 2(m + n)$  locations.

It is evident that for large  $s$  and  $t$  this procedure, which uses the conventional (i.e., slow) Fourier transform, will be much slower than a comparable solver using FFT. Now we will establish restrictions on  $s$  and  $t$  for this procedure to be competitive with a fast Helmholtz solver. The operation count is

$$\psi(n, m, s, t) = 3n \cdot (s + t) + 4mn. \quad (4.4)$$

Then  $\psi = \theta$  for  $s + t = \frac{3}{2} m \log_2 n/2 + \frac{11}{3} m$ . For example, take  $m = n$  and  $s + t = 10n$ . This gives the ratio  $\theta(n)/\psi(n) = \frac{9}{2} \log_2(n/2) + \frac{15}{34}$ , which is equal to 1.1 for  $n = 64$ . The corresponding ratio for execution times in numerical experiments (see Section 7) is very close to it.

## 5. AN IMPLICIT CAPACITANCE MATRIX METHOD USING A HELMHOLTZ SOLVER THAT EMPLOYS SPARSITY

We recall from Section 3 that the main computational effort (more than 90 percent) in an implicit capacitance matrix method goes for computing vectors  $y = C^T Cx = D^T G (U^T A)^T (U^T A) G D x$  during the conjugate gradient iterations. Moreover, a dominant part of this computation, also over 90 percent, is spent on a fast Helmholtz solver. We recall also that while distributing the discrete dipoles  $D$  and using the five-point stencils  $U^T A$ , only the mesh points from a close neighborhood of the  $p$  irregular mesh points are involved in the computation. The values at the rest of the mesh points are set to zero. In the second and fifth steps of the sequence of equations  $Bx_2 = x$ , and  $Bx_5 = x_4$ , the sources  $x$  and  $x_4$  are sparse ( $s$ ). Moreover, we need the solution  $x_2$  and  $x_5$  only at a few,  $t$ , mesh points.

A straightforward count for the Dirichlet problem gives  $s < 5p$  for  $x_4$  ( $t < 5p$  for  $x_2$ ) and  $s = 3p$  for  $x_1$  ( $t = 3p$  for  $x_5$ ). For the Neumann problem the corresponding values differ for  $x_1$  ( $s = p$ ) and  $x_5$  ( $t = p$ ), as we have here a single layer of charges instead of dipoles.

We observe that the coordinates of the mesh points involved in computation are often repeated as we go from one irregular mesh point to the next. For example, for  $U^T A$  only the layer of mesh points inside  $\Omega$  at a distance not larger than  $2h$  from the boundary  $\partial\Omega$  is used in computation. This gives the final value  $s + t < 5p$  for the Dirichlet problem and  $s + t \approx 3p$  for the Neumann problem. A comparison of the amount of computational effort, formula (4.4), shows that the use of Helmholtz solver that employs sparsity instead of a conventional one using FFT should be favorable here also.

When we use a two-dimensional array of entries the summation over the same coordinates (here double indexes) comes in a natural way. On the contrary, while using the Helmholtz solver described in Section 4, we work only with vectors of values and of coordinates of the entries. That is why we must construct an algorithm to recognize entries with the same coordinates in an effective way. To perform such a search in each conjugate gradient iteration would be costly. Therefore, in our computer implementation we preprocess the information about the irregular mesh points and their neighbors. It is performed only once at a cost proportional to the execution of  $p^2$  logical IF statements. This constitutes only a small part of the total computational effort; less than 2.5 percent for meshes  $n \geq 64$ . Additional storage for two vectors of an approximate total length of  $5p$  is required.

We now give a brief summary of the implicit capacitance algorithm:

1. Compute  $g = -ZGf$
2. Solve  $C^T C\mu = C^T g$  for  $\mu$
3. Compute  $u = G(f + D\mu)$ .

The capacitance matrix equation for the Dirichlet problem, in its normal form (Step 2) can be solved at the cost of  $2k \cdot (15np + 4mn)$  operations, where  $k$  is the

number of conjugate gradient iterations proportional to  $\log_2 p$ , while using only  $32p$  storage locations. For the Neumann problem, the corresponding values are  $2k \cdot (9np + 4mn)$  and  $25p$ .

Thus, if  $m = n$ , the cost of the main part of our Helmholtz solver is proportional to  $np \log_2 p$ . This is well confirmed by the experiments reported in Table II.

Assume for the moment that we solve only the Laplace equation. Then the total storage requirements for the present algorithm are proportional to  $p$ . This allows us to use a very fine mesh or to employ a computer with a small core storage.

In a general case, i.e., when  $f \neq 0$ , we must also compute the term  $u_s$  denoted as the space potential in Section 2. There we have both  $s$  and  $t$  equal to almost  $m \cdot n$ , and the sparse Helmholtz solver is quite ineffective (the operation count is proportional to  $m \cdot n^2$ ); hence in most cases it cannot be recommended. On the other hand, a standard fast solver requires  $m \cdot n$  storage locations, i.e., much more than needed for the rest of our algorithm. To resolve the last difficulty we have designed a fast Helmholtz solver that requires  $2nm^{1/2}$  storage locations, which is described in Section 6.

In the only process where a repetitive use of a Helmholtz solver is needed, Step 2, both the sources are sparse and the targets are needed at a few mesh points, as we have already seen. Therefore, a certain increase of computation time in Step 1 and in Step 3, in order to save storage (see Section 6), plays a lesser role in the total computational effort. Moreover, we can easily use LCM for the two-dimensional arrays needed in Step 1 and Step 3. LCM here will be accessed infrequently and therefore at a comparatively low cost. For CDC 7600 computers we can take advantage of the use of an inexpensive routine, MOVLEV, to manipulate  $n$  consecutive locations of the storage between SCM and LCM.

There are also applications where the solution or its gradient are required only on  $\partial\Omega_h$ . In those cases  $t \ll m \cdot n$  in Step 3. In order to retain flexibility we implemented the Helmholtz solver that employs sparsity together with the fast solver that uses FFT. Four different options were used for sparse and nonsparse  $s$  and  $t$ .

## 6. A METHOD OF SOLVING HELMHOLTZ'S EQUATION IN A RECTANGLE USING $2n^{3/2}$ STORAGE LOCATIONS

Standard fast methods for solving Helmholtz's equation on a rectangular region require  $n \cdot m$  storage locations. This poses a limitation on the size of the mesh if only fast core memory is to be used. If a disk is used to store the  $n \cdot m$  array for large  $n$  and  $m$ , then the standard fast solver requires frequent transfers to and from the disk of parts of the array. The long access times make this procedure very expensive.

We now propose a method that somewhat lessens the restriction on the mesh at the expense of an increase in the number of operations.

Let us consider an infinite parallel strip with periodic boundary conditions. Denote by  $n$  the number of mesh points across the strip and impose free-space boundary conditions. Denote then by  $m$  the number of mesh points along the strip. After the change of basis by the FFT on  $m$  lines of length  $n$  and a suitable permutation, we obtain  $n$  tridiagonal systems of equations of order  $m$  having the following special structure; see Proskurowski and Widlund [9]:



$$\begin{bmatrix} \mu & -1 & & & & \\ -1 & \lambda & -1 & & & \\ & & & \bigcirc & & \\ & & & & & \\ \bigcirc & & & & & \\ & & & & -1 & \lambda & -1 \\ & & & & & -1 & \mu \end{bmatrix} \cdot \hat{x} = \hat{f}. \quad (6.1)$$

where  $\mu = \lambda/2 + (\lambda^2/4 - 1)^{1/2}$ , i.e.  $\lambda = \mu + \mu^{-1}$ . For  $|\lambda| > 2$  we have the following explicit formula for the solution of our problem:

$$\hat{x}_i = \sigma \sum_j \mu^{-|i-j|} \cdot \hat{f}_j, \quad i, j = 1, \dots, m, \quad (6.2a)$$

where  $\sigma = (\mu - \mu^{-1})^{-1}$ , as can be verified by inspection. For  $|\lambda| = 2$  we choose

$$\hat{x}_i = -\frac{1}{2} \sum_j (\lambda/2)^{|i-j|+1} \cdot |i-j| \cdot \hat{f}_j, \quad (6.2b)$$

and for  $|\lambda| < 2$  we choose

$$\hat{x}_i = -\frac{1}{2} \sum_j [\sin(|i-j|\phi)/\sin\phi] \cdot \hat{f}_j, \quad (6.2c)$$

where  $\lambda = 2 \cos \phi$ ; see Proskurowski and Widlund [9]. Note that when  $|\lambda| \rightarrow 2$ , the expression in (6.2c) converges to the one in (6.2b).

We divide the strip lengthwise into  $k$  equal boxes and find the solution  $\hat{x}$  on the  $(k+1)$  lines connecting the boxes, in accordance with formulas (6.2). Then by taking the inverse FFT on those  $(k+1)$  lines, we obtain the solution  $x$  on them. We remark that this in itself is a cheap method of computing the solution to Helmholtz's equation on a sparse set of lines. Moreover, the summation in eqs. (6.2) needs to be taken only for those  $j$  for which  $\hat{f}_j$  is nonzero, in the case of sparse data  $f$ .

Finally, we consider each box of the size  $n \cdot m/k$  as a separate problem with Dirichlet boundary conditions across the strip (the values of  $x$  on  $(k+1)$  lines computed within the machine accuracy) and with periodic conditions on the shorter edges. By a suitable choice of the value  $m/k$ , we can use a standard fast method for each separate box.

Note that in order to save storage the procedures of computing  $\hat{f}$  by taking FFT on each line and of computing  $x$  on  $(k+1)$  lines must be carried out simultaneously in the same loop for all  $i = 1, \dots, m$ :

1. Take the FFT  $f^{(i)} \rightarrow \hat{f}^{(i)}$ .
2. For all chosen lines  $l = 1, \dots, k+1$ , add the terms  $\hat{f}_j^{(i)}$ ,  $j = 1, \dots, m$ , to  $\hat{x}^{(i)}$  in accordance with eq. (6.2)
3. Store  $\hat{x}^{(i)}$ ; release  $\hat{f}^{(i)}$ .

Note also that the data  $f$  are required twice (from a function subprogram or from a disk): in Step (1) above and while solving problems in each individual box. In both cases values of  $n$  consecutive locations of the  $n \cdot m$  array are used simultaneously, which simplifies the access from and to a disk. No intermediate

results are stored in an auxiliary storage in contrast to a straightforward extension of standard fast solvers to fine meshes.

Storage requirements for this algorithm are as follows:  $m$  temporary locations for  $\mu^{-i}$ ,  $i = 1, \dots, m - 1$  and  $\sigma$ , and  $n$  locations for  $\hat{f}$ , which all can be released after the first step is performed;  $n(k + 1)$  locations for the solution  $x$  on the  $(k + 1)$  chosen lines, and  $n \cdot m/k$  locations needed for solving  $k$  problems on small boxes. Thus totally we need  $n \cdot (m/k + k + 1)$  memory locations. The operations count is consequently:

- (1) FFT on  $m$  lines of length  $n$   $\frac{9}{4}mn \log_2(n/2) + \frac{11}{2}mn$ ,
- (2) computing  $\hat{x}$  on  $(k + 1)$  lines  $(k + 1)mn$ ,
- (3) inverse FFT on  $(k + 1)$  lines  $\frac{9}{4}(k + 1)n \log_2(n/2) + \frac{11}{2}n(k + 1)$ ;
- (4) solving the Dirichlet problem in  $k$  boxes.

$$k(\frac{9}{2}n(m/k) \log_2(m/2k) + 15(m/k) \cdot n).$$

Thus the total operation count is

$$\psi(m, n, k) = mn [\frac{9}{4} \log_2(n/2) + \frac{9}{2} \log_2(m/2k) + k + \frac{43}{2}]$$

plus a lower-order term. For comparison, the operation count for the standard fast solver is  $\theta(m, n) = mn (\frac{9}{2} \log_2(n/2) + 15)$ .

For  $k = m^{1/2}$  the storage requirements are minimal and equal to  $n(2m^{1/2} + 1)$ . For simplicity, take  $m = n$ . Then for the optimal choice of  $k$ , equal to  $n^{1/2}$ , we have the following increase in computational effort:

$n$	$\psi(n, n)/\theta(n, n)$	$n^2$	$2n^{3/2}$	$32p$
$2^8$	1.6	$2^{16}$	$2^{13}$	$2^{14}$
$2^{10}$	1.8	$2^{20}$	$2^{16}$	$2^{16}$
$2^{12}$	2.2	$2^{24}$	$2^{19}$	$2^{18}$

We have also tabulated the storage requirement for standard and present fast solvers, and the storage requirements for the capacitance matrix iterations, i.e.,  $n^2$ ,  $2n^{3/2}$ , and  $32p$  ( $p = 2n$ ), respectively. The values of the last two columns are of the same order of magnitude for mesh sizes considered.

This solver has not yet been tested experimentally.

## 7. NUMERICAL EXPERIMENTS

In this section we report the results of a series of numerical experiments that were carried out on the CDC 7600 computer at the Lawrence Berkeley Laboratory. In our experiments we have used programs that are now obtainable from LBL's computer library [7].

For studying our capacitance matrix methods we have chosen problems with no discretization error and regions that are circular.

We use the following notation:

- Variant 1 The capacitance matrix  $C$  is generated and factored, the capacitance matrix system is solved by Gaussian eliminations.
- Variant 2:  $C$  is explicitly generated; the linear system with  $C$  is solved by the conjugate gradient method.
- Variant 3 An implicit capacitance method, an FFT solver is used in each conjugate gradient iteration.
- Variant 4 An implicit capacitance method, a solver that employs sparsity is used in each conjugate gradient iteration.

A comparison of the performance of all four variants is presented in Table I for a comparatively crude mesh,  $64 \times 64$  points. Storage requirements for arrays in this case are roughly 24,000 locations for Variants 1 and 2, while only 5,000 to 8,000 locations are needed for Variants 3 and 4; see Table IB, where a more detailed comparison of storage requirements is reported. If the capacitance matrix fits into the core memory, Variant 2 is the fastest one. The only exception is Variant 4 for the Neumann problem, where the number of sources and targets for which the solution is required is quite low. If a repetitive use of a solver for a given geometry is needed [for instance, in nonseparable problems (see Concus et al. [3]), or in eigenvalue problems (see Proskurowski [6])], Variant 1 is recommended. Variant 1 is also superior when a very high accuracy of the solution is required. On the other hand, one must keep in mind that most often the discretization error is larger than  $1 \times 10^{-6}$ . To obtain the solution with normalized  $l_2$ -norm of the error,  $\|\epsilon\|_2 = (1/n^{1/2})(\sum \epsilon_i^2)^{1/2}$ , of  $1 \times 10^{-3}$  only six iterations for Variants 2, 3, and 4 were needed with a subsequent saving in execution time from  $0.85 \div 1.24$  to  $0.65 \div 0.71$  second for all these solvers.

In Table III we compare the separable solvers (on rectangular regions) we have used; one uses FFT and the other one employs sparsity. The number of sources

Table IA. CPU Time on CDC 7600 with FTN4 Compiler (opt = 2) for All Solvers  
(The region was a circle, the mesh  $64 \times 64$ . There were 132 boundary mesh points and 1,789 mesh points inside  $\Omega$ )

	Variant 1 Dirichlet	Variant 2 Dirichlet	Variant 3 Dirichlet	Variant 4	
				Dirichlet	Neumann
Generation of $C$	0.455*	0.445	—	—	—
Factorization of $C$	0.602	—	—	—	—
Total execution time in seconds	1.180*	0.855	1.243	1.147	0.602
Number of iterations	—	12	12	12	8
$l_2$ -norm of conjugate gradient residuals	—	$0.3 \times 10^{-6}$	$0.8 \times 10^{-6}$	$1.0 \times 10^{-6}$	$0.1 \times 10^{-6}$
$l_2$ -norm of error	$1.2 \times 10^{-12}$	$0.8 \times 10^{-6}$	$2.3 \times 10^{-6}$	$3.2 \times 10^{-6}$	$1.2 \times 10^{-6}$
Time per iteration	—	0.033	0.090	0.085	0.061
Time to solve an additional problem	0.138	0.410	1.243	1.116	0.582

\* Using single-layer Ansatz, one can generate  $C$  in 0.165 seconds, decreasing the total execution time to 0.905 second.

Table IB. Approximate Storage Requirements for Arrays in All Solvers

Mesh	$64 \times 64$	$128 \times 128$	$256 \times 256$
Variants 1 and 2	24,000	—	—
Variant 3	5,300	18,700	70,000*
Variant 4 as implemented	8,000	24,400	81,000†
Variant 4 using the solver of Section 6	5,000	11,200	25,000

\* Exceeds SCM of the CDC 7600 and therefore not used. The use of LCM would be costly here.

† Using LCM, see Tables II and III

Table II. Influence of Mesh Refinements on Number of Iterations and Execution Time  
 (CDC 7600 with FTN4 (opt = 2) compiler The region was a circle.)

	Variant 3			Variant 4		
	Dirichlet problem			Dirichlet problem		Neumann problem
Mesh	$64 \times 64$	$128 \times 128$		$64 \times 64$	$128 \times 128$	$256 \times 256$
Number of boundary points	132	268		132	268	540
Number of mesh points inside $\Omega$	1789	7209		1789	7209	28913
Number of conjugate gradient iterations	12	14	13	15	16	8
$\ell_2$ -norm of error	$2.3 \times 10^{-6}$	$2.0 \times 10^{-6}$	$1.0 \times 10^{-6}$	$1.0 \times 10^{-6}$	$1.3 \times 10^{-6}$	$1.2 \times 10^{-6}$
Total execution time in seconds	1.243	6.473	1.232	5.470	24.158	0.602
Share of separable solvers	1.155	6.087	1.095	5.038	22.941	0.536
Time per iteration	0.090	0.432	0.085	0.324	1.290	0.061

 Table III CPU Time on CDC 7600 with FTN4 Compiler (opt = 2) for  
 Solvers on Rectangular Regions

	Mesh	Execution time in seconds	Number of sources and targets
Solver using FFT	$64 \times 64$	0.043	—
	$128 \times 128$	0.203	—
	$256 \times 256$	1.235*	—
Solver employing sparsity	$64 \times 64$	0.027†	392
	$64 \times 64$	0.0385	600
	$128 \times 128$	0.156	1224
	$256 \times 256$	0.622	2472

\* Computed while using LCM, which slows down the solver by 25–30 percent.

† This result is for an application to a Neumann problem, all the rest for Dirichlet problems

$s$  and targets  $t$  given here is typical for Variant 4 and corresponds to experiments shown in Table II. Results of experiments with refined mesh confirm our conjecture that the execution time for the solver that employs sparsity is proportional to  $n(s + t)$  and thus is essentially of the order of  $n^2$  for Variant 4, while the corresponding growth factor for Variant 3 is  $mn \log_2 n$ .

Table IV. Influence of Mesh Refinements on Number of Conjugate Gradient Iterations for Circular Region with Circular Hole in Its Center

Mesh	$16 \times 16$	$32 \times 32$	$64 \times 64$
Number of boundary points	44	84	168
Number of mesh points inside $\Omega$	100	412	1680
Number of conjugate gradient iterations	16	16	16
$\ell_2$ -norm of conjugate gradient residuals	$1.7 \times 10^{-6}$	$1.7 \times 10^{-6}$	$1.6 \times 10^{-6}$

The results given in Table II show that the number of conjugate gradient iterations remains almost constant when the mesh is refined; more precisely, it grows as  $\log_2 p$ . We can also see clearly that the execution time per iteration is proportional to  $n^2$  for Variant 4 and to  $mn \log_2 n$  for Variant 3. Hence the total computational effort behaves similarly; when the number of inner points grows by a factor of 4.0, the total CPU time grows by a factor 4.4 for Variant 4 and by 5.2 for Variant 3, for the mesh sizes considered.

We have also run experiments on very fine meshes for the Laplace equation using Variant 4. The solution was obtained on a sparse set of mesh points close to  $\partial\Omega$ . We report here one of these experiments, with  $n = 650$ ,  $p = 1468$ ,  $s + t = 6728$ , and 212,201 mesh points inside  $\Omega$ . The number of conjugate gradient iterations was 18 (the  $\ell_2$ -norm of the error  $3 \times 10^{-6}$ ) and the total execution time was 157.6 second on the CDC 7600 with the FTN4 (opt = 2) compiler. Storage space needed for arrays here was about 47,000 locations.

Additionally, we chose a circular region with a circular hole in its center. The diameter of the hole was one-fourth of the diameter of the outer circle. The rate of convergence of the conjugate gradient iterations was in this case slower than for plain circles. Nevertheless, the influence of mesh refinement was insignificant; see Table IV. This seems to support our conviction that such behavior is to be expected for regions with smooth boundaries, i.e., without cusps, slits, etc.

#### ACKNOWLEDGMENTS

I would like to thank Olof Widlund for his consistent support and advice throughout my work, Alexandra Banegas and Dianne P. O'Leary for making available early versions of their programs, and Paul Concus for comments on the manuscript.

#### REFERENCES

1. BANEGAS, A. Fast Poisson solvers for problems with sparsity. *Math. Comp.* 32 (1978), 441-446.
2. BUZBEE, B. L., AND DORR, F. W. The direct solution of the biharmonic equation on rectangular regions and the Poisson equation on irregular regions. *SIAM J. Numer. Anal.* 11 (1974), 753-763.
3. CONCUS, P., GOLUB, G. H., AND O'LEARY, D. P. A generalized conjugate gradient method for the solution of elliptic partial differential equations. In *Proc. Symp. on Sparse Matrix Computation*, J. R. Bunch and D. J. Rose, Eds., Academic Press, New York, 1976.
4. GEORGE, J. A. The use of direct methods for the solution of the discrete Poisson equation on non-rectangular regions. Rep. 159, Computer Science Dep., Stanford U., Stanford, Calif., 1970.
5. HAYES, R. M. Iterative methods of solving linear problems in Hilbert space, contributions to the solution of systems of linear equations and the determination of eigenvalues, O. Taussky, Ed.,

- Nat. Bur. Stand. (Appl. Math. Ser. 39), pp. 71-103, 1954
- 6 O'LEARY, D P, AND WIDLUND, O ERDA-NYU Rep. To appear.
  - 7 PROSKUROWSKI, W Four Fortran programs for numerically solving Helmholtz's equation in an arbitrary bounded planar region Rep. 7516, Lawrence Berkeley Lab., Berkeley, Calif., 1978.
  - 8 PROSKUROWSKI, W. On the numerical solution of the eigenvalue problem of the Laplace operator by a capacitance matrix method. *Computing* 20 (1978), 139-151
  9. PROSKUROWSKI, W., AND WIDLUND, O On the numerical solution of Helmholtz's equation by the capacitance matrix method. *Math. Comp.* 30 (1976), 433-468.
  - 10 SHIEH, A. Fast Poisson solver on nonrectangular domains. New York U., Ph.D. Thesis, June 1976.
  11. WIDLUND, O. On the use of fast methods for separable finite difference equations for the solution of general elliptic problems. In *Sparse Matrices and Their Applications*, D.J. Rose and R A Willoughby, Eds , Plenum Press, New York, 1972
  12. WIDLUND, O. Capacitance matrix methods for Helmholtz's equation on general bounded regions. In Proc. Meeting in Oberwolfach, Lecture Notes in Mathematics, vol. 631, Springer-Verlag, New York, 1978, 209-219.

Received April 1977