# On a class of preconditioners for solving the Helmholtz equation ☆

## Y.A. Erlangga [*], C. Vuik, C.W. Oosterlee

*Department of Applied Mathematical Analysis, Delft University of Technology, Mekelweg 4, 2628 CD, Delft, The Netherlands*

Available online 5 March 2004

## Abstract

In 1983, a preconditioner was proposed [J. Comput. Phys. 49 (1983) 443] based on the Laplace operator for solving the discrete Helmholtz equation efficiently with CGNR. The preconditioner is especially effective for low wavenumber cases where the linear system is slightly indefinite. Laird [Preconditioned iterative solution of the 2D Helmholtz equation, First Year's Report, St. Hugh's College, Oxford, 2001] proposed a preconditioner where an extra term is added to the Laplace operator. This term is similar to the zeroth order term in the Helmholtz equation but with reversed sign. In this paper, both approaches are further generalized to a new class of preconditioners, the so-called "shifted Laplace" preconditioners of the form $\Delta\phi - \alpha k^2\phi$ with $\alpha \in \mathbb{C}$. Numerical experiments for various wavenumbers indicate the effectiveness of the preconditioner. The preconditioner is evaluated in combination with GMRES, Bi-CGSTAB, and CGNR.

© 2004 IMACS. Published by Elsevier B.V. All rights reserved.

*Keywords:* Helmholtz equation; Krylov subspace; Preconditioner

## 1. Introduction

In this paper, the time-harmonic wave equation in 2D heterogeneous media is solved numerically. The underlying equation governs wave propagations and scattering phenomena arising in acoustic problems in many areas, e.g., in aeronautics, marine technology, geophysics, and optical problems. In particular, we look for solutions of the Helmholtz equation discretized by using finite difference discretizations. Since the number of gridpoints per wavelength should be sufficiently large to result in acceptable solutions, for very high wavenumbers the discrete problem becomes extremely large, prohibiting the use of direct

---

methods. Iterative methods are the interesting alternative. However, Krylov subspace methods are not competitive without a good preconditioner. In this paper, we consider a class of preconditioners to improve the convergence of the Krylov subspace methods.

Various authors contributed to the development of powerful preconditioners for Helmholtz problems. The work in [2] and the follow-up investigation in [11] can be considered as the start for the class of preconditioners we are interested in. A generalization has been recently proposed in [14]. In [2,11,14], the preconditioners are constructed based on the Laplace operator. In [14], this operator is perturbed by a real-valued linear term. This surprisingly straightforward idea leads to very satisfactorily convergence. Furthermore, the preconditioning matrix allows the use of SSOR, ILU, or multigrid to approximate the inversion within an iteration.

In this paper, we will generalize the approach in [2,11,14]. We give theoretical and numerical evidence that introducing a *complex* perturbation to the Laplace operator can result in a better preconditioner than using a real-valued perturbation. We call the resulting class of preconditioners "shifted Laplace" preconditioners. This class of preconditioners is simple to construct and is easy to extend to inhomogeneous media.

There are various other types of preconditioners for general indefinite linear systems, e.g., [7,9,15,17]. In particular for Helmholtz problems, [9] proposed a class of preconditioners (so-called AILU) based on a parabolic factorization of the Helmholtz operator. In [15] another approach is pursued by perturbing the real part of the matrix to make it less indefinite. An interesting alternative is also described in [17], where a preconditioner based on the separation of variables is proposed. This preconditioner effectively accelerates the convergence for high wavenumbers.

This paper is organized as follows. In Section 2 we describe the mathematical model and the discretization used to solve wave propagation problems. Iterative methods used to solve the resulting linear system and the preconditioner will be discussed in Sections 3 and 4, respectively. In Section 5, we present the shifted Laplace preconditioners and show theoretically the convergence of this type of preconditioner. Numerical results are then presented in Section 6.

## 2. Mathematical model

We solve wave propagations in a two-dimensional medium with inhomogeneous properties in a unit (scaled) domain governed by the Helmholtz equation

$$\Delta\phi + k^2(x,y)\phi = f, \quad \Omega = [0,1]^2, \tag{1}$$

where $\Delta \equiv \partial^2/\partial x^2 + \partial^2/\partial y^2$, the Laplace operator, and $k(x,y) \in \mathbb{R}$ is the wavenumber, which depends on the spatial position in the domain. We consider a so-called "open problem", i.e., outgoing waves penetrate at least at one boundary without (spurious) reflections. To satisfy this condition, a radiation-type condition is imposed. Several formulations have been developed to model the nonreflecting condition at the boundary [1,4,5]. In this paper, the first order Sommerfeld condition is chosen of the form

$$\frac{\partial\phi}{\partial n} - \mathrm{i}k\phi = 0, \quad \text{on a part of } \Gamma = \partial\Omega, \tag{2}$$

with $n$ an outward direction normal to the boundary. Even though may not be sufficiently accurate for inclined outgoing waves [5], it is state-of-the-art in industrial codes, easy to implement in our discretization, and requires only a few gridpoints (as compared to, e.g., perfectly matched layer [3]).

In the implementation, a sufficiently large computational domain may help in reducing the effect of spurious reflections due to inaccurate boundary conditions.

To find numerical solutions of (1), the equation is discretized using the second-order difference scheme, in $x$-direction:

$$\frac{\partial^2 \phi}{\partial x^2} = \frac{1}{\Delta x^2}(\phi_{i-1} - 2\phi_i + \phi_{i+1}) + \mathcal{O}(\Delta x^2), \tag{3}$$

and similar in $y$-direction. The first-order derivative in (2) is discretized with the first-order forward scheme

$$\frac{\partial \phi}{\partial n} = \frac{1}{\Delta n}(\phi_{i-1} - \phi_i). \tag{4}$$

Substituting (3) and (4) into (1) and (2), one obtains a linear system

$$Ap = b, \quad A \in \mathbb{C}^{N \times N}, \tag{5}$$

where $A$ is a large, sparse symmetric matrix, with $N$ the number of gridpoints. Matrix $A$ is complex-valued and indefinite for large values of $k$. Throughout this paper, we say "$A$ is indefinite" if $A$ has eigenvalues with both positive and negative real part [7].

## 3. Krylov subspace method

For a large, sparse matrix, Krylov subspace methods are very popular. The methods are developed based on a construction of iterants in the subspace

$$\mathcal{K}^j(A, r_0) = \text{span}\{r_0, Ar_0, A^2 r_0, \dots, A^{j-1} r_0\}, \tag{6}$$

where $\mathcal{K}^j(A, r_0)$ is the $j$th Krylov subspace associated with $A$ and $r_0$ (see, e.g., [19]).

The basic algorithm within this class is the conjugate gradient method (CG) which has the nice properties that it uses only three vectors in memory and minimizes the error in the $A$-norm. However, the algorithm mainly performs well if the matrix $A$ is symmetric, and positive definite. In cases where one of these two properties is violated, CG may break down. For indefinite linear systems, CG can be applied to the normal equations since the resulting linear system becomes (positive) definite. Upon application of CG to the normal equations, CGNR [19] results. Using CGNR, the iterations are guaranteed to converge. The drawback is that the condition number of the normal equations equals the square of the condition number of $A$, slowing down the convergence drastically.

Some algorithms with short recurrences but without the minimizing property are constructed based on the bi-Lanczos algorithm [19]. Within this class, BiCG [6] exists and its modifications: CGS [21] and Bi-CGSTAB [22]. In BiCG, the Krylov subspace is constructed from the orthogonalization of two residual vectors based on actual matrix $A$ and its transpose $A^{\mathrm{T}}$. Accordingly, one extra matrix/vector multiplication and one transpose operation are needed. In CGS, the extra transpose operation is avoided. One can accelerate the convergence by squaring the polynomial. Whenever the convergence is smooth, CGS converges twice as fast as BiCG. However, if the BiCG iteration diverges, CGS also diverges twice as fast as BiCG. To stabilize CGS, rather than taking the square of the polynomial, another polynomial can be chosen and multiplied with the polynomial of $A$. This results in Bi-CGSTAB. In many cases, Bi-CGSTAB exhibits a smooth convergence behavior and often converges faster than CGS. Also within this class are QMR [8] and COCG [23].

MINRES [16] can also be used to solve indefinite symmetric linear systems, as well as its generalization to the nonsymmetric case, GMRES [20,19]. Both algorithms have the minimization property but GMRES uses long recurrences. GMRES has the advantage that theoretically the algorithm does not break down unless convergence has been reached. The main problem in GMRES is that the amount of storage increases as the iteration number increases. Therefore, the application of GMRES may be limited by the computer storage. To remedy this problem, a restarted version, GMRES($m$), can be utilized [20]. Since restarting removes the previous convergence history, GMRES($m$) is not guaranteed to converge. There is no specific rule to determine the restart parameter $m$. In cases characterized by superlinear convergence, $m$ should often be chosen very large which makes restarting much less attractive. Another way to remedy the storage problem in GMRES is by including a so-called "inner iteration" as in GMRESR [24] and FGMRES [18].

Since the convergence theory of GMRES is well established, in our numerical experiments mainly full GMRES is used. Of course, experiments then become very restrictive (in this paper, up to $k = 30$) and for large problems, restarts become necessary. Bi-CGSTAB, which requires much less storage but some more matrix/vector multiplications than GMRES, is also used for comparison. For completeness, since the underlying theory for the preconditioners is developed based on the normal equations [14], we also include the convergence results using CGNR in the last experiment.

## 4. Preconditioner

To improve the convergence of Krylov subspace methods, a preconditioner should be incorporated.

By left preconditioning, one solves a linear system premultiplied by a preconditioning matrix $M^{-1}$,

$$M^{-1}Ap = M^{-1}b. \tag{7}$$

Often, right preconditioning is used, i.e.,

$$AM^{-1}\tilde{p} = b, \tag{8}$$

where $\tilde{p} = Mp$. Both preconditionings show typically a very similar convergence behavior. However, for left preconditioning GMRES computes the residuals based on the preconditioned system. In contrast, for right preconditioning GMRES computes the actual residuals. This difference may affect the stopping criterion to be used (see discussions in [19]).

The best choice for $M^{-1}$ is the inverse of $A$, which is impractical. If $A$ is SPD, one can approximate $A^{-1}$ by one iteration of SSOR or multigrid. However, most practical wave problems result in an indefinite linear system, for which SSOR or multigrid are not guaranteed to converge (and do not converge).

In general, one can distinguish two approaches for constructing preconditioners: matrix-based and operator-based. Within the first class lie, e.g., incomplete LU (ILU) factorizations. Several ILU techniques have been developed with different choices of the tolerated fill-in in the sparsity pattern of $A$, e.g., zero fill-in ILU, or ILU with drop tolerance [15,13]. Another different approach but falling into this category is the approximate inverse (see, e.g., [19]). An example of an operator-based preconditioner is analytic ILU (AILU) [9], which is based on the continuous Helmholtz operator.

In the next sections, we will briefly discuss some preconditioners for Helmholtz problems.

### 4.1. ILU preconditioner

An ILU preconditioner can be constructed by performing Gauss elimination and dropping some elements based on certain criteria. One can, e.g., drop all elements except for those in the same diagonals as the original matrix. This leads to ILU(0). ILU($p$) allows fill-in in $p$ additional diagonals. One can also drop elements which are smaller than a specified value, giving ILU(*tol*). In applications involving $M$-matrices, this class of preconditioners is sufficiently effective. However, preconditioners from this class are not effective for general indefinite problems. Reference [9] shows some results in which ILU-type preconditioners are used to solve the Helmholtz equation using QMR. For high wavenumbers $k$, ILU(0) converges slowly, while ILU(*tol*) encounters storage problems (and also slow convergence). For sufficiently high wavenumbers $k$, the cost to construct the ILU(*tol*) factors may become very high.

Instead of constructing the ILU factors from $A$, the Helmholtz operator $\mathcal{L}_h = \Delta + k^2$ can be used to set up *ILU-like* factors in so-called analytic ILU (AILU) [9]. Starting with the Fourier transform of the analytic operator in one direction, one constructs parabolic factors of the Helmholtz operator consisting of a first order derivative in one direction and a nonlocal operator. To remove the nonlocal operator, a localized approximation is proposed, involving optimization parameters. Finding a good approximation for inhomogeneous problems is the major difficulty in this type of preconditioner. This is because the method is sensitive with respect to small changes in these parameters. The optimization parameters depend on $k(x, y)$.

### 4.2. Shifted Laplace preconditioner

Another approach is found in *not* looking for an approximate inverse of the discrete indefinite operator $A$, but merely looking for a form of $M$, for which $M^{-1}A$ has satisfactory properties for Krylov subspace acceleration. A first effort to construct a preconditioner in such a way is in [2]. An easy-to-construct $M = \Delta$ preconditioner is incorporated for CGNR. One SSOR iteration is used whenever operations involving $M^{-1}$ are required. The subsequent work on this preconditioner with multigrid was done in [11].

Instead of the Laplace operator as the preconditioner, [14] investigates possible improvements if an extra term $-k^2$ is added to the Laplace operator. So, the Helmholtz equation with reversed sign is proposed as the preconditioner $M$. This preconditioner is then used in CGNR. One multigrid iteration is employed whenever $M^{-1}$ must be computed. GMRES can also solve the preconditioned linear system efficiently in less arithmetic operations. In the case of large $k$, the storage problem of full GMRES can be overcome, e.g., by applying GMRES($m$) or GMRESR. Bi-CGSTAB does not always perform satisfactorily as in [14]. (See also results in Section 6.)

In the next section, we concentrate on this type of preconditioners and present a generalization.

## 5. Spectral properties of shifted Laplace preconditioners

In this section we provide some analysis to understand the performance of the shifted Laplace preconditioners. The analysis is based on eigenvalue properties of the preconditioned system. It is often that the eigenvalue distribution can help in understanding the behavior of CG-like iterations. Since the spectra of $M^{-1}A$ and $AM^{-1}$ are identical, we concentrate on left preconditioning.

## 5.1. Real shifted–Laplace preconditioner

The preconditioners in [2,14] can be motivated as follows. Consider the continuous 1D Helmholtz equation, subject to discretization. For simplicity, suppose that both boundary conditions are either Dirichlet or Neumann conditions.

We first consider the eigenvalues for the 1D Helmholtz operator without any preconditioning. Eigenvalues of this standard problem, denoted by $\lambda^s$, are found to be

$$\lambda_n^s = k_n^2 - k^2, \quad k_n = n\pi, \, n \in \mathbb{N} \setminus \{0\}. \tag{9}$$

In (9), $k_n$ is the natural frequency of the system. (We use $n$ to indicate the eigenmodes). If one considers the modulus of the eigenvalues (which in this case is simply their absolute value), it is easily seen that $|\lambda|$ becomes unbounded if either $n$ or $k$ are large. If the $l_2$-condition number $\kappa = |\lambda_{\max}/\lambda_{\min}|$ is used to evaluate the quality of eigenvalue clustering, one concludes also that for any sufficiently small $\lambda_{\min}$ the condition number is extremely large. Now, suppose an operator of the form

$$\frac{\mathrm{d}^2}{\mathrm{d}x^2} - \alpha k^2, \quad \alpha \in \mathbb{R}, \, \alpha \geqslant 0 \tag{10}$$

is used as a preconditioner, which is later on discretized, and with the same boundary conditions as those in the Helmholtz equation are imposed. The following generalized eigenvalue problem is obtained, i.e.,

$$\left( \frac{\mathrm{d}^2}{\mathrm{d}x^2} + k^2 \right) \phi_v = \lambda \left( \frac{\mathrm{d}^2}{\mathrm{d}x^2} - \alpha k^2 \right) \phi_v, \quad x \in [0,1] \subseteq \mathbb{R}. \tag{11}$$

For (11), we find the eigenvalues to be

$$\lambda_n = \frac{k_n^2 - k^2}{k_n^2 + \alpha k^2} = \frac{1 - (k/k_n)^2}{1 + \alpha (k/k_n)^2}, \quad n \in \mathbb{N} \setminus \{0\}. \tag{12}$$

For $n \to \infty$, $\lambda_n \to 1$, i.e., the eigenvalues are bounded above by one. Examining the low eigenmodes, for $k_n \to 0$, we obtain $\lambda \to -1/\alpha$. This eigenvalue remains below one unless $\alpha \leqslant 1$. The maximum eigenvalue can thus be written as

$$|\lambda_{\max}| = \max \left( \left| \frac{1}{\alpha} \right|, 1 \right), \quad \alpha \geqslant 0 \in \mathbb{R}. \tag{13}$$

To estimate the minimum eigenvalue, one can use a simple but rough analysis as follows. It is assumed that the minimum eigenvalue is very close (but not equal) to zero. This assumption implies a condition $k_j \approx k$ as obtained from (12). To be more precise, let $k_j = k + \varepsilon$, where $\varepsilon$ is any small number. If this relation is substituted into (12), and if higher order terms are neglected, and $\varepsilon k \ll k^2$ is assumed, then we find

$$\lambda_{\min} = \frac{2}{1 + \alpha} \left( \frac{\varepsilon}{k} \right). \tag{14}$$

From (14), the minimum eigenvalue can be very close to zero as $\alpha$ goes to infinity. The condition number of the preconditioned Helmholtz operator now reads

$$\kappa = \begin{cases} \frac{1}{2}(1 + \alpha)k/\varepsilon & \text{if } \alpha \geqslant 1, \\ \frac{1}{2\alpha}(1 + \alpha)k/\varepsilon & \text{if } 0 \leqslant \alpha \leqslant 1. \end{cases} \tag{15}$$

If $\alpha \geqslant 1$, $\kappa$ is a monotonically increasing function with respect to $\alpha$. The best choice is $\alpha = 1$, which gives minimal $\kappa$. If $0 \leqslant \alpha \leqslant 1$, $\kappa$ is a monotonically decreasing function with respect to $\alpha$. $\kappa$ is minimal in this range if $\alpha = 1$. In the limit sense we find that

$$\lim_{\alpha \downarrow 1} \kappa = \lim_{\alpha \uparrow 1} \kappa = k/\varepsilon, \tag{16}$$

which is the minimum value of $\kappa$ for $\alpha \geqslant 0 \in \mathbb{R}$.

## 5.2. Generalization to complex $\alpha$

The analysis on 1D shifted Laplace preconditioners for $\alpha \in \mathbb{R}$ gives $\alpha = 1$ as the optimum case. The nice property of the real shifted Laplace operator, at least in 1D, is that the eigenvalues have an upper bound. However, this property does not guarantee that the eigenvalues are favourably distributed. There is still the possibility that one or some eigenvalues can be very close to zero. For example, setting $\alpha = 1$ gives the minimal $\kappa$ but, at the same time, results in $\lambda_{\min}$ which is not better than that for $\alpha = 0$. We can improve the preconditioner by still preserving the upper boundedness and at the same time shifting the minimum eigenvalue as far as possible from zero. In this section, we generalize $\alpha$ to be complex-valued.

Consider the minimum eigenvalue $\lambda_{\min}$ obtained from the 1D problem (14). We may shift this eigenvalue away from zero by adding some real values to $\lambda$. In general, this addition will shift all eigenvalues, which is undesirable. An alternative is by multiplying the eigenvalues by a factor. From (12) the relation between eigenvalues for $\alpha = 0$ and $\alpha = 1$ can be derived, i.e.,

$$(\lambda_{\alpha=1})_n = \frac{1}{1 + (k/k_n)^2} (\lambda_{\alpha=0})_n. \tag{17}$$

Eq. (17) indicates that $\lambda_{\alpha=0}$ is scaled by a factor $1/(1 + (k/k_n)^2)$ to obtain $\lambda_{\alpha=1}$. Similarly, using (14), we obtain the following relation:

$$(\lambda_{\alpha=1})_{\min} = \frac{1}{2} (\lambda_{\alpha=0})_{\min}. \tag{18}$$

Since the eigenvalues of a general matrix may be complex, relation (17) can be considered as a particular case of scaling of the eigenvalues along the real axis in the complex plane. Our attempt to improve the clustering is by possibly introducing additional shift along the imaginary axis which moves the small eigenvalues further from zero. For that purpose, we introduce a complex coefficient of the form $\alpha + i\beta$, and consider a more general complex-valued shifted Laplace operator

$$\frac{\mathrm{d}^2}{\mathrm{d}x^2} - (\alpha + i\beta)k^2, \quad \alpha \in \mathbb{R}, \ \alpha \geqslant 0, \ \beta \in \mathbb{R}. \tag{19}$$

Eigenvalues of the premultiplied equation, denoted by $\lambda^c$, are

$$\lambda^c = \frac{k_n^2 - k^2}{k_n^2 + (\alpha + i\beta)k^2} \implies |\lambda^c|^2 = \frac{(k_n^2 - k^2)^2}{(k_n^2 + \alpha k^2)^2 + \beta^2 k^4}. \tag{20}$$

Evaluating $\lambda_{\max}$ and $\lambda_{\min}$ as in (13) and (14) one finds

$$|\lambda_{\max}^c|^2 = \max\left(\frac{1}{\alpha^2 + \beta^2}, 1\right), \qquad |\lambda_{\min}^c|^2 = \frac{4}{(1 + \alpha)^2 + \beta^2}\left(\frac{\varepsilon}{k}\right). \tag{21}$$
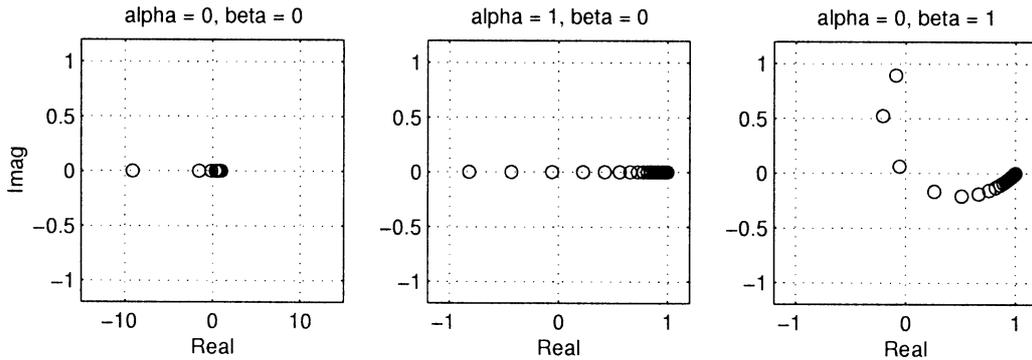
Fig. 1. Generalized eigenvalues of the continuous 1D Helmholtz equation, $k = 10$.

These results give the following condition numbers

$$\kappa^2 = \begin{cases} \frac{1}{4}\left(1 + \frac{1+2\alpha}{\alpha^2+\beta^2}\right)(k/\varepsilon)^2, & \alpha^2 + \beta^2 \leqslant 1, \\ \frac{1}{4}\left((1+\alpha)^2 + \beta^2\right)(k/\varepsilon)^2, & \alpha^2 + \beta^2 \geqslant 1. \end{cases} \tag{22}$$

Since $\alpha^2 + \beta^2$ is nonnegative, for any given $\alpha$ taking the circle $\alpha^2 + \beta^2 = 1$ in the first expression in (22) provides the smallest $\kappa^2$. Likewise, for any given $\alpha \geqslant 0$, $\kappa^2$ is minimal for the second expression in (22) whenever $\alpha^2 + \beta^2 = 1$. (One can verify that there is no other circle giving $\kappa^2$ lower than that on the circle with radius one. This can be seen, e.g., by introducing condition $\alpha^2 + \beta^2 = 1 + \varepsilon_1$, $\varepsilon_1 \geqslant 0$). With condition $\alpha^2 + \beta^2 = 1$, $\kappa$ is minimal if one takes $\alpha = 0$, implying $\beta = 1$. This combination gives the lowest condition number possible for the shifted–Laplace preconditioner for the 1D model problem.

Fig. 1 shows spectra of the preconditioned systems of the 1D Helmholtz problem using $M_{\alpha=0,\beta=0}$, $M_{\alpha=1,\beta=0}$, and $M_{\alpha=0,\beta=1}$ for our 1D problem. For simplicity, we denote these preconditioners as $M_0$, $M_1$, and $M_i$, respectively. Fig. 1 shows that the preconditioner $M_i$ clusters the eigenvalues stronger than $M_1$ and pushes the eigenvalues in the negative real plane towards the imaginary axis. This clustering may improve the performance of the preconditioned iterative methods. However, with this preconditioner there is still a possibility that some eigenvalues lie very close to zero causing unsatisfactory numerical performance. To estimate the position of these minimum eigenvalues, we consider the real part of (20). Similar as in (14), one finds that

$$\text{Re}(\lambda_{\min}^c) = \varepsilon/k. \tag{23}$$

This estimate is the same as the estimate for $M_1$ and smaller than that for $M_0$. However, the modulus $|\lambda_{\min}^c| = \sqrt{2}(\varepsilon/k) > |\lambda_{\min}^{\alpha=1}| = \varepsilon/k$ because of the imaginary shift (see Fig. 2). Because of the same upper bound as $M_1$, $M_i$ may perform better than $M_0$ and $M_1$.

In Fig. 2, a comparison of the modulus of eigenvalues for $k = 10$ is shown, indicating boundedness of eigenvalues of $M_1$ and $M_0$ near $|\lambda| = 0$. The right-hand figure zooms in to show the minimum $|\lambda|$. Evidently, $M_i$ has small eigenvalues with the modulus slightly larger than $M_1$, but smaller than $M_0$.
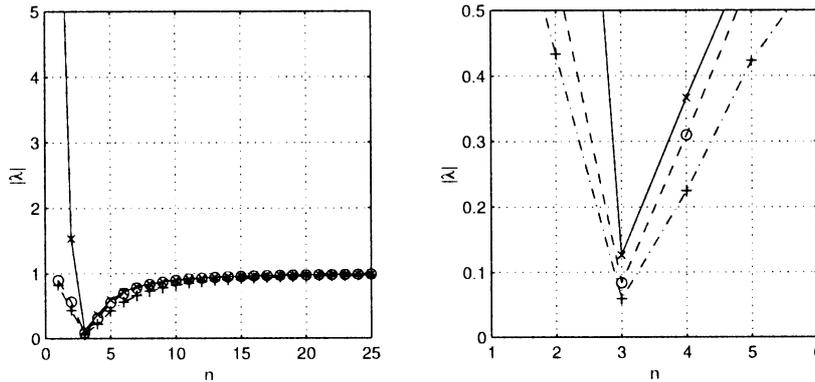
Fig. 2. The modulus of eigenvalues of the continuous 1D Helmholtz equation. $k = 10$ and $h^{-1} = 100$ for various preconditioners: $M_0(\times)$, $M_1(+)$, $M_i(\circ)$.

### 5.3. Spectrum of the discrete Helmholtz equation

We extend the analysis to the discrete formulation of (1). Suppose that the Helmholtz equation is discretized using (3), we arrive at the linear system, $Ap = b$. Matrix $A$ can be splitted into two parts: the Laplace component $B$ and the additional diagonal term $k^2 I$ so that $A = B + k^2 I$ and therefore

$$(B + k^2 I)p = b. \tag{24}$$

In the present analysis, we assume only Dirichlet or Neumann conditions at the boundaries in order to keep the matrix $A$ real-valued. Since $A$ is symmetric, the eigenvalues are all real-valued. This assumption of course is irrelevant for exterior problems we are interested in. However, this simplifies the analysis. The inclusion of a Sommerfeld condition will lead to a different result. Our numerical experiments on exterior problems, however, show consistency with the analysis based on the interior problem.

We precondition (24) using $M = B - (\alpha + i\beta)k^2 I$, with the same boundary conditions as for $A$. This gives

$$\left(B - (\alpha + i\beta)k^2 I\right)^{-1}\left(B + k^2 I\right)p = \left(B - (\alpha + i\beta)k^2 I\right)^{-1}b. \tag{25}$$

The generalized eigenvalue problem of (25) is accordingly

$$\left(B + k^2 I\right)p_v = \lambda_v\left(B - (\alpha + i\beta)k^2 I\right)p_v. \tag{26}$$

Both systems (25) and (26) are indefinite if $k^2$ is larger than the smallest eigenvalue of $B$. In such the case, the convergence is difficult to estimate. Therefore, the subsequent analysis will be based on the normal equations formulation of the preconditioned matrix system (as in [14]).

Denote the ordered eigenvalues of $B$ as $0 < \mu_1 \leqslant \mu_2 \leqslant \cdots \leqslant \mu_N$. We find the eigenvalues of the four following cases:

$$\lambda\left(A^* A\right) = \left(\mu_j - k^2\right)^2, \tag{27}$$

$$\lambda\left((M_0^{-1}A)^*(M_0^{-1}A)\right) = \left(\frac{\mu_j - k^2}{\mu_j}\right)^2 = \left(1 - \frac{k^2}{\mu_j}\right)^2, \tag{28}$$

$$\lambda\left((M_1^{-1}A)^*(M_1^{-1}A)\right) = \left(\frac{\mu_j - k^2}{\mu_j + k^2}\right)^2 = \left(1 - \frac{2k^2}{\mu_j + k^2}\right)^2, \tag{29}$$

$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big) = \left(\frac{\mu_j - k^2}{\mu_j + \mathrm{i}k^2}\right)\overline{\left(\frac{\mu_j - k^2}{\mu_j + \mathrm{i}k^2}\right)} = 1 - \frac{2\mu_j k^2}{\mu_j^2 + k^4}. \tag{30}$$

We first consider the case where $k$ is small such that $0 < k^2 < \mu_1$, where $\mu_1$ the smallest eigenvalue of $B$. Using (27)–(30), we find the minimal and maximal eigenvalues as follows:

$$\lambda\big((A^*A)\big)_{\min} = (\mu_1 - k^2)^2,$$
$$\lambda\big((A^*A)\big)_{\max} = (\mu_N - k^2)^2, \tag{31}$$

$$\lambda\big((M_0^{-1}A)^*(M_0^{-1}A)\big)_{\min} = \left(1 - \frac{k^2}{\mu_1}\right)^2,$$
$$\lambda\big((M_0^{-1}A)^*(M_0^{-1}A)\big)_{\max} = \left(1 - \frac{k^2}{\mu_N}\right)^2, \tag{32}$$

$$\lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\min} = \left(1 - \frac{2k^2}{\mu_1 + k^2}\right)^2,$$
$$\lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\max} = \left(1 - \frac{2k^2}{\mu_N + k^2}\right)^2, \tag{33}$$

$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\min} = 1 - \frac{2\mu_1 k^2}{\mu_1^2 + k^4},$$
$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\max} = 1 - \frac{2\mu_N k^2}{\mu_N^2 + k^4}. \tag{34}$$

Since $k^2/\mu_1 < 1$, one easily sees that

$$\lambda\big((M_0^{-1}A)^*(M_0^{-1}A)\big)_{\min} > \lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\min}.$$

As $n \to \infty$, one finds also that

$$\lim_{\mu_N \to \infty} \lambda\big((M_0^{-1}A)^*(M_0^{-1}A)\big)_{\max} = \lim_{\mu_N \to \infty} \lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\max} = 1.$$

With respect to the $l_2$-condition number, it becomes evident that for $k < \sqrt{\mu_1}$, preconditioning with $M_0$ gives a lower condition number than preconditioning with $M_1$. Hence, for small $k$, $M_0$ is more effective than $M_1$.

For $M_i$, one can compute that

$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\min}\Big/\lambda\big((M_0^{-1}A)^*(M_0^{-1}A)\big)_{\min} = \frac{(\mu_1 + k^2)^2}{\mu_1^2 + k^4} > 1,$$

$$\lim_{\mu_N \to \infty} \lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\max} = 1.$$

So, for $k < \sqrt{\mu_1}$, compared to $M_i$, preconditioning with $M_0$ still gives a better condition number. In cases where $k$ is small, $M_0$ is more effective than $M_1$ and $M_i$.

We consider now the case where $k$ is large, such that $\mu_1 < k^2 < \mu_N$. For the standard $A^*A$, one finds that

$$\lambda\big(A^*A\big)_{\min} = \big(\mu_{m_1} - k^2\big)^2, \quad \text{where } \big|\mu_{m_1} - k^2\big| \leqslant \big|\mu_j - k^2\big|, \ \forall j,$$

$$\lambda\big(A^*A\big)_{\max} = \big(\mu_N - k^2\big)^2. \tag{35}$$

The eigenvalues are unbounded either for large $\mu_N$ or large $k$.

For the preconditioned system $(M_0^{-1}A)^*(M_0^{-1}A)$ one finds

$$\lambda\big(\big(M_0^{-1}A\big)^*\big(M_0^{-1}A\big)\big)_{\min} = \left(\frac{\mu_{m_2} - k^2}{\mu_{m_2}}\right)^2, \quad \text{where } \left|\frac{\mu_{m_2} - k^2}{\mu_{m_2}}\right| \leqslant \left|\frac{\mu_j - k^2}{\mu_j}\right|, \ \forall j,$$

$$\lambda\big(\big(M_0^{-1}A\big)^*\big(M_0^{-1}A\big)\big)_{\max} = \max\left(\left(\frac{\mu_N - k^2}{\mu_N}\right)^2, \left(\frac{\mu_1 - k^2}{\mu_1}\right)^2\right). \tag{36}$$

In this case, there will be a possible boundedness for large $\mu_N$, i.e., for $\mu_N \to \infty$, $\lambda_N = 1$ as long as $k$ is finite (because $\lim_{k\to\infty}((\mu_j - k^2)/(\mu_j))^2 = \infty$). Furthermore, $\lim_{\mu_1\to 0}((\mu_1 - k^2)/(\mu_1))^2 = \infty$. Therefore, $\lambda_{\max}$ can become extremely large, which makes $M_0$ less favorable for preconditioning.

For the preconditioned system $(M_1^{-1}A)^*(M_1^{-1}A)$, one finds that

$$\lambda\big(\big(M_1^{-1}A\big)^*\big(M_1^{-1}A\big)\big)_{\min} = \left(\frac{\mu_{m_3} - k^2}{\mu_{m_3} + k^2}\right)^2, \quad \text{where } \left|\frac{\mu_{m_3} - k^2}{\mu_{m_3} + k^2}\right| \leqslant \left|\frac{\mu_j - k^2}{\mu_j + \mu_{m_3}}\right|, \ \forall j,$$

$$\lambda\big(\big(M_1^{-1}A\big)^*\big(M_1^{-1}A\big)\big)_{\max} = \max\left(\left(\frac{\mu_N - k^2}{\mu_N + k^2}\right)^2, \left(\frac{\mu_1 - k^2}{\mu_1 + k^2}\right)^2\right). \tag{37}$$

From (37), it is found that

$$\lim_{\mu_N\to\infty} \left(\frac{\mu_N - k^2}{\mu_N + k^2}\right)^2 = \lim_{\mu_1\to 0} \left(\frac{\mu_1 - k^2}{\mu_1 + k^2}\right)^2 = \lim_{k\to\infty} \left(\frac{\mu_j - k^2}{\mu_j + k^2}\right)^2 = 1. \tag{38}$$

The preconditioned system $M_1^{-1}A$ is always bounded above by one, i.e., the eigenvalues are always clustered. If $\lambda_{\min}$ in the case of preconditioning with $M_0$ and $M_1$ is of the same order of magnitude, then boundedness in case of $M_1$ provides a better condition number than $M_0$. For large $k$, $M_1$ is more effective than $M_0$.

Finally, we are looking at the complex shifted preconditioned system with $M_i$. One finds that

$$\lambda\big(\big(M_i^{-1}A\big)^*\big(M_i^{-1}A\big)\big)_{\min} = \frac{(\mu_{m_4} - k^2)^2}{\mu_{m_4}^2 + k^4}, \quad \text{where } \left|\frac{(\mu_{m_4} - k^2)^2}{\mu_{m_4}^2 + k^4}\right| \leqslant \left|\frac{(\mu_j - k^2)^2}{\mu_j^2 + k^4}\right|, \ \forall j,$$

$$\lambda\big(\big(M_i^{-1}A\big)^*\big(M_i^{-1}A\big)\big)_{\max} = \max\left(1 - \frac{2\mu_1 k^2}{\mu_1^2 + k^4}, 1 - \frac{2\mu_N k^2}{\mu_N^2 + k^4}\right). \tag{39}$$

The following results follow from (39):

$$\lim_{\mu_N\to\infty} \lambda\big(\big(M_i^{-1}A\big)^*\big(M_i^{-1}A\big)\big)_{\max} = \lim_{\mu_1\to 0} \lambda\big(\big(M_i^{-1}A\big)^*\big(M_i^{-1}A\big)\big)_{\max}$$

$$= \lim_{k\to\infty} \lambda\big(\big(M_i^{-1}A\big)^*\big(M_i^{-1}A\big)\big)_{\max} = 1. \tag{40}$$

Hence, the eigenvalues of $(M_i^{-1}A)^*(M_i^{-1}A)$ are always bounded above by one. Typically, preconditioning with $M_i$ gives a better condition number than with $M_0$.

To compare $M_i$ with $M_1$ we need to estimate the lower bound. In doing this, we assume that $\lambda_{\min} \approx 0$ implying $\mu_m = k^2 + \varepsilon$, $\varepsilon > 0$. After substituting this relation to (39), one finds that

$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\min} = \frac{1}{2}\frac{\varepsilon^2}{k^4}. \tag{41}$$

For $M_1$ we find that

$$\lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\min} = \frac{1}{4}\frac{\varepsilon^2}{k^4}. \tag{42}$$

Therefore,

$$\lambda\big((M_i^{-1}A)^*(M_i^{-1}A)\big)_{\min} = 2\lambda\big((M_1^{-1}A)^*(M_1^{-1}A)\big)_{\min}. \tag{43}$$

With respect to the $l_2$-condition number, one finds that

$$\kappa\big((M_i^{-1}A)^*(M_i^{-1}A)\big) = 2\left(\frac{k^4}{\varepsilon^2}\right) < \kappa\big((M_1^{-1}A)^*(M_1^{-1}A)\big) = 4\left(\frac{k^4}{\varepsilon^2}\right).$$

Considering the above result, we conclude that $M_i$ is more effective as the preconditioner than $M_1$.

**Remark.** For an interior problem where the resulting linear system is real-valued, using complex shift preconditioner requires more arithmetic operations. In this situation, it is possible that the gain in the convergence speed-up is overruled by the extra costs of the complex arithmetic operations.

## 6. Numerical results

We provide some numerical results for solving Eq. (1), and present three cases as the model problems

(i) a 2-D closed-off problem with Dirichlet conditions at all boundaries,
(ii) a 2-D open problem in a homogeneous medium with Sommerfeld conditions on a part of the boundary, and
(iii) a 2-D open problem in an inhomogeneous medium.

For all cases, we solve the resulting linear system with full GMRES and compare three preconditioners $M_0$, $M_1$, and $M_i$. We set the maximum number of GMRES iterations to 150. Storing 150 vectors, however, is too expensive, requiring a restart parameter. The storage issue is the main drawback of using GMRES. Therefore, for the third problem the GMRES convergence is compared to that of CGNR and Bi-CGSTAB. The iteration is terminated at the $k$th step if $\|b - Ap_k\|_2/\|b\|_2 < 10^{-6}$.

For the preconditioner solves, a direct method is used. In practice this process is very costly. Since the matrix $M$ is complex symmetric and both the real and imaginary parts are positive definite (or CSPD), the $LDL^T$ factorization can always be done (without requiring pivoting) and is unique [12]. This allows us, e.g., to approximate $M$ using ILU and then to use backward-forward substitution subsequently. As alternatives, we can also approximate $M^{-1}$ using a few steps of SSOR or multigrid [10]. In this paper, we do not implement these cheaper processes. Rather, we compute $M^{-1}$ exactly. We expect that having exact solution of the preconditioning step will provide us the detailed insights in the convergence, the lowest iteration numbers and therefore, it can be used as reference for approximation methods for $M^{-1}$.

## 6.1. Closed-off problem

We consider a problem in a rectangular homogeneous medium governed by

$$(\Delta + k^2)\phi = (k^2 - 5\pi^2) \sin(\pi x) \sin(2\pi y), \quad x = [0, 1], \ y = [0, 1],$$
$$\phi = 0, \quad \text{at the boundaries.} \tag{44}$$

The exact solution of (44) is $\phi = \sin(\pi x) \sin(2\pi y)$. Different grid resolutions are used to solve the problem with various wavenumbers $k = 2, 5, 10, 15, 20, 30, 40$. $k = 2$ resembles the definite problem. In Fig. 3, spectra of the preconditioned system for $k = 5$, a "slightly" indefinite problem, are shown. All spectra are bounded above by one.

Table 1 shows the computational performance in terms of number of iterations and computational time to reach the specified convergence. For low frequencies, all preconditioners show a very satisfactorily comparable performance. $M_0$ becomes less effective for increasing values of $k$, where the number of iterations increases somewhat faster than for $M_1$ or $M_i$. For large $k$, preconditioning with $M_i$ gives the fastest convergence. This behavior agrees with the theory. However, preconditioning with $M_i$ is expensive. As Problem 1 only requires real arithmetic operations, using $M_i$ destroys the cheap operations. Furthermore, the computational time shown in Table 1 is practically unacceptable due to the exact inverse of $M$. In real applications, some cheaper approximate methods for the preconditioner will be implemented.
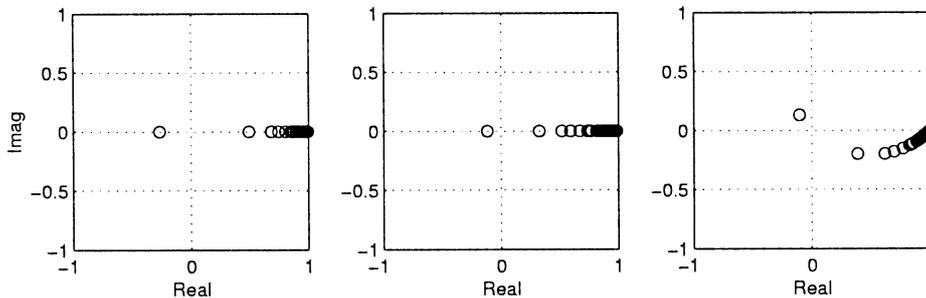


Fig. 3. Some extreme eigenvalues of the preconditioned systems of Problem 1 with $k = 5$ and gridsize $h^{-1} = 20$.

Table 1
Computational performance of GMRES for 2-D closed-off problem. The preconditioner is the shifted Laplace operator. 30 gridpoints per wavelength are used

| $k$ | $M_0$ | | $M_1$ | | $M_i$ | |
|---|---|---|---|---|---|---|
| | Iter | Time(s) | Iter | Time(s) | Iter | Time(s) |
| 2 | 5 | 0.02 | 5 | 0.02 | 5 | 0.05 |
| 5 | 8 | 0.24 | 10 | 0.31 | 9 | 0.68 |
| 10 | 13 | 1.75 | 16 | 2.16 | 15 | 8.45 |
| 15 | 18 | 7.01 | 22 | 8.44 | 20 | 39.26 |
| 20 | 26 | 21.49 | 29 | 24.32 | 26 | 194.86 |
| 30 | 57 | 170.04 | 60 | 178.28 | 49 | 1190.32 |
| 40 | 103 | 729.54 | 99 | 709.94 | 80 | 6623.48 |

## 6.2. 2-D open homogeneous problem

The second problem represents an open problem allowing waves to penetrate the boundaries. We first look at a homogeneous medium in which waves created at the upper surface propagate. We consider

$$
\begin{aligned}
\Delta\phi + k^2\phi &= f, & \Omega &= [0,1]^2, \\
f &= \delta\left(x - \frac{1}{2}\right)\delta(y), & x &= [0,1], \ y = 0, \\
\phi &= 0, & y &= 0, \\
\frac{\partial\phi}{\partial n} - \mathrm{i}k\phi &= 0, & x &= 0,1, \ y = 1,
\end{aligned}
\tag{45}
$$

with $k$ constant in $\Omega$. The performance of GMRES with preconditioners $M_0$, $M_1$, and $M_i$ is compared. In the construction of the preconditioning matrix, the same boundary conditions as in (45) are used.

Table 2 shows the number of GMRES iterations to solve Problem 2. For all frequencies, $M_i$ outperforms $M_0$ and $M_1$. $M_0$ still performs reasonably well compared to $M_i$. This is not explained by the theory and may be due to the influence of Sommerfeld boundary conditions imposed in constructing the preconditioning matrix, which is not taken into account in our analysis.

Fig. 4 shows the updated residual computed at each iteration for $k = 20$. The residual curve indicates slow convergence for the first few iterations and a convergence improvement later on, indicating a superlinear convergence. The slow convergence part is mainly due to the small eigenvalues. Once they are removed from the spectrum the convergence rate increases.
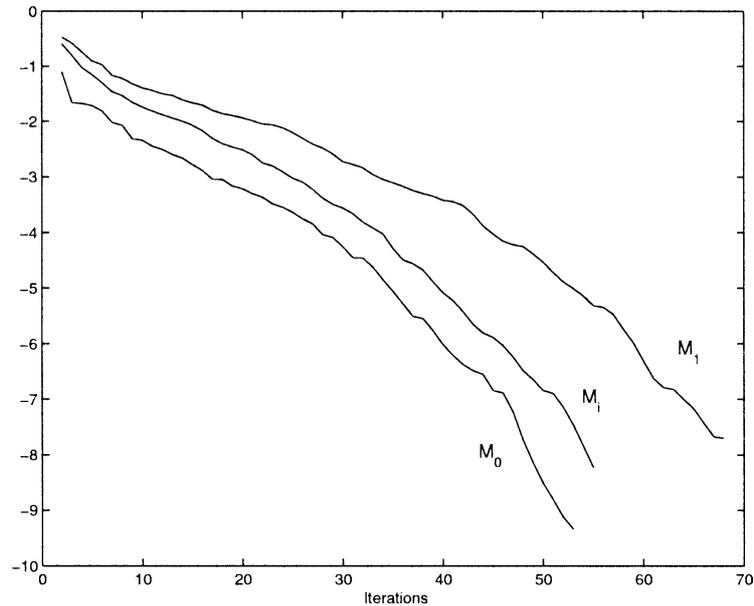


Fig. 4. Relative residual of preconditioned GMRES iterations, $k = 20$. $r_k = M^{-1}(b - Ap_k)$.

Table 2
Computational performance of GMRES to solve Problem 2. The preconditioner is the shifted Laplace preconditioners. 30 gridpoints per wavelength are used

| $k$ | $M_0$ | | $M_1$ | | $M_i$ | |
|---|---|---|---|---|---|---|
| | Iter | Time(s) | Iter | Time(s) | Iter | Time(s) |
| 2 | 6 | 0.05 | 7 | 0.06 | 6 | 0.07 |
| 5 | 10 | 0.80 | 13 | 0.96 | 11 | 0.85 |
| 10 | 20 | 11.72 | 25 | 15.78 | 22 | 14.70 |
| 15 | 33 | 64.87 | 41 | 92.31 | 37 | 84.43 |
| 20 | 52 | 358.03 | 67 | 449.42 | 54 | 401.06 |
| 30 | 102 | 3382.82 | 136 | 4059.72 | 97 | 2819.53 |

Table 3
Computational performance of GMRES, CGNR, and Bi-CGSTAB to solve Problem 3. The preconditioner is the shifted Laplace operator. 30 gridpoints per $k_{\mathrm{ref}}$ are used

| $k_{\mathrm{ref}}$ | GMRES | | | CGNR | | | Bi-CGSTAB | | |
|---|---|---|---|---|---|---|---|---|---|
| | $M_0$ | $M_1$ | $M_i$ | $M_0$ | $M_1$ | $M_i$ | $M_0$ | $M_1$ | $M_i$ |
| 2 | 7 | 9 | 8 | 10 | 13 | 11 | 5 | 6 | 5 |
| 5 | 14 | 19 | 16 | 22 | 31 | 26 | 10 | 13 | 10 |
| 10 | 34 | 42 | 36 | 77 | 83 | 65 | 50 | 53 | 26 |
| 15 | 64 | 82 | 63 | 210 | 160 | 122 | 206 | 115 | 36 |
| 20 | 107 | 136 | 91 | – | – | 185 | 448 | 159 | 50 |
| 30 | >150 | >150 | 140 | – | – | – | – | – | 70 |

### 6.3. 2-D open inhomogeneous problem

In this example we repeat the computation of Problem 2 but now in an in-homogeneous medium. The wavenumber varies inside the domain according to

$$k = \begin{cases} k_{\mathrm{ref}} & 0 \leqslant y \leqslant 1/3, \\ 1.5k_{\mathrm{ref}} & 1/3 \leqslant y \leqslant 2/3, \\ 2.0k_{\mathrm{ref}} & 2/3 \leqslant y \leqslant 1.0. \end{cases} \tag{46}$$

The number of gridpoints used is $5 \times k_{\mathrm{ref}}$ (i.e., approximately 30 gridpoints per reference wavelength) in the $x$ and $y$ directions. Numerical results are presented in Table 3. Here, we compute the solutions using full GMRES, and compare the computational performances with CGNR and Bi-CGSTAB.

In this harder problem, $M_i$ again outperforms $M_0$ and $M_1$ indicated by the smaller number of iterations required to reach convergence. Compared to $M_0$, $M_1$ shows a less satisfactorily performance, and based on our computational restrictions restart is needed. For GMRES, restarting is needed for $k > 20$.

From Table 3, we also see that the preconditioned Bi-CGSTAB does not perform well for $M_0$ and $M_1$, as already indicated in [14]. However, the convergence with $M_i$ as the preconditioner is still satisfactory. Compared to GMRES, Bi-CGSTAB preconditioned by $M_i$ shows better convergence performance (despite of requiring two preconditioning steps within one iteration). If $M_i$ is used as the preconditioner, Bi-CGSTAB can be the alternative to replace full GMRES.

From Table 3, one also concludes that CGNR may not be a good iterative method to solve the Helmholtz problem with the shifted Laplace preconditioners. This is mainly due to the squaring of the

original eigenvalues in the case of the Normal Equations, causing too many small eigenvalues. With such a spectrum, CG often exhibits very slow convergence. However, since our analysis for the preconditioners is based on the normal equations, the results of CGNR are included and confirm our analysis.

## 7. Conclusion

In this paper, a class of preconditioners based on the shifted Laplace operator for the Helmholtz equation has been presented and analyzed. We find that the complex shifted–Laplace operator leads to the most effective preconditioning matrix within this class of preconditioners. Numerical experiments have been presented to show the effectiveness of the preconditioner. This preconditioner is easy to construct and to extend to inhomogeneous medium cases. Our numerical experiments show that for the latter, this preconditioner performs effectively. With respect to storage and CPU time requirements, we advocate the complex shifted preconditioner in combination with Bi-CGSTAB.

## References

[1] A. Bamberger, P. Joly, J.E. Roberts, Second-order absorbing boundary conditions for the wave equation: A solution for corner problem, SIAM J. Numer. Anal. 27 (2) (1990) 323–352.
[2] A. Bayliss, C.I. Goldstein, E. Turkel, An iterative method for Helmholtz equation, J. Comput. Phys. 49 (1983) 443–457.
[3] J.-P. Berenger, A perfectly matched layer for the absorption of electromagnetic waves, J. Comput. Phys. 114 (1994) 185–200.
[4] R.W. Clayton, B. Engquist, Absorbing boundary conditions for wave-equation migration, Geophysics 45 (5) (1980) 895–904.
[5] B. Engquist, A. Majda, Absorbing boundary conditions for the numerical simulation of waves, Math. Comp. 31 (1977) 629–651.
[6] R. Fletcher, Conjugate gradient methods for indefinite systems, in: G.A. Watson (Ed.), Proc. of the Dundee Biennal Conference on Numerical Analysis, 1974, Springer, New York, 1975, pp. 73–89.
[7] R.W. Ereund, Preconditioning of symmetric, but highly indefinite linear systems, Numer. Anal. Manus. 97-3-03, Bell Labs., Murray Hill, NJ, 1997.
[8] R.W. Freund, N.M. Nachtigal, QMR: A quasi minimum residual method for non-Hermitian linear systems, Numer. Math. 60 (1991) 315–339.
[9] M.J. Gander, F. Nataf, AILU for Helmholtz problems: A new preconditioner based on the analytic parabolic factorization, J. Comput. Acoustics 9 (4) (2001) 1499–1509.
[10] D. Lahaye, H. De Gersem, S. Vandewalle, K. Hameyer, Algebraic multigrid for complex symmetric systems, IEEE Trans. Magn. 36 (4) (2000) 1535–1538.
[11] J. Gozani, A. Nachshon, E. Turkel, Conjugate gradient coupled with multi-grid for an indefinite problem, in: R. Vichnevetsky, R.S. Tepelman (Eds.), in: Advances in Computer Methods for Partial Differential Equations, vol. V, IMACS, New Brunswick, NJ, 1984, pp. 425–427.
[12] N.J. Higham, Factorizing complex symmetric matrices with positive definite real and imaginary parts, Math. Comp. 67 (3) (1998) 1591–1599.
[13] R. Kechroud, A. Soulaimani, Y. Saad, Preconditioning techniques for the solution of the Helmholtz equation by the finite element method, in: Workshop in Wave Phenomena in Physics and Engineering: New Models, Algorithms and Applications, May 18–21, 2003, Springer, Berlin, 2003.
[14] A.L. Laird, Preconditioned iterative solution of the 2D Helmholtz equation, First Year's Report, St. Hugh's College, Oxford, 2001.
[15] M.M.M. Made, Incomplete factorization-based preconditionings for solving the Helmholtz equation, Internat. J. Numer. Methods Engng. 50 (2001) 1077–1101.

[16] C.C. Paige, M.A. Saunders, Solution of sparse indefinite systems of linear equations, SIAM J. Numer. Anal. 12 (4) (1975) 617–629.

[17] R.-E. Plessix, W.A. Mulder, Separation of variables as a preconditioner for an iterative Helmholtz solver, Appl. Numer. Math. 44 (2003) 385–400.

[18] Y. Saad, A flexible inner-outer preconditioned GMRES algorithm, SIAM J. Sci. Statist. Comput. 14 (1993) 461–469.

[19] Y. Saad, Iterative Methods for Sparse Linear Systems, PWS, Boston, 1996.

[20] Y. Saad, M.H. Schultz, GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems, SIAM J. Sci. Statist. Comput. 7 (12) (1986) 856–869.

[21] P. Sonneveld, CGS, a fast Lanczos-type solver for nonsymmetric linear systems, SIAM J. Sci. Statist. Comput. 10 (1) (1986) 36–52.

[22] H.A. van der Vorst, Bi-CGSTAB: A fast and smoothly converging variant of BI-CG for the solution of nonsymmetric linear systems, SIAM J. Sci. Statist. Comput. 13 (2) (1992) 631–644.

[23] H.A. van der Vorst, J.B.M. Melissen, A Petrov–Galerkin type method for solving $Ax = b$, where $A$ is symmetric complex, IEEE Trans. Magnetics 26 (2) (1990) 706–708.

[24] H.A. van der Vorst, C. Vuik, GMRESR: A family of nested GMRES methods, Numer. Linear Algebra Appl. 1 (4) (1994) 369–386.