

**An optimal block iterative method and  
preconditioner for banded matrices  
with applications to PDEs on irregular domains**

Martin J. Gander, Sébastien Loisel, and Daniel B. Szyld

Report 10-05-21  
May 2010

This report is available in the World Wide Web at  
<http://www.math.temple.edu/~szyld>



# AN OPTIMAL BLOCK ITERATIVE METHOD AND PRECONDITIONER FOR BANDED MATRICES WITH APPLICATIONS TO PDES ON IRREGULAR DOMAINS\*

MARTIN J. GANDER<sup>†</sup>, SÉBASTIEN LOISEL<sup>‡</sup>, AND DANIEL B. SZYLD<sup>‡</sup>

**Abstract.** Classical Schwarz methods and preconditioners subdivide the domain of a partial differential equation into subdomains and use Dirichlet transmission conditions at the artificial interfaces. Optimized Schwarz methods use Robin (or higher order) transmission conditions instead, and the Robin parameter can be optimized so that the resulting iterative method has an optimized convergence factor. The usual technique used to find the optimal parameter is Fourier analysis; but this is only applicable to certain regular domains, for example, a rectangle, and with constant coefficients. In this paper, we present a completely algebraic version of the optimized Schwarz method, including an algebraic approach to find the optimal operator or a sparse approximation thereof. This approach allows us to apply this method to any banded or block banded linear system of equations, and in particular to discretizations of partial differential equations in two and three dimensions on irregular domains. With the computable optimal operator, we prove that the optimized Schwarz method converges in no more than two iterations for the case of two subdomains. Similarly, we prove that when we use an optimized Schwarz preconditioner with this optimal parameter, the underlying minimal residual Krylov subspace method (e.g., GMRES) converges in no more than two iterations. Very fast convergence is attained even when the optimal transmission operator is approximated by a sparse matrix. Numerical examples illustrating these results are presented.

**AMS subject classifications.** 65F08, 65F10, 65N22, 65N55

**Key words.** Linear systems, banded matrices, block matrices, Schwarz methods, optimized Schwarz methods, iterative methods, preconditioners.

**1. Introduction.** Finite difference or finite element discretizations of partial differential equations usually produce matrices which are banded, or block banded (e.g., block tridiagonal, or block pentadiagonal). In this paper, we present a novel iterative method for such block and banded matrices, guaranteed to converge in at most two steps. Similarly, its use as a preconditioner for minimal residual methods also achieves convergence in two steps. The formulation of this method proceeds by appropriately replacing a small block of the matrix in the iteration operator. As we will show, approximations of this replacement also produce very fast convergence. The method is based on an algebraic rendition of optimized Schwarz methods.

Schwarz methods are important tools for the numerical solution of partial differential equations. They are based on a decomposition of the domain into subdomains, and on the (approximate) solution of the (local) problems in each subdomain. In the classical formulation, Dirichlet boundary conditions at the artificial interfaces are used; see, e.g., [21], [24], [27], [30]. In optimized Schwarz methods, Robin and higher order boundary conditions are used in the artificial interfaces, e.g., of the form  $\partial_n u(x) + \alpha u(x)$ . By optimizing the parameter  $\alpha$ , one can obtain optimized convergence of the Schwarz methods; see, e.g., [4], [5], [6], [10], [11], [12], [16]. The tools usually employed for the study of optimized Schwarz methods and its parameter estimation are based on Fourier analysis. This limits the applicability of the technique to certain classes of differential equations, and simple domains, e.g., rectangles or spheres.

Algebraic analyses of classical Schwarz methods were shown to be useful in their understanding and extensions; see, e.g., [2], [8], [14], [23], [29]. In particular, it follows that the classical additive and multiplicative Schwarz iterative methods and preconditioners can be regarded as the classical block Jacobi or block Gauss-Seidel methods, respectively, with the addition of overlap; see section 2. In our approach, we consider the restricted version of the Schwarz methods, briefly described in section 3.

Inspired in part by the earlier work on algebraic Schwarz methods, in this paper, we mimic the philosophy of optimized Schwarz methods when solving block banded linear systems; see also [18], [19], [20]. Our

---

\*This version dated May 21, 2010

<sup>†</sup>Section de Mathématiques, Université de Genève, CP 64, CH-1211 Geneva, Switzerland ([gander@math.unige.ch](mailto:gander@math.unige.ch)).

<sup>‡</sup>Department of Mathematics, Temple University (038-16), 1805 N. Broad Street, Philadelphia, Pennsylvania 19122-6094, USA ([loisel,szylid@temple.edu](mailto:{loisel,szylid}@temple.edu)).

approach consists of optimizing the block which would correspond to the artificial interface (called transmission matrix), so that the spectral radius of the iteration operator is reduced; see section 4. With the optimal transmission operator, we show that the new method is guaranteed to converge in no more than two steps. In section 5 we use the same approach for multiplicative Schwarz methods. We show that the same minimization as in the additive case is what is needed.

When we use our optimal approach to precondition a minimal residual Krylov subspace method, such as GMRES, the preconditioned iterations are also guaranteed to converge in no more than two steps; see section 6.

We can approximate the computation of the optimal transmission matrix in at least two general ways. We can approximate some inverses appearing in the expression of the optimal matrix, e.g., by using an incomplete LU factorization. We can also restrict the minimization of the norm of some factors of the iteration operator to blocks with certain sparsity patterns, for example, scalar, diagonal, or tridiagonal matrices. Since the new method is applicable to any (block) banded matrix, in particular we can use it to solve systems arising from the discretization of PDEs on unstructured meshes, and/or on irregular domains, and in section 9 we show experiments for these cases as well. The numerical results both for the iterative methods and the preconditioner show again convergence in two steps in the optimal case, while very fast convergence is also achieved with simple minimizations with very few parameters.

For a model problem, we compare our algebraic results to those that can be obtained with Fourier analysis on the discretized differential equation; see section 10.

We end the paper with an experimental study showing that the ILU approximations can produce very good results (section 11), and present some concluding remarks.

**2. Classical methods with overlap.** Our aim is to solve a linear system of equations of the form  $Au = f$ , where the  $n \times n$  matrix  $A$  is either banded, or block-banded, or, more generally, it has the following form

$$A = \begin{bmatrix} A_{11} & A_{12} & A_{13} & & \\ A_{21} & A_{22} & A_{23} & & \\ & A_{32} & A_{33} & A_{34} & \\ & & A_{42} & A_{43} & A_{44} \end{bmatrix}. \quad (2.1)$$

In most practical cases, where  $A$  corresponds to a discretization of a differential equation, one has that  $A_{13} = A_{42} = O$ , i.e., they are zero blocks. Each block  $A_{ij}$  is of order  $n_i \times n_j$ ,  $i, j = 1, \dots, 4$ , and  $\sum_i n_i = n$ . We have in mind the situation where  $n_1 \gg n_2$  and  $n_4 \gg n_3$ , as illustrated, e.g., in Figure 2.1.

Consider first the following two diagonal blocks (without overlap)

$$A_1 = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad A_2 = \begin{bmatrix} A_{33} & A_{34} \\ A_{43} & A_{44} \end{bmatrix}, \quad (2.2)$$

which are square, but not necessarily of the same size; cf. the example in Figure 2.1. The Block Jacobi preconditioner (or block diagonal preconditioner is)

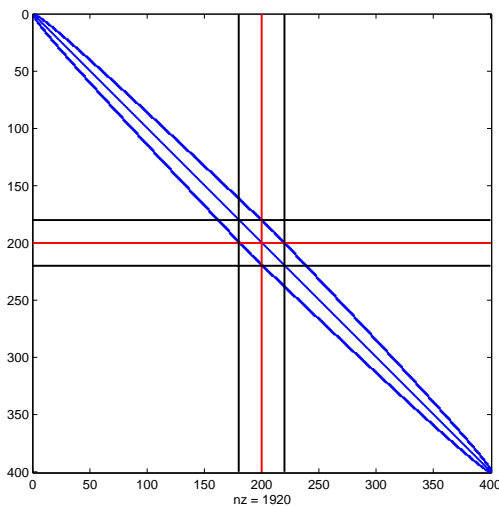
$$M^{-1} = M_{BJ}^{-1} = \begin{bmatrix} A_1^{-1} & O \\ O & A_2^{-1} \end{bmatrix} = \sum_{i=1}^2 R_i^T A_i^{-1} R_i, \quad (2.3)$$

where the restriction operators are

$$R_1 = [ I \quad O ] \text{ and } R_2 = [ O \quad I ],$$

which have order  $(n_1 + n_2) \times n$  and  $(n_3 + n_4) \times n$ , respectively. The transpose of these operators,  $R_i^T$  are prolongation operators. The standard block Jacobi method, using these two blocks has an iteration operator of the form

$$T = T_{BJ} = I - M_{BJ}^{-1}A = I - \sum R_i^T A_i^{-1} R_i A.$$

FIG. 2.1. A  $400 \times 400$  band matrix partitioned into  $4 \times 4$  blocks.

The iterative method is then, for a given initial vector  $u^0$ ,  $u^{k+1} = Tu^k + M^{-1}f$ ,  $k = 0, 1, \dots$ , and its convergence is linear with an asymptotic convergence factor  $\rho(T)$ , the spectral radius of the iteration operator; see, e.g., the classical reference [31].

Similarly, the Block Gauss-Seidel iterative method for a system with a coefficient matrix (2.1) is defined by an iteration matrix of the form

$$T = T_{GS} = (I - R_2^T A_2^{-1} R_2 A)(I - R_1^T A_1^{-1} R_1 A) = \prod_{i=2}^1 (I - R_i^T A_i^{-1} R_i A),$$

where the corresponding preconditioner can thus be written as

$$M_{GS}^{-1} = [I - (I - R_2^T A_2^{-1} R_2 A)(I - R_1^T A_1^{-1} R_1 A)]A^{-1}. \quad (2.4)$$

Consider now the same blocks (2.2), with *overlap*, and using the same notation we write the new blocks with overlap as

$$A_1 = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ & A_{32} & A_{33} \end{bmatrix}, \quad A_2 = \begin{bmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} & A_{34} \\ A_{42} & A_{43} & A_{44} \end{bmatrix}. \quad (2.5)$$

The corresponding restriction operators are again

$$R_1 = [ I \quad O ] \quad \text{and} \quad R_2 = [ O \quad I ], \quad (2.6)$$

which have now order  $(n_1 + n_2 + n_3) \times n$  and  $(n_2 + n_3 + n_4) \times n$ , respectively. With this notation, the additive and multiplicative Schwarz preconditioners (with or without overlap) are

$$M_{AS}^{-1} = \sum_{i=1}^2 R_i^T A_i^{-1} R_i \quad \text{and} \quad M_{MS}^{-1} = [I - (I - R_2^T A_2^{-1} R_2 A)(I - R_1^T A_1^{-1} R_1 A)]A^{-1}, \quad (2.7)$$

respectively; see, e.g., [27], [30]. By comparing (2.3) and (2.4) with (2.7) one concludes that the classical Schwarz preconditioners can be regarded as Block Jacobi or Block Gauss-Seidel with the addition of overlap.

**3. Restricted Schwarz methods.** From the preconditioners (2.7), one can write explicitly the iteration operators for the additive and multiplicative Schwarz iterations as

$$T_{AS} = I - \sum_{i=1}^2 R_i^T A_i^{-1} R_i A \quad (3.1)$$

$$\text{and } T_{MS} = \prod_{i=2}^1 (I - R_i^T A_i^{-1} R_i A),$$

respectively. The additive Schwarz iteration (with overlap) associated with the iteration operator in (3.1) is usually not convergent; this is because it holds that with overlap  $\sum R_i^T R_i > I$ . The standard approach is to use a damping parameter  $0 < \gamma < 1$  so that the iteration operator  $T_R(\gamma) = I - \gamma \sum_{i=1}^2 R_i^T A_i^{-1} R_i A$  is such that  $\rho(T_R(\gamma)) < 1$ ; see, e.g., [27], [30]. We will not pursue this strategy here. Instead we consider the Restricted Additive Schwarz (RAS) iterations [3], [9].

The RAS method consists of using the local solvers with the overlap (2.5), with the corresponding restriction operators  $R_i$ , but use the prolongations  $\tilde{R}_i^T$  without the overlap, which are defined as

$$\tilde{R}_1 = \begin{bmatrix} I & O \\ O & O \end{bmatrix} \quad \text{and} \quad \tilde{R}_2 = \begin{bmatrix} O & O \\ O & I \end{bmatrix}, \quad (3.2)$$

having the same order as the matrices  $R_i$  in (2.6), and where the identity in  $\tilde{R}_1$  is of order  $n_1 + n_2$  and that in  $\tilde{R}_2$  of order  $n_3 + n_4$ . These restriction operators select the variables without the overlap. Note that we have now  $\sum \tilde{R}_i^T R_i = I$ . In this way, there is no ‘‘double counting’’ of the variables on the overlap, and, under certain hypothesis, there is no need to use a relaxation parameter to obtain convergence; see [9], [11] for details. Thus, the RAS iteration operator is

$$T_{RAS} = I - \sum \tilde{R}_i^T A_i^{-1} R_i A. \quad (3.3)$$

Similarly, one can have Restricted Multiplicative Schwarz (RMS) [3], [22], and the iteration operator is

$$T = T_{RMS} = \prod_{i=2}^1 (I - \tilde{R}_i^T A_i^{-1} R_i A) = (I - \tilde{R}_2^T A_2^{-1} R_2 A)(I - \tilde{R}_1^T A_1^{-1} R_1 A), \quad (3.4)$$

although in this case the  $\tilde{R}_i^T$  are not necessary to avoid double counting. We include this method for completeness.

**4. Replacing the transmission matrices. Additive case.** Our proposed new method consists of replacing the transmission matrices  $A_{33}$  (lowest right corner) in  $A_1$  and  $A_{22}$  (upper left corner) in  $A_2$  so that the modified operators of the form (3.3) and (3.4) have small spectral radii, and thus, the corresponding iterative methods have fast convergence. Let the replaced blocks in  $A_1$  and in  $A_2$  be

$$S_1 = A_{33} + D_1, \quad \text{and} \quad S_2 = A_{22} + D_2, \quad (4.1)$$

respectively, and let us call the modified matrices  $\tilde{A}_i$ , i.e., we have

$$\tilde{A}_1 = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ & A_{32} & S_1 \end{bmatrix}, \quad \tilde{A}_2 = \begin{bmatrix} S_2 & A_{23} \\ A_{32} & A_{33} & A_{34} \\ A_{42} & A_{43} & A_{44} \end{bmatrix}; \quad (4.2)$$

cf. (2.5). With this notation, our proposed modified RAS iteration operator is thus

$$T_{MRAS} = I - \sum \tilde{R}_i^T \tilde{A}_i^{-1} R_i A, \quad (4.3)$$

and we want to study modifications  $D_i$  so that  $\|T_{MRAS}\| \ll 1$  for some suitable norm. This would imply of course that  $\rho(T_{MRAS}) \ll 1$ . Finding the appropriate modifications  $D_i$  is analogous to finding the appropriate parameter  $\alpha$  in optimized Schwarz methods; see our discussion in section 1 and references therein.

To that end, we first introduce some notation. Let  $E_3 = E_{n_3}$  be the  $(n_1 + n_2 + n_3) \times n_3$  matrix given by  $E_3^T = [ O \ O \ I ]$ , and let  $E_1 = E_{n_2}$  be the  $(n_2 + n_3 + n_4) \times n_2$  matrix given by  $E_1^T = [ I \ O \ O ]$ . Let

$$A_1^{-1}E_3 =: B_3^{(1)} = \begin{bmatrix} B_{31} \\ B_{32} \\ B_{33} \end{bmatrix}, \quad A_2^{-1}E_1 =: B_1^{(2)} = \begin{bmatrix} B_{11} \\ B_{12} \\ B_{13} \end{bmatrix}, \quad (4.4)$$

i.e., the last block column of  $A_1^{-1}$ , and the first block column of  $A_2^{-1}$ , respectively. Furthermore, we denote

$$\tilde{B}_3^{(1)} = \begin{bmatrix} B_{31} \\ B_{32} \end{bmatrix}, \quad \tilde{B}_1^{(2)} = \begin{bmatrix} B_{12} \\ B_{13} \end{bmatrix}, \quad (4.5)$$

i.e., pick the first two blocks of  $B_3^{(1)}$  of order  $(n_1 + n_2) \times n_3$ , and the last two blocks of  $B_1^{(2)}$  of order  $(n_3 + n_4) \times n_2$ . Finally, let

$$\bar{E}_1^T = [ I \ O ] \text{ and } \bar{E}_2^T = [ O \ I ], \quad (4.6)$$

which have order  $(n_1 + n_2) \times n$  and  $(n_3 + n_4) \times n$ , respectively.

LEMMA 4.1. *The new iteration matrix (4.3) has the form*

$$T = T_{MRAS} = \left[ \begin{array}{c|c} O & K \\ \hline L & O \end{array} \right], \quad (4.7)$$

where

$$K = \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1} [\bar{E}_1^T - D_1^{-1}A_{34}\bar{E}_2^T], \quad L = \tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1} [-D_2^{-1}A_{21}\bar{E}_1^T + \bar{E}_2^T]. \quad (4.8)$$

*Proof.* We can write

$$\tilde{A}_1 = A_1 + E_3D_1E_3^T, \quad \tilde{A}_2 = A_2 + E_1D_2E_1^T.$$

Using the Sherman-Morrison-Woodbury formula (see, e.g., [15]) we can explicitly write  $\tilde{A}_i^{-1}$  in terms of  $A_i^{-1}$  as follows

$$\tilde{A}_1^{-1} = A_1^{-1} - A_1^{-1}E_3(D_1^{-1} + E_3^T A_1^{-1}E_3)^{-1}E_3^T A_1^{-1} =: A_1^{-1} - C_1, \quad (4.9)$$

$$\tilde{A}_2^{-1} = A_2^{-1} - A_2^{-1}E_1(D_2^{-1} + E_1^T A_2^{-1}E_1)^{-1}E_1^T A_2^{-1} =: A_2^{-1} - C_2, \quad (4.10)$$

and observe that  $E_3^T A_1^{-1}E_3 = B_{33}$ ,  $E_1^T A_2^{-1}E_1 = B_{11}$ .

Let us first consider the term with  $i = 1$  in (4.3). We begin by noting that from (2.1) it follows that  $R_1A = [ A_1 \mid E_3A_{34} ]$ . Thus,

$$\tilde{A}_1^{-1}R_1A = (A_1^{-1} - C_1) [ A_1 \mid E_3A_{34} ] = [ I - C_1A_1 \mid A_1^{-1}E_3A_{34} - C_1E_3A_{34} ]. \quad (4.11)$$

We look now at each part of (4.11). First from (4.9), we have that  $C_1A_1 = B_3^{(1)}(D_1^{-1} + B_{33})^{-1}E_3^T$ . Then we see that  $C_1E_3A_{34} = B_3^{(1)}(D_1^{-1} + B_{33})^{-1}E_3^T A_1^{-1}E_3A_{34} = B_3^{(1)}(D_1^{-1} + B_{33})^{-1}B_{33}A_{34}$ , and therefore

$$\begin{aligned} A_1^{-1}E_3A_{34} - C_1E_3A_{34} &= B_3^{(1)}[I - (D_1^{-1} + B_{33})^{-1}B_{33}]A_{34} = \\ &= B_3^{(1)}[(D_1^{-1} + B_{33})^{-1}(D_1^{-1} + B_{33} - B_{33})^{-1}D_1^{-1}A_{34}] = B_3^{(1)}(D_1^{-1} + B_{33})^{-1}D_1^{-1}A_{34}. \end{aligned}$$

Putting this together we have

$$\tilde{A}_1^{-1}R_1A = \left[ I - B_3^{(1)}(D_1^{-1} + B_{33})^{-1}E_3^T \mid B_3^{(1)}(D_1^{-1} + B_{33})^{-1}D_1^{-1}A_{34} \right]. \quad (4.12)$$

It is important to note that the lower blocks in this expression, corresponding to the overlap, will not be considered once it is multiplied by  $\tilde{R}_1^T$ . An analogous calculation produces

$$\tilde{A}_2^{-1}R_2A = \left[ B_1^{(2)}(D_2^{-1} + B_{11})^{-1}D_2^{-1}A_{21} \mid I - B_1^{(2)}(D_2^{-1} + B_{11})^{-1}E_1^T \right], \quad (4.13)$$

and again, one should note that the upper blocks would be eliminated with the multiplication by  $\tilde{R}_2^T$ . We remark that the number of columns of the blocks in (4.12) and (4.13) are not the same. Indeed, the first block in (4.12) is of order  $(n_1 + n_2) \times (n_1 + n_2 + n_3)$ , and the second is of order  $(n_1 + n_2) \times n_4$ , while the first block in (4.13) is of order  $(n_3 + n_4) \times n_1$  and the second is of order  $(n_3 + n_4) \times (n_2 + n_3 + n_4)$ .

We apply the restrictions (3.2) to (4.12) and (4.13), and collect terms to form (4.3). First notice that the identity matrix in (4.3) and the identity matrices in (4.12) and (4.13) cancel each other. We thus have

$$\begin{aligned} T &= \begin{bmatrix} \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1}E_3^T & \left| -\tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1}D_1^{-1}A_{34} \right. \\ -\tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1}D_2^{-1}A_{21} & \left| \tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1}E_1^T \right. \end{bmatrix} \\ &= \begin{bmatrix} \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1} \left[ E_3^T & \left| -D_1^{-1}A_{34} \right. \right] \\ \tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1} \left[ -D_2^{-1}A_{21} & \left| E_1^T \right. \right] \end{bmatrix} = \begin{bmatrix} \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1} \left[ \tilde{E}_3^T - D_1^{-1}A_{34}\tilde{E}_4^T \right] \\ \tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1} \left[ -D_2^{-1}A_{21}\tilde{E}_1^T + \tilde{E}_2^T \right] \end{bmatrix}, \end{aligned} \quad (4.14)$$

where the last equality follows from enlarging  $E_3 = [O \ O \ I]^T$  to  $\tilde{E}_3 = [O \ O \ I \ O]^T$  and  $E_1 = [I \ O \ O]^T$  to  $\tilde{E}_1 = [O \ I \ O \ O]^T$ , and introducing  $\tilde{E}_4 = [O \ O \ O \ I]^T$  and  $\tilde{E}_2 = [I \ O \ O \ O]^T$ . A careful look at the form of the matrix (4.14) reveals the block structure (4.7) with (4.8).  $\square$

Recall that our goal is to find the appropriate matrices  $D_1, D_2$  in (4.1) to obtain a small  $\rho(T_{MRAS})$ . Given the form (4.7) we obtained, it would suffice to minimize  $\|K\|$  and  $\|L\|$ . As it turns out, even in simple cases, the best possible choices of the matrices  $D_1, D_2$ , produce a matrix  $T = T_{MRAS}$  with  $\|T\| > 1$  (although  $\rho(T) < 1$ ); see for example the case reported below in Figure 9.3. Thus, another strategy is needed. We proceed by considering  $T^2$ , which can easily be computed from (4.7) to obtain

$$T^2 = \begin{bmatrix} KL & \left| O \right. \\ O & \left| LK \right. \end{bmatrix}. \quad (4.15)$$

**THEOREM 4.2.** *The asymptotic convergence factor of the modified RAS method given by (4.3) is bounded by the product of the following two norms*

$$\|(I + D_1B_{33})^{-1} [D_1B_{12} - A_{34}B_{13}]\|, \quad \|(I + D_2B_{11})^{-1} [D_2B_{32} - A_{21}B_{31}]\|. \quad (4.16)$$

*Proof.* We consider  $T^2$  as in (4.15). Using (4.8), (4.6), and (4.5), we can write

$$\begin{aligned} KL &= \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1} [\tilde{E}_1^T - D_1^{-1}A_{34}\tilde{E}_2^T] \tilde{B}_1^{(2)}(D_2^{-1} + B_{11})^{-1} [-D_2^{-1}A_{21}\tilde{E}_1^T + \tilde{E}_2^T] \\ &= \tilde{B}_3^{(1)}(D_1^{-1} + B_{33})^{-1} [B_{12} - D_1^{-1}A_{34}B_{13}] (D_2^{-1} + B_{11})^{-1} [-D_2^{-1}A_{21}\tilde{E}_1^T + \tilde{E}_2^T] \\ &= \tilde{B}_3^{(1)}(I + D_1B_{33})^{-1} [D_1B_{12} - A_{34}B_{13}] (D_2^{-1} + B_{11})^{-1} [-D_2^{-1}A_{21}\tilde{E}_1^T + \tilde{E}_2^T], \end{aligned} \quad (4.17)$$

and similarly

$$LK = \tilde{B}_1^{(2)}(I + D_2B_{11})^{-1} [D_2B_{32} - A_{21}B_{31}] (D_1^{-1} + B_{33})^{-1} [\tilde{E}_1^T - D_1^{-1}A_{34}\tilde{E}_2^T].$$

Furthermore, let us consider the following products, which are present in  $KLKL$  and in  $LKLK$ ,

$$KL\tilde{B}_3^{(1)} = \tilde{B}_3^{(1)}(I + D_1B_{33})^{-1} [D_1B_{12} - A_{34}B_{13}] (I + D_2B_{11})^{-1} [D_2B_{32} - A_{21}B_{31}], \quad (4.18)$$

$$LK\tilde{B}_1^{(2)} = \tilde{B}_1^{(2)}(I + D_2B_{11})^{-1} [D_2B_{32} - A_{21}B_{31}] (I + D_1B_{33})^{-1} [D_1B_{12} - A_{34}B_{13}]. \quad (4.19)$$

These factors are present when considering the powers  $T^{2k}$ , and therefore, asymptotically their norm provides the convergence factor in which  $T^2$  goes to zero. Thus, the asymptotic convergence factor is bounded by the product of the two norms (4.16).  $\square$

As it can be appreciated, we have the same two factors in (4.18)-(4.19), and in order to find a small  $\|T^2\|$ , we could look for matrices  $D_1$  and  $D_2$  in certain sets of matrices, say  $\mathcal{Q}_1$ , and  $\mathcal{Q}_2$ , to minimize the norm of these factors, i.e., solving the nonlinear problems

$$\min_{D_1 \in \mathcal{Q}_1} \|(I + D_1B_{33})^{-1} [D_1B_{12} - A_{34}B_{13}]\|, \quad \min_{D_2 \in \mathcal{Q}_2} \|(I + D_2B_{11})^{-1} [D_2B_{32} - A_{21}B_{31}]\|. \quad (4.20)$$



Our approach is to consider instead, as a first approximation, matrices obtained by minimizing the following linear problems

$$\min_{D_1 \in \mathcal{Q}_1} \|D_1 B_{12} - A_{34} B_{13}\|, \quad \min_{D_2 \in \mathcal{Q}_2} \|D_2 B_{32} - A_{21} B_{31}\|. \quad (4.21)$$

If  $n_3 = n_2$  and either  $B_{12}$  or  $B_{32}$  are nonsingular (assumptions which hold in most practical problems), then we can solve, e.g., for  $S_1 = A_{33} - D_1$  in the expression

$$S_1 B_{12} - (A_{33} B_{12} + A_{34} B_{13}) = O, \quad (4.22)$$

cf. (4.21), i.e., we have

$$S_1 = (A_{33} B_{12} + A_{34} B_{13}) B_{12}^{-1} = A_{33} + A_{34} B_{13} B_{12}^{-1} \quad (4.23)$$

thus obtaining  $T^2 = O$ , and consequently, the proposed iterative method MRAS would converge in no more than two iterations. We call the matrices  $D_1$  and  $D_2$  thus obtained the optimal transmission matrices.

**5. Replacing the transmission matrices. Multiplicative case.** In this section we study the idea of using the modified matrices (4.2) for the restricted multiplicative Schwarz iterations, obtained by modifying the iteration operator (3.4), i.e., we have

$$T = T_{MRMS} = \prod_{i=2}^1 (I - \tilde{R}_i^T \tilde{A}_i^{-1} R_i A) = (I - \tilde{R}_2^T \tilde{A}_2^{-1} R_2 A) (I - \tilde{R}_1^T \tilde{A}_1^{-1} R_1 A) \quad (5.1)$$

and its associated preconditioner.

From (4.12), (4.13), and (4.8), we see that

$$(I - \tilde{R}_1^T \tilde{A}_1^{-1} R_1 A) = \left[ \begin{array}{c|c} O & K \\ \hline O & I \end{array} \right], \quad (I - \tilde{R}_2^T \tilde{A}_2^{-1} R_2 A) = \left[ \begin{array}{c|c} I & O \\ \hline L & O \end{array} \right].$$

As a consequence, putting together (5.1), we have the following structure

$$T = T_{MRMS} = \left[ \begin{array}{c|c} O & K \\ \hline O & LK \end{array} \right],$$

and from it, we can obtain the following result on its eigenvalues.

PROPOSITION 5.1. *Let*

$$T_{MRAS} = \left[ \begin{array}{c|c} O & K \\ \hline L & O \end{array} \right], \quad T_{TRMS} = \left[ \begin{array}{c|c} O & K \\ \hline O & LK \end{array} \right].$$

If  $\lambda \in \sigma(T_{MRAS})$ , then  $\lambda^2 \in \sigma(T_{MRMS})$ .

*Proof.* Let  $[x, v]^T$  be the eigenvector of  $T_{MRAS}$  corresponding to  $\lambda$ , i.e.,

$$\left[ \begin{array}{c|c} O & K \\ \hline L & O \end{array} \right] \begin{bmatrix} x \\ v \end{bmatrix} = \lambda \begin{bmatrix} x \\ v \end{bmatrix}.$$

Thus,  $Kv = \lambda x$ , and  $Lx = \lambda v$ . Then,  $LKv = \lambda Lx = \lambda^2 v$ , and the eigenvector for  $T_{MRMS}$  corresponding to  $\lambda^2$  is  $[0, v]^T$ .  $\square$

REMARK 5.2. *We note that the structure of (4.7) is the structure of a standard Block Jacobi iteration matrix for a “consistently ordered matrix” (see, e.g., [31], [32]), but our matrix is not of a Block Jacobi iteration. We note then that a matrix of this form has the property that if  $\mu \in \sigma(T)$ , then  $-\mu \in \sigma(T)$ ; see, e.g., [25, p. 120, Prop. 4.12]. This is consistent with our calculations of the spectra of the iteration matrices.*

Note that for consistently ordered matrices  $\rho(T_{GS}) = \rho(T_J)^2$ ; see, e.g., [31, Corollary 4.26]. Our generic block matrix  $A$  is not consistently ordered, but in Proposition 5.1 we proved a similar result.

Observe that in Proposition 5.1 we only provide half of the eigenvalues of  $T_{MRMS}$ ; the other eigenvalues are zero. Thus we have that  $\rho(T_{MRMS}) = \rho(T_{MRAS})^2$ , indicating a much faster asymptotic convergence of the multiplicative version.

Note also that for the “optimal” case, i.e., when we produce  $KL = O$ , we still need 2 iterations for convergence of the optimized restricted multiplicative Schwarz method, since in that case,

$$T_{MRMS} = \left[ \begin{array}{c|c} O & K \\ \hline O & O \end{array} \right] \neq O;$$

but we have that  $T_{MRMS}^2 = O$ .<sup>1</sup>

In general, we have that

$$T_{MRMS}^2 = \left[ \begin{array}{c|c} O & O \\ \hline (LK)^2 & O \end{array} \right]$$

and this implies that the *same* joint minimization (4.21) applies here as well, but asymptotically one expects to have half the number of iterations as for MRAS.

**6. Optimal preconditioner.** We consider in this section the solution of the block banded linear systems using a preconditioned minimal residual method such as GMRES or MINRES; see, e.g., [25], [26]. Based on the iteration matrix (4.3) our proposed Modified Restricted Additive Schwarz preconditioner is

$$M^{-1} = M_{MRAS}^{-1} = \sum \tilde{R}_i^T \tilde{A}_i^{-1} R_i. \quad (6.1)$$

Similarly, we can have a Modified Restricted Multiplicative Schwarz preconditioner corresponding to the iteration matrix (5.1). We show next that if we use the optimal transmission matrix, then the preconditioned problems can be solved in at most two iterations.

**PROPOSITION 6.1.** *Consider a linear system with coefficient matrix of the form (2.1), and a minimal residual method for its solution with either the MRAS preconditioner (6.1), or with the MRMS preconditioner, with  $\tilde{A}_i$  of the form (4.2) and  $S_i$  ( $i = 1, 2$ ) solving the appropriate equation such as (4.22). Then, the preconditioned minimal residual method converges in at most two iterations.*

*Proof.* We can write  $T^2 = (I - M^{-1}A)^2 = p_2(M^{-1}A) = 0$ , where  $p_2(z) = (1 - z)^2$  is a particular polynomial of degree 2 with  $p_2(0) = 1$ . Thus, the minimal residual polynomial  $q_2(z)$  of degree 2 also satisfies  $q_2(z) = 0$ .  $\square$

**7. Multiple diagonal blocks.** Our analysis so far has been restricted to the case of two (overlapping) blocks. We show in this section that our analysis applies to multiple (overlapping) blocks. To that end, we use a standard trick of Schwarz methods to handle the case of multiple subdomains, if they can be colored with two colors.

Let  $Q_1, \dots, Q_p$  be restriction matrices, defined by taking rows of the  $n \times n$  identity matrix  $I$ ; cf. (2.6). Let  $\tilde{Q}_1^T, \dots, \tilde{Q}_p^T$  be the corresponding prolongation operators, such that

$$I = \sum_{k=1}^p \tilde{Q}_k^T Q_k;$$

cf. (3.2). Given the stiffness matrix  $A$ , we say that the domain decomposition is two-colored if

$$Q_i^T A Q_j = O \text{ for all } |i - j| > 1.$$

<sup>1</sup>In fact, it suffices to apply only one factor of (5.1) in the second iteration, i.e., one and a half iterations suffice for convergence.

In this situation, if  $p$  is even, we can define

$$R_1 = \begin{bmatrix} Q_1 \\ Q_3 \\ \vdots \\ Q_{p-1} \end{bmatrix} \quad \text{and} \quad R_2 = \begin{bmatrix} Q_2 \\ Q_4 \\ \vdots \\ Q_p \end{bmatrix}. \quad (7.1)$$

We make similar definitions if  $p$  is odd, and also assemble  $\tilde{R}_1$  and  $\tilde{R}_2$  in a similar fashion.

The rows and columns of the matrices  $R_1$  and  $R_2$  could be permuted in such a way that (2.6) holds. Therefore, all the arguments of the previous sections as in the rest of the paper hold, *mutatis mutandis*. It is computationally more convenient to work with the matrices  $R_1$  and  $R_2$  as defined by (7.1). We now outline the computations needed to obtain the optimal transmission matrices (or their approximations).

The matrices  $A_1$  and  $A_2$ , defined by  $A_i = R_i A R_i^T$ , similar to (2.5), are block diagonal. The matrices  $E_1$  and  $E_3$  are defined by

$$E_i = \tilde{R}_i R_{3-i}^T,$$

and the matrices  $B_3^{(1)}$  and  $B_1^{(2)}$  are defined by

$$B_3^{(1)} = A_1^{-1} E_3 \quad \text{and} \quad B_1^{(2)} = A_2^{-1} E_1;$$

cf. (4.4). Since the matrices  $A_1$  and  $A_2$  are block diagonal, we have retained the parallelism of the  $p$  subdomains.

We must now define the finer structures, such as  $B_{12}$  and  $A_{34}$ . We say that the  $k$ th row is in the kernel of  $X$  if  $X e_k = 0$ , where  $e_k = [0, \dots, 0, 1, 0, \dots, 0]^T$  is the usual basis vector. Likewise, we say that the  $k$ th column is in the kernel of  $X$  if  $e_k^T X = 0$ . We define the matrix  $B_{12}$  to be the rows of  $B_1^{(2)}$  that are not in the kernel of  $R_1 \tilde{R}_2^T$ . We define the matrix  $B_{13}$  to be the rows of  $B_1^{(2)}$  that are in the kernel of  $R_1 \tilde{R}_2^T$ , and we make similar definitions for  $B_{32}$  and  $B_{31}$ ; cf. (4.4). The matrix  $A_{34}$  is the submatrix of  $A_2$  whose rows are not in the kernel of  $R_1 \tilde{R}_2^T$ , and whose columns are in the kernel of  $R_1 \tilde{R}_2^T$ . We make similar considerations for the other blocks  $A_{ij}$ .

This derivation allows us to define transmission matrices defined by (4.21) in the case of multiple diagonal overlapping blocks.

**8. Approximations of the transmission matrices.** Let us look at a practical way of finding the optimal transmission matrices such as (4.23). We consider the practical case when  $A_{13} = A_{42} = O$ . Then, from the definition (4.4), we have that  $A_{43} B_{12} + A_{44} B_{13} = O$  or  $B_{13} = -A_{44}^{-1} A_{43} B_{12}$ . Thus,

$$S_1 = A_{33} - A_{34} A_{44}^{-1} A_{43}. \quad (8.1)$$

Therefore, we need  $n_3$  solves with the local matrix  $A_{44}$  (i.e., local but stripped of the overlap). In the case of multiple diagonal blocks this matrix is block diagonal, and each diagonal block can be solved separately.

One alternative is to approximate the inverse  $A_{44}^{-1}$  (and similarly  $A_{11}^{-1}$ ), or equivalently, approximate the solution of the corresponding linear systems defined by

$$A_{44} X = A_{43}. \quad (8.2)$$

There are of course many computationally attractive ways to do this, including incomplete LU factorizations (ILU) of the block  $A_{44}$  [25] (or of each diagonal block in it in the case of multiple blocks) or the use of sparse approximate inverse factorizations [1]. In our experiments in this paper we use ILU to approximate the solution of systems like (8.2). An experimental study showing the effectiveness of ILU in this context is presented later in section 11.

A second alternative, which can be used in conjunction with an approximation of the inverse of the “atomic subdomain” matrix, is to solve the minimization problem (4.21) with a small set  $\mathcal{Q}_1$ . For example, we can consider scalar, diagonal, tridiagonal matrices, or for that matter, any prescribed sparsity pattern. In the case of scalar matrices  $D_1 = \beta I$ , we only need to solve a one parameter linear least squares problem. In the case of diagonal matrices, we need to solve  $n_3$  one parameter linear least squares problems. We also considered  $\mathcal{Q}_1$  to be the set of tridiagonal matrices. We call these three approaches the O0s, O0, and O2 methods, respectively, since they have the same spirit of the OO0 and OO2 methods in optimized Schwarz methods; see, e.g., [10], [11], and further section 10.

We choose the Frobenius norm for the minimization problems, and thus obtain in the scalar case, e.g., from (4.21),

$$\begin{aligned} \beta_0 &= \arg \min_{\beta} \|\beta B_{12} - A_{34} B_{13}\|_F = \arg \min_{\beta} \|\beta \text{vec}(B_{12}) - \text{vec}(A_{34} B_{13})\|_2 \\ &= \text{vec}(B_{12})^T \text{vec}(A_{34} B_{13}) / \text{vec}(B_{12})^T \text{vec}(B_{12}), \end{aligned} \quad (8.3)$$

where the Matlab `vec` command produces here an  $n_3 \cdot n_2$  vector with the matrix entries. In the diagonal case, we look for a diagonal matrix  $D = \text{diag}(d_1, \dots, d_{n_3})$  such that

$$\min_D \|DB_{12} - A_{34} B_{13}\|_F,$$

which can be decoupled as  $n_3$  problems for each nonzero of  $D$  using each column of  $B_{12}$  and  $A_{34} B_{13}$  to obtain

$$d_i = \arg \min_d \|d(B_{12})_i - (A_{34} B_{13})_i\|_F = (B_{12})_i^T (A_{34} B_{13})_i / (B_{12})_i^T (B_{12})_i, \quad (8.4)$$

where we have used the notation  $X_i$  to denote the  $i$ th column of  $X$ . Observe that the cost of obtaining  $\beta_0$  in (8.3) and that of obtaining the  $n_3$  values of  $d_i$  in (8.4) is essentially the same.

**REMARK 8.1.** *We will shortly validate these various choices of  $D_i$  with numerical experiments. We note however that an important open question is to find boundary matrices  $D_i$  that work well in an efficient manner, hopefully without having to compute a Schur complement or its ILU approximation.*

**9. Numerical Experiments.** For our numerical experiments, we consider the discretization of a two-dimensional advection diffusion equation of the form

$$\eta u - \nabla \cdot (a \nabla u) + b \cdot \nabla u = f, \quad (9.1)$$

where

$$a = a(\underline{x}), \quad b = \begin{bmatrix} b_1(\underline{x}) \\ b_2(\underline{x}) \end{bmatrix}, \quad \eta = \eta(\underline{x}) \geq 0, \quad (9.2)$$

with  $b_1 = y - 1/2$ ,  $b_2 = -(x - 1/2)$ ,  $\eta = x^2 \cos(x + y)^2$ ,  $a = (x + y)^2 e^{x-y}$ . We consider two domains, a square, and an L-shaped region. For our first set of experiments, our domain is  $(0, 1)^2$  partitioned in two vertical halves. We illustrate the field  $b(x, y)$  in Figure 9.1.

We use finite differences with  $h = 1/20$  in each direction resulting in a banded matrix with  $n = 400$  and a semiband of size 20. We preprocess the matrix using the reverse Cuthill-McKee algorithm (see, e.g., [13]). This results in the matrix depicted in Figure 2.1. In the same figure, we show the partition used, i.e., with  $n_1 = n_4 = 180$  and  $n_2 = n_3 = 20$ .

In Figure 9.2, we present the convergence history of our new block iterative method for this advection diffusion problem (on the square and two blocks) with  $f = 0$ , both using the optimal transmission matrices (making  $T^2 = 0$ ), and using several approximations to them. We consider the one parameter minimization of (4.21) for scalar matrices  $D_i = \beta I$  (O0s), diagonal (O0), and tridiagonal matrices (O2). The number of iterations to reduce the norm of the error below  $10^{-8}$  are 2 (as predicted by our theory), 40, 32, and 27,

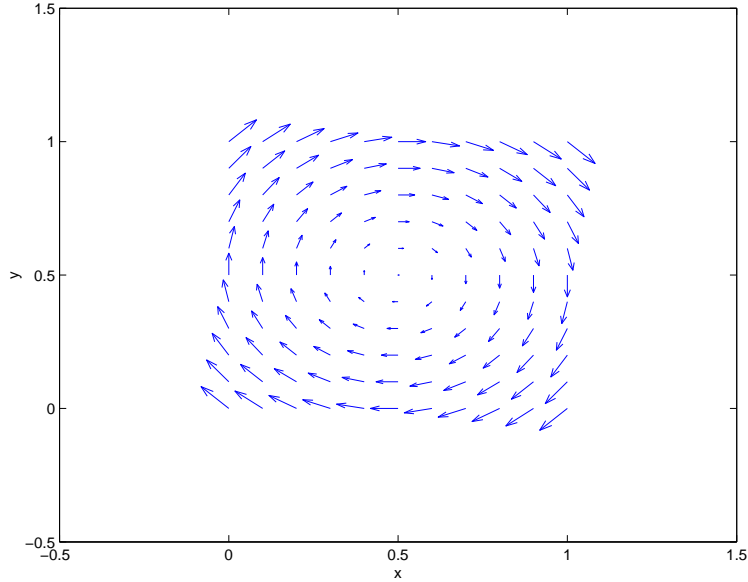


FIG. 9.1. The rotating field  $b(x, y)$  in the advection diffusion example.

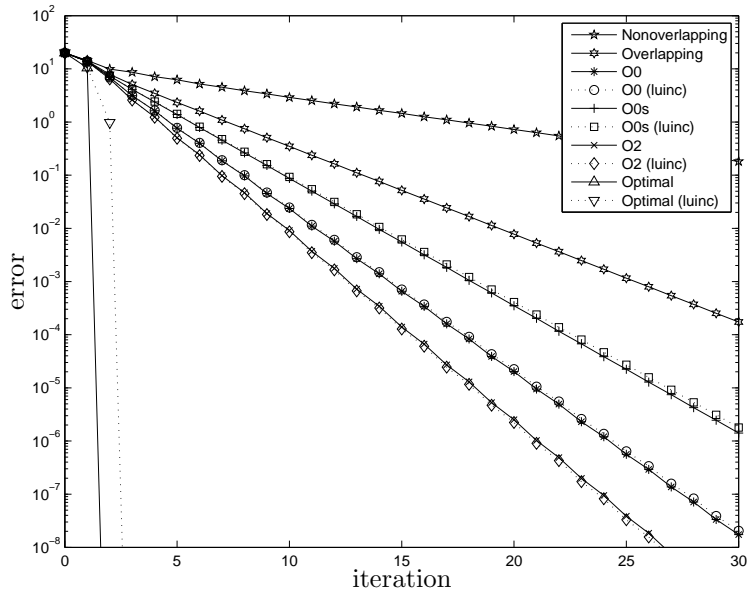


FIG. 9.2. Convergence history of Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative method, with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), and a tridiagonal matrix (called O2). All cases with the new method are also run with ILU approximations. Square domain, Two blocks.

respectively. The optimization problem (4.21) requires the evaluation of the matrix  $B_{ij}$ , which are obtained by solving linear problems in blocks of  $A$ , cf. (4.4). In a practical algorithm, it will be preferable to compute  $B_{ij}$  approximately. To simulate this, we have used an incomplete LU factorization (ILU) of blocks of  $A$  with threshold  $\tau = 1/n_3$  to compute approximations  $B_{ij}^{(ILU)}$  to the matrices  $B_{ij}$ . From those matrices, we have then computed O0 (luinc), O0s (luinc), O2 (luinc) and Optimal (luinc) transmission conditions  $D_1$  and  $D_2$ . For completeness, we also include the classical Schwarz method, both with nonoverlapping and overlapping blocks.

Note that except for the optimal case, where convergence takes three iterations, all other runs using the ILU approximations are almost identical to those with the exact solution of (8.2). We also show in the same plot the convergence history of Block Jacobi (without overlap), and RAS (overlapping). In these cases, the number of iterations to reduce the norm of the error below  $10^{-8}$  are 146 and 56, respectively. In order to allow the reader to reproduce our runs, in all experiments we use as initial approximation  $u^0$  the vector of all ones. Results with other initial vectors, have produced similar qualitative results.

For the simple case of  $\eta = 0$ ,  $b = 0$ ,  $a = 1$ , i.e., the Laplacian, for the case of the substitution  $D_i = \beta I$ , we computed the value of  $\|T\|$ ,  $\|T^2\|$ , (the 2-norm) and  $\rho(T)$ , for varying values of the parameter  $\beta$ . They are reported in Figure 9.3, together with the value obtained by the minimization (4.21). It can be appreciated that while  $\rho(T) < 1$  for all values of  $\beta \in [0, 1]$ ,  $\|T\| > 1$  for most of those values. Furthermore the curve for  $\|T^2\|$  is pretty close to that of  $\rho(T)$  for a wide range of values of  $\beta$ .

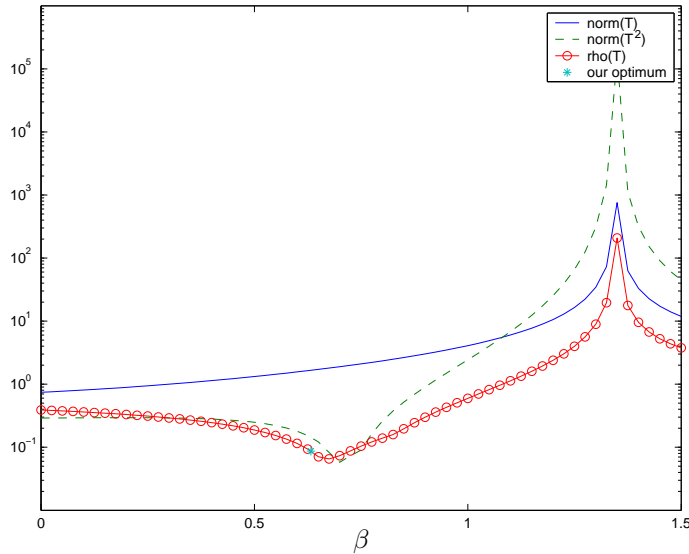


FIG. 9.3.  $\|T\|$ ,  $\|T^2\|$ , and  $\rho(T)$ , for  $D_i = \beta I$  for varying  $\beta$ . Laplacian.

In Figure 9.4 we show convergence curves for the same cases, when we consider four blocks of the matrix each of order 100 (with the same overlap 20). Observe that neither Block Jacobi, nor the method with the one-parameter approximation of the optimal transmission matrix (O0s) converge in this case. Note also that the use of the ILU in the approximation of the optimal transmission matrix results in convergence to an error below  $10^{-8}$  in five iterations.

We ran the same experiment on an L-shaped domain, consisting of the square, where the upper right quarter is cut-off. The order of the matrix  $A$  is now 300, and we kept the same overlap of 20, which is given by the band size. Thus, in the two block-case, each atomic block is of order 170, and in the four-block case, it is of order 95. Results for two and four blocks are reported in Figures 9.5 and 9.6. The rapid convergence, compared to the square subdomain, is due in part to the fact that the overlap is relatively larger in this case.

**9.1. Experiments for the new iterative method. Multiplicative case.** We run the new modified multiplicative restricted Schwarz method for the same advection diffusion problem (9.1)-(9.2), and  $f = 0$ , together with Block Gauss-Seidel (without overlap) and MRAS (overlapping), both in the case of a square domain and the L-shaped domain. We considered in addition to the optimal transmission matrices (which make  $T^2 = O$ ), the same approximations to them as in the additive case, together with the use of ILU to approximate the solution of (8.2) for all these cases. The four sets of experiments are presented in Figures 9.7-9.10. Here we can make similar observations as in the experiments in the additive case, except that as expected, the multiplicative method is faster; cf. Proposition 5.1 and the comments following it.

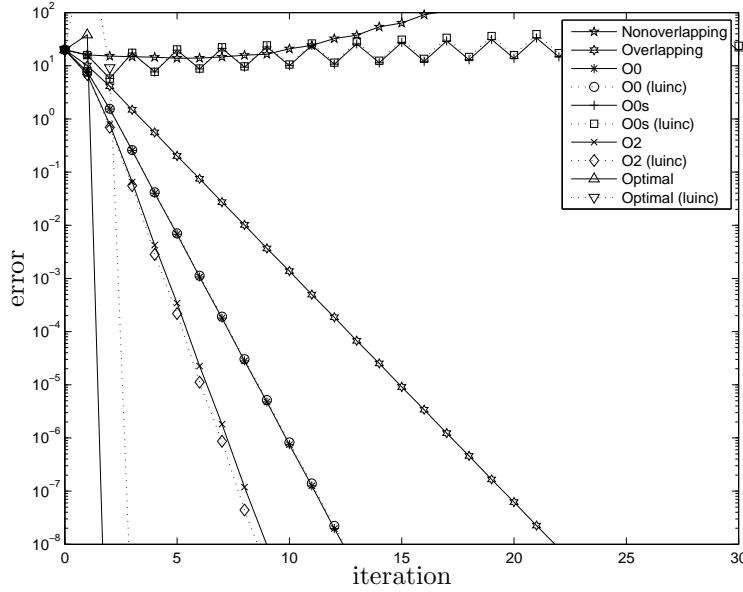


FIG. 9.4. Convergence history of Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative method, with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), and a tridiagonal matrix (called  $O2$ ). All cases with the new method are also run with ILU approximations. Square domain, four blocks.

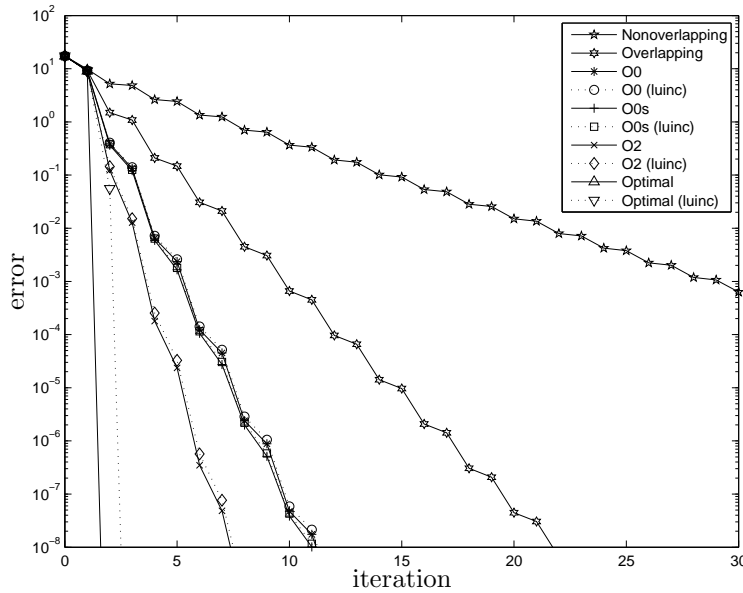


FIG. 9.5. Convergence history of Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative method, with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), and a tridiagonal matrix (called  $O2$ ). All cases with the new method are also run with ILU approximations. L-shaped domain, two blocks.

**9.2. Experiments with preconditioned GMRES.** We consider in this section the use of the new block iterative method as a preconditioner. We run preconditioned GMRES for the same advection diffusion problem (9.1)-(9.2) for the square and the L-shaped domains. In both cases, we treated two and four blocks. We considered both the additive preconditioner (6.1), as well as the multiplicative one. We report the eight sets of experiments in Figures 9.11-9.18. As before, our experiments are with the optimal transmission

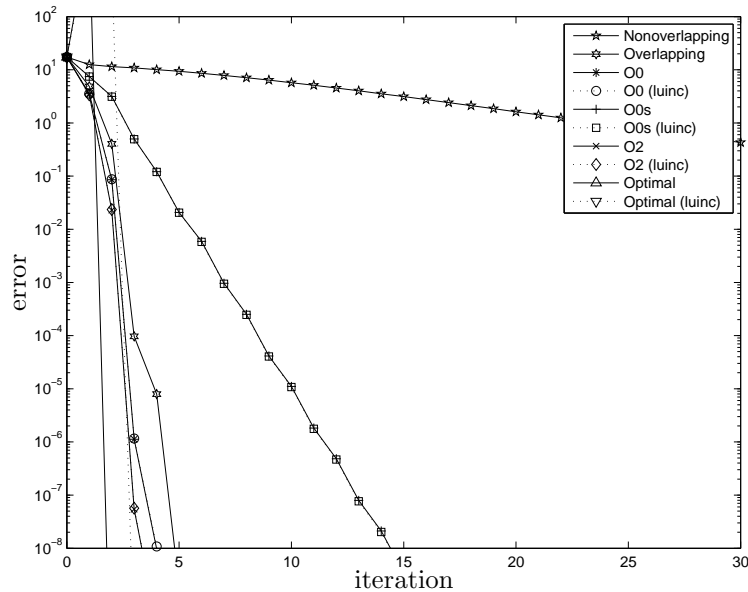


FIG. 9.6. Convergence history of Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative method, with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), and a tridiagonal matrix (called O2). All cases with the new method are also run with ILU approximations. L-shaped domain, four blocks.

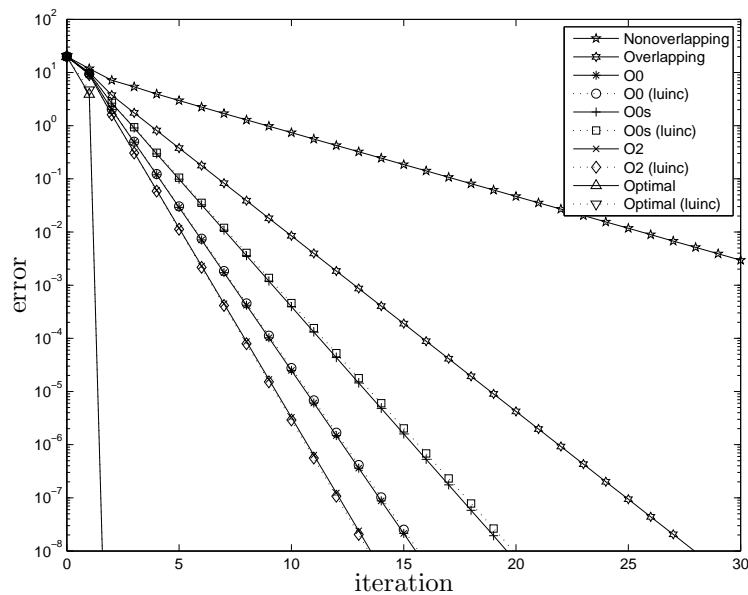


FIG. 9.7. Convergence history of the advection diffusion problem: Block Gauss-Seidel (non-overlapping), Multiplicative RAS (overlapping), with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), and a tridiagonal matrix (called O2). All cases with the new method are also run with ILU approximations. Square domain, two blocks.

matrices, and with the approximations O0s, O0, and O2, with and without the use of ILU. For completeness, we also used the Block Jacobi (or Block Gauss-Seidel) preconditioner without overlap and with overlap (RAS or MRAS). In the figures, we show the preconditioned residual norm  $\|M^{-1}(f - Au_k)\|$ , for the four different preconditioners  $M^{-1}$ . Note that even in the cases where the iterative method fails to converge, we have a viable preconditioner.



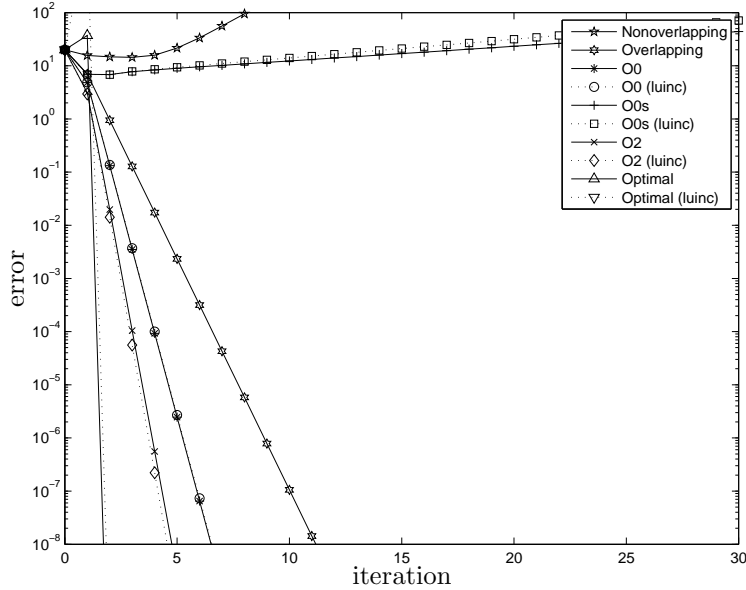


FIG. 9.8. Convergence history of the advection diffusion problem: Block Gauss-Seidel (non-overlapping), Multiplicative RAS (overlapping), with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), and a tridiagonal matrix (called  $O2$ ). All cases with the new method are also run with ILU approximations. Square domain, four blocks.

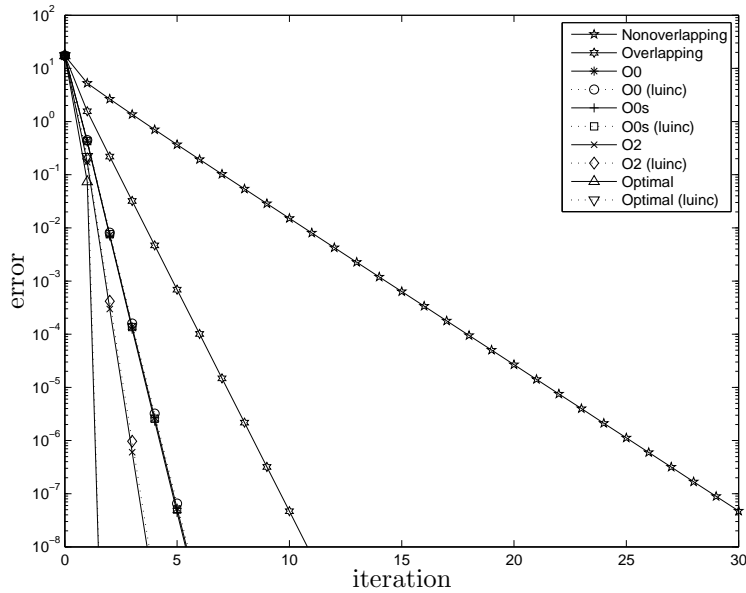


FIG. 9.9. Convergence history of the advection diffusion problem: Block Gauss-Seidel (non-overlapping), Multiplicative RAS (overlapping), with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), and a tridiagonal matrix (called  $O2$ ). All cases with the new method are also run with ILU approximations. L-shaped domain, two blocks.

We remark that Proposition 6.1 applies and therefore with the optimal transmission matrices one expects convergence of a preconditioned minimal residual method in at most two iterations. This is reflected in our numerical experiments.

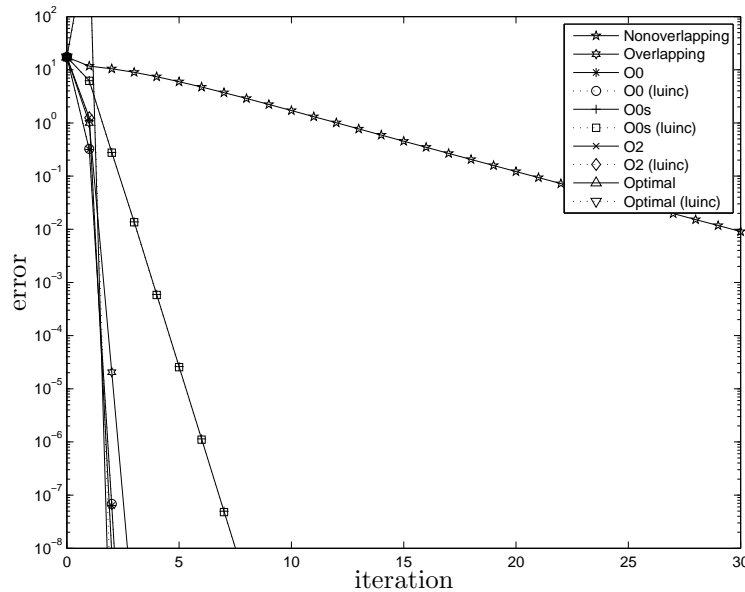


FIG. 9.10. Convergence history of the advection diffusion problem: Block Gauss-Seidel (non-overlapping), Multiplicative RAS (overlapping), with an optimal transmission matrix, its approximation with a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), and a tridiagonal matrix (called  $O2$ ). All cases with the new method are also run with ILU approximations.  $L$ -shaped domain, four blocks.

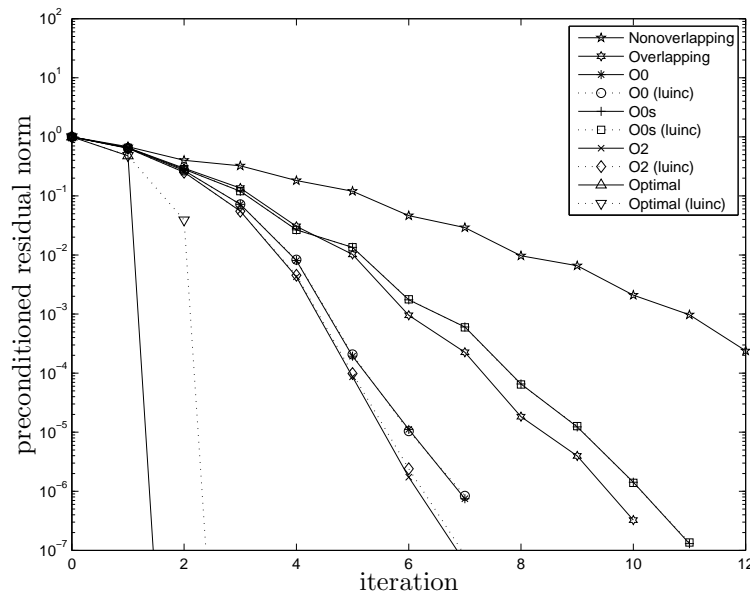


FIG. 9.11. Convergence history of GMRES preconditioned with Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative additive preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), or a tridiagonal matrix (called  $O2$ ). ILU approximations are also reported. Square domain, two blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

## 10. Asymptotic convergence factor estimates for a model problem using Fourier analysis.

In this section we consider a problem on a simple domain, so we can use Fourier analysis to calculate the optimal parameters as is usually done in optimized Schwarz methods; see, e.g., [10]. We use this analysis to compute the asymptotic convergence factor of the optimized Schwarz iterative method, and compare it to what we obtain with our algebraic counterpart.

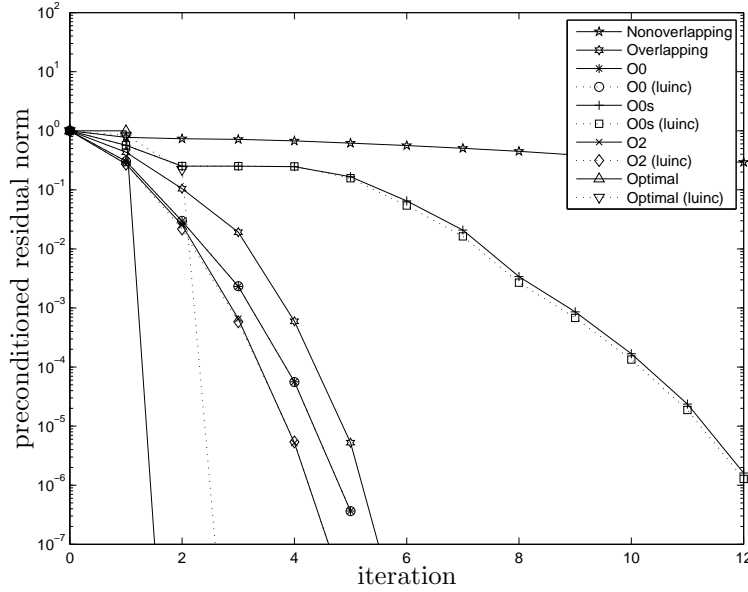


FIG. 9.12. Convergence history of GMRES preconditioned with Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative additive preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), or a tridiagonal matrix (called  $O2$ ). ILU approximations are also reported. Square domain, four blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

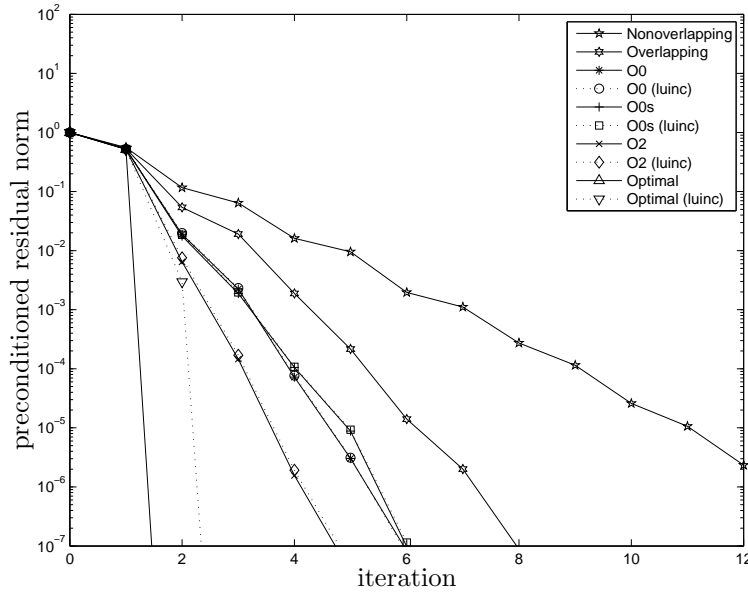


FIG. 9.13. Convergence history of GMRES preconditioned with Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative additive preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called  $O0s$ ), a diagonal matrix (called  $O0$ ), or a tridiagonal matrix (called  $O2$ ). ILU approximations are also reported. L-shaped domain, two blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

The model problem we consider is  $-\Delta u = f$  in the (horizontal) strip  $\Omega = \mathbb{R} \times (0, L)$ , with Dirichlet conditions  $u = 0$  on the boundary  $\partial\Omega$ , i.e., at  $x = 0, L$ . We discretize the continuous operator on a grid whose interval is  $h$  in both the  $x$  and  $y$  directions; i.e., with vertices at  $(jh, kh)$ . We assume that  $h = L/(m + 1)$ , so that there are  $m$  degrees of freedom along the  $y$  axis, given by  $y = h, 2h, \dots, mh$ . The stiffness matrix is

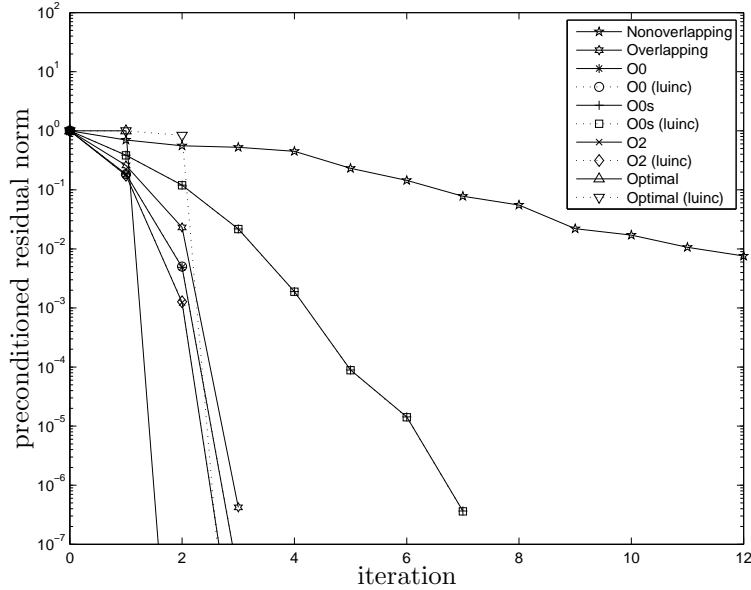


FIG. 9.14. Convergence history of GMRES preconditioned with Block Jacobi (non-overlapping), RAS (overlapping), and the new block iterative additive preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), or a tridiagonal matrix (called O2). ILU approximations are also reported. L-shaped domain, four blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

infinite and block-tridiagonal of the following form,

$$A = \begin{bmatrix} \ddots & & & & & & & & \\ & \ddots & & & & & & & \\ & & -I & E & -I & & & & \\ & & & -I & E & -I & & & \\ & & & & & & \ddots & & \\ & & & & & & & \ddots & \\ & & & & & & & & \ddots \end{bmatrix},$$

where  $I$  is the  $m \times m$  identity matrix and  $E$  is the  $m \times m$  tridiagonal matrix  $E = \text{tridiag}(-1, 4, -1)$ . This is the stiffness matrix obtained when we discretize with the Finite Element Method using piecewise linear elements. Since the matrix is infinite, we must specify the space that it acts on. We look for solutions in the space  $\ell^2(\mathbb{Z})$  of square-summable sequences. In particular, a solution to  $Au = b$  must vanish at infinity. This is similar to solving the Laplace problem in  $H_0^1(\Omega)$ , where the solution also vanishes at infinity.

We use the subdomains  $\Omega_1 = (-\infty, h) \times (0, L)$  and  $\Omega_2 = (0, \infty) \times (0, L)$ , leading to the decomposition

$$\begin{bmatrix} A_{11} & A_{12} & O & O \\ A_{21} & A_{22} & A_{23} & O \\ O & A_{32} & A_{33} & A_{34} \\ O & O & A_{43} & A_{44} \end{bmatrix} = \begin{bmatrix} \ddots & & & & & & & & \\ & \ddots & & & & & & & \\ & & -I & E & -I & & & & \\ & & & -I & E & -I & & & \\ & & & & & & \ddots & & \\ & & & & & & & \ddots & \\ & & & & & & & & \ddots \end{bmatrix}, \tag{10.1}$$

i.e., we have in this case  $n_2 = n_3 = m$ .

In optimized Schwarz Methods, one uses either Robin conditions (OO0) on the artificial interface, or a second order tangential condition (OO2). If we discretize these transmission conditions using the piecewise linear spectral element method (i.e., by replacing the integrals with quadrature rules), we get that  $S_i = \frac{1}{2}E + pI$  for the OO0 iteration, where the scalar  $p$  is typically optimized by considering a continuous version of the problem, and using Fourier transforms; see, e.g., [10]. Likewise, for the OO2 iteration, we get that

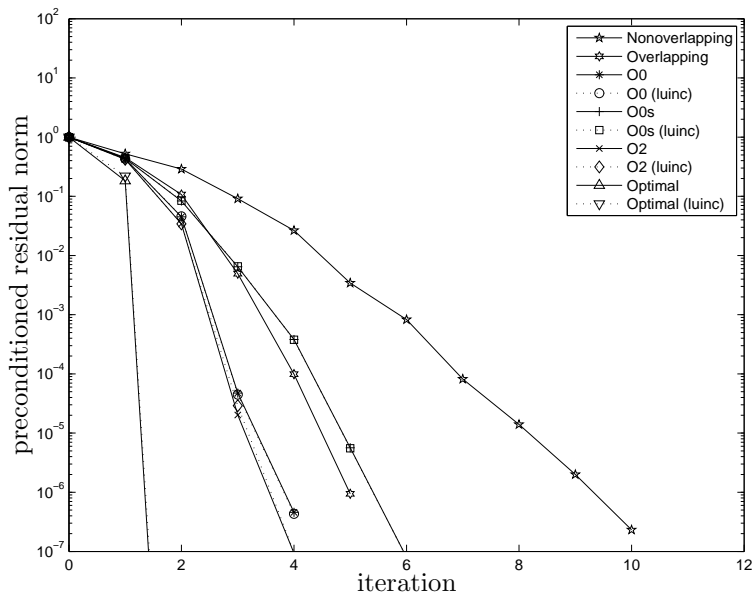


FIG. 9.15. Convergence history of GMRES preconditioned with Block Gauss-Seidel (non-overlapping), MRAS (overlapping), and the new block iterative multiplicative preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), or a tridiagonal matrix (called O2). ILU approximations are also reported. Square domain, two blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

$S_i = \frac{1}{2}E + pI + qJ$ , where  $J$  is the tridiagonal matrix  $\text{tridiag}(1, 0, 1)$ , and where  $p$  and  $q$  are optimized using a continuous version of the problem. In the current paper, we have also proposed the choices  $S_i = E - \beta I$  (O0) and  $S_i = E - \beta I + \gamma J$  (O2). The O2 and OO2 methods are related via

$$4 - \beta = 2 + p \text{ and } \gamma - 1 = q - 1/2.$$

However, the O0 method is new and is not directly comparable to the OO0 method, since the off-diagonal entries of  $E - \beta I$  cannot match the off-diagonal entries of  $E/2 + pI$ .<sup>2</sup>

We now obtain an estimate of the convergence factor for the proposed new method.

LEMMA 10.1. *Let  $A$  be given by (10.1). For  $S_1 = A_{33} + D_1$  with  $D_1 = -\beta I$  and  $\beta \in \mathbb{R}$ , the convergence factor estimate (4.16) is*

$$\|(I + D_1 B_{33})^{-1}(D_1 B_{12} - A_{34} B_{13})\| = \max_{k=1, \dots, m} \left| \frac{-\beta + e^{-w(k)h}}{1 - \beta e^{-w(k)h}} \right| e^{-2w(k)h}, \quad (10.2)$$

where  $w(k) = w(k, L, h)$  is the unique positive solution of the relation

$$\cosh(w(k)h) = 2 - \cos\left(\frac{\pi kh}{L}\right). \quad (10.3)$$

Note that  $w(k)$  is a monotonically increasing function of  $k \in [1, m]$ .

*Proof of Lemma 10.1.* Let  $F$  be the symmetric orthogonal matrix whose entries are

$$F_{jk} = \sqrt{\frac{m+1}{2}} \sin(\pi jk/(m+1)). \quad (10.4)$$

<sup>2</sup>In OO0, the optimized  $p$  is positive because it represents a Robin transmission condition. The best choice of  $p$  is small and hence the corresponding row sums of  $\tilde{A}_i$  are almost zero, but positive. We have chosen  $D_i = -\beta I$  in order to achieve similar properties for the rows our  $\tilde{A}_i$ .

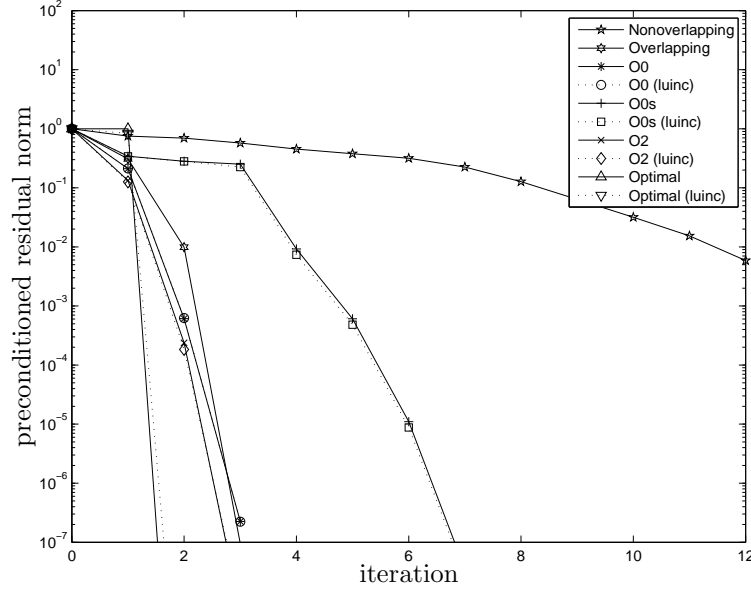


FIG. 9.16. Convergence history of GMRES preconditioned with Block Gauss-Seidel (non-overlapping), MRAS (overlapping), and the new block iterative multiplicative preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), or a tridiagonal matrix (called O2). ILU approximations are also reported. Square domain, four blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

Consider the auxiliary problem

$$\begin{bmatrix} A_{11} & A_{12} & O \\ A_{21} & A_{22} & A_{23} \\ O & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} C_1 \\ C_2 \\ C_3 \end{bmatrix} = \begin{bmatrix} O \\ O \\ F \end{bmatrix} \quad (10.5)$$

Observe that since  $F^2 = I$ , we have that

$$\begin{bmatrix} B_{31} \\ B_{32} \\ B_{33} \end{bmatrix} = \begin{bmatrix} C_1 \\ C_2 \\ C_3 \end{bmatrix} F. \quad (10.6)$$

We can solve (10.5) for the unknowns  $C_1, C_2, C_3$ . By considering (10.1), and since  $A_{34} = [-I \ O \ O \ \dots]$ , we see that

$$\begin{bmatrix} A_{11} & A_{12} & O & O \\ A_{21} & A_{22} & A_{23} & O \\ O & A_{32} & A_{33} & -I \end{bmatrix} \begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ F \end{bmatrix} = \begin{bmatrix} O \\ O \\ O \end{bmatrix}.$$

Hence, we are solving the discrete problem

$$\begin{cases} (L_h u)(x, y) = 0 \text{ for } x = \dots, -h, 0, h \text{ and } y = h, 2h, \dots, mh; \\ u(2h, y) = \sqrt{\frac{m+1}{2}} \sin(\pi ky/L) \text{ for } y = h, 2h, \dots, mh; \text{ and} \\ u(x, 0) = u(x, L) = 0 \text{ for } x = \dots, -h, 0, h; \end{cases} \quad (10.7)$$

where the discrete Laplacian  $L_h$  is given by

$$(L_h u)(x, y) = 4u(x, y) - u(x-h, y) - u(x+h, y) - u(x, y-h) - u(x, y+h).$$

The two basic solutions to the difference equation are:

$$u_{\pm}(x, y) = e^{\pm w(k)x} \sin(\pi ky/L),$$

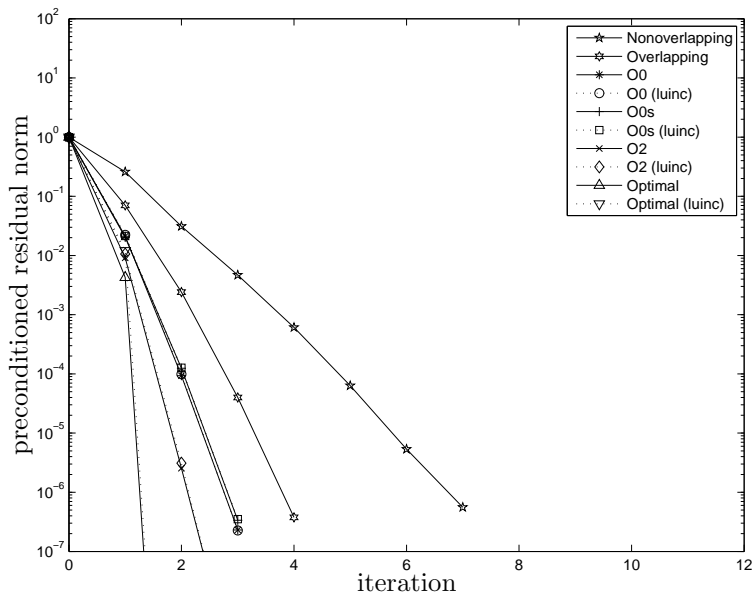


FIG. 9.17. Convergence history of GMRES preconditioned with Block Gauss-Seidel (non-overlapping), MRAS (overlapping), and the new block iterative multiplicative preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), or a tridiagonal matrix (called O2). ILU approximations are also reported. L-shaped domain, two blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

where  $w(k)$  is the unique positive solution of (10.3).

The subdomain  $\Omega_1$  does not contain the  $x = \infty$  boundary, but it does contain the  $x = -\infty$  boundary. Since we are looking for solutions that vanish at infinity (which, for  $\Omega_1$ , means  $x = -\infty$ ), the unique solution for the given Dirichlet data at  $x = 2h$  is therefore

$$u(x, y) = \left( \sqrt{\frac{m+1}{2}} e^{-2w(k)h} \right) e^{w(k)x} \sin(\pi ky/L).$$

Using (10.4), this gives the formula

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \end{bmatrix} = \begin{bmatrix} \vdots \\ \frac{FD(3h)}{FD(2h)} \\ \frac{FD(2h)}{FD(h)} \end{bmatrix},$$

where  $D(\xi)$  is the diagonal  $m \times m$  matrix whose  $(k, k)$ th entry is  $e^{-w(k)\xi}$ . Hence, from (10.6),

$$\begin{bmatrix} B_{31} \\ B_{32} \\ B_{33} \end{bmatrix} = \begin{bmatrix} \vdots \\ \frac{FD(3h)F}{FD(2h)F} \\ \frac{FD(h)F}{FD(h)F} \end{bmatrix}.$$

In other words, the matrix  $F$  diagonalizes all the  $m \times m$  blocks of  $B_3^{(1)}$ . Observe that  $F$  also diagonalizes  $J$  and  $E = 4I - J$ , and hence all the blocks of  $A$ ; see the right-hand-side of (10.1). A similar reasoning shows

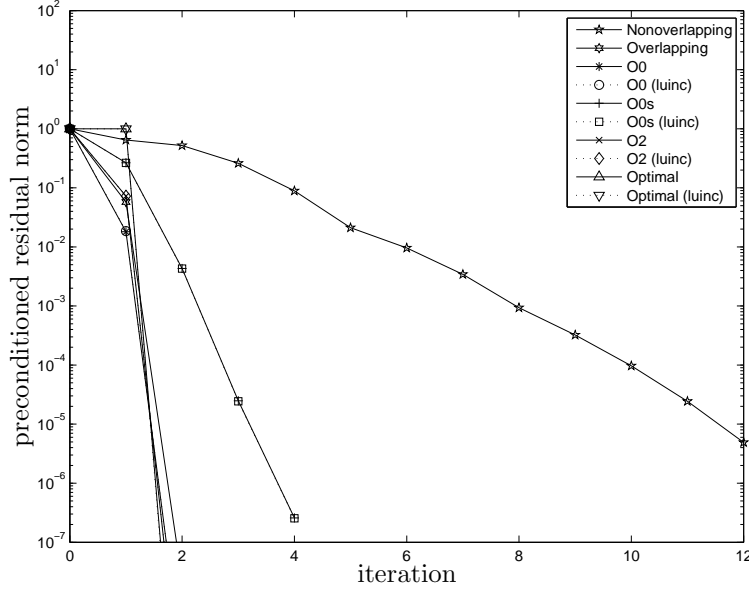


FIG. 9.18. Convergence history of GMRES preconditioned with Block Gauss-Seidel (non-overlapping), MRAS (overlapping), and the new block iterative multiplicative preconditioner, with optimal transmission matrices, and with the addition of a scalar matrix ( $\beta I$ , called O0s), a diagonal matrix (called O0), or a tridiagonal matrix (called O2). ILU approximations are also reported. L-shaped domain, four blocks. Preconditioned residual  $\|M^{-1}(f - Au_k)\|$ , for the different  $M^{-1}$ .

that  $F$  diagonalizes also the  $m \times m$  blocks of  $B_1^{(2)}$ :

$$\begin{bmatrix} B_{11} \\ B_{12} \\ B_{13} \end{bmatrix} = \begin{bmatrix} \frac{FD(h)F}{FD(2h)F} \\ \frac{FD(2h)F}{FD(3h)F} \\ \vdots \end{bmatrix}.$$

Hence, the convergence factor estimate (4.16) for our model problem, is given by

$$\|(I + D_1 B_{33})^{-1}(D_1 B_{12} - A_{34} B_{13})\| = \|F(I - \beta D(h))^{-1}(-\beta D(2h) + D(3h))F\|,$$

which leads to (10.2).  $\square$

LEMMA 10.2. In the limit as  $h$  tends to zero, the optimized parameter  $\beta$  behaves asymptotically like

$$\beta_{opt} = 1 - c \left(\frac{h}{L}\right)^{2/3} + O(h), \quad (10.8)$$

where  $c = (\pi/2)^{2/3} \approx 1.35$ . The resulting convergence factor is

$$\rho_{opt} = 1 - \left(\frac{32\pi h}{L}\right)^{1/3} + O(h^{2/3}). \quad (10.9)$$

*Proof.* We begin this proof with some notation and a few observations. Let

$$r(w, h, \beta) = \frac{-\beta + e^{-wh}}{1 - \beta e^{-wh}} e^{-2wh}.$$

Then, the convergence factor estimate (10.2) is bounded by and very close to (cf. the argument in [17])

$$\rho(L, h, \beta) = \max_{w \in [w(1), w(m)]} |r(w, h, \beta)|,$$



where  $w(k) = w(k, h, L)$  is given by (10.3). Clearly,  $r(w, h, 0) < r(w, h, \beta)$  whenever  $\beta < 0$ , hence we will optimize over the range  $\beta \geq 0$ . Conversely, we must have  $1 - \beta e^{-wh} > 0$  for every  $w \in [w(1), w(m)]$ , to avoid an explosion in the denominator of (10.2). By using the value  $w = w(1)$ , we find that

$$0 \leq \beta < 2 - \cos(\pi h/L) + \sqrt{(3 - \cos(\pi h/L))(1 - \cos(\pi h/L))} = 1 + \pi \frac{h}{L} + O(h^2). \quad (10.10)$$

We are therefore optimizing  $\beta$  in a closed interval  $[0, \beta_{\max}] = [0, 1 + \pi h/L + O(h^2)]$ .

We divide the rest of the proof into seven steps.

Step 1: we show that  $\beta_{opt}$  is obtained by solving an equioscillation problem. We define the set  $W(\beta) = W(L, h, \beta) = \{w > 0 \text{ such that } |r(w, h, \beta)| = \rho(L, h, \beta)\}$ , and we now show that, if  $\rho(\beta) = \rho(L, h, \beta)$  is minimized at  $\beta = \beta_{opt}$ , then  $\#W(\beta_{opt}) > 1$ . By the Envelope Theorem, if  $W(\beta) = \{w^*\}$  is a singleton, then  $\rho(\beta) = \rho(L, h, \beta)$  is a differentiable function of  $\beta$  and its derivative is  $\frac{\partial}{\partial \beta} |r(w^*, h, \beta)|$ . Since

$$\frac{\partial}{\partial \beta} r(w, h, \beta) = \frac{e^{-2wh} - 1}{(\beta e^{-wh} - 1)^2} e^{-2wh}, \quad (10.11)$$

we get

$$0 = \frac{d\rho}{d\beta}(\beta_{opt}) = \frac{e^{-2w^*h} - 1}{(\beta_{opt} e^{-w^*h} - 1)^2} e^{-2w^*h} \operatorname{sgn}(r),$$

which is impossible. Therefore,  $\#W(\beta_{opt}) \geq 2$ ; i.e.,  $\beta_{opt}$  is obtained by equioscillating  $r(w)$  at at least two distinct points of  $W(\beta)$ .

Step 2: We find the critical values of  $r(w) = r(w, h, \beta)$  as a function of  $w$  alone. By differentiating, we find that the critical points are  $\pm w_{\min}$ , where

$$w_{\min} = w_{\min}(\beta, h) = -\ln \left( \frac{1}{4} \frac{3 + \beta^2 - \sqrt{9 - 10\beta^2 + \beta^4}}{\beta} \right) h^{-1}. \quad (10.12)$$

In the situation  $\beta > 1$ ,  $w_{\min}$  is complex, and hence there is no critical point. In the situation  $\beta = 1$ , we have  $w_{\min} = 0$ , which is outside of the domain  $[w(1), w(m)]$  of  $r(w)$ . Since  $r(w)$  is differentiable over its domain  $[w(1), w(m)]$ , its extrema must be either at critical points, or at the endpoints of the interval; i.e.,

$$W(\beta_{opt}) \subset \{w(1), w_{\min}, w(m)\}.$$

We now compute  $\beta_{opt}$  by assuming  $W(\beta_{opt}) = \{w(1), w_{\min}\}$ . We will subsequently verify this assumption.

Step 3: We now consider the solution(s) of the equioscillation problem  $W(\beta_{opt}) = \{w(1), w_{\min}\}$ , and we show that  $r(w(1)) > 0$  and  $r(w_{\min}) < 0$ , and that  $r(w(1)) + r(w_{\min}) = 0$ . Since  $W(\beta) = \{w(1), w_{\min}\}$ , we must have that  $w_{\min} > w(1)$ . If we had that  $r(w(1)) = r(w_{\min})$ , the mean value theorem would yield another critical point in the interval  $(w(1), w_{\min})$ . Therefore, it must be that  $r(w(1)) + r(w_{\min}) = 0$ . We now check that,  $r(w_{\min}) < 0$ . Indeed,  $r(w)$  is negative when  $w$  is large, and  $r(+\infty) = 0$ . If we had  $r(w_{\min}) > 0$ , there would be a  $w' > w_{\min}$  such that  $r(w') < 0$  is minimized, creating another critical point. Since  $w_{\min}$  is the only critical point, it must be that  $r(w_{\min}) < 0$ . Hence,  $r(w(1)) > 0$ .

Step 4: We show that there is a unique solution to the equioscillation problem  $W(\beta_{opt}) = \{w(1), w_{\min}\}$ , characterized by  $r(w(1)) + r(w_{\min}) = 0$ . From (10.11), we see that  $\frac{\partial r}{\partial \beta}(w(1), h, \beta) < 0$ , and likewise,

$$\frac{\partial(r(w_{\min}(\beta, h), h, \beta))}{\partial \beta} = \frac{\partial r}{\partial \beta}(w_{\min}, h, \beta) + \overbrace{\frac{\partial r}{\partial w}(w_{\min}, h, \beta)}^{=0} \frac{\partial w_{\min}}{\partial \beta}(\beta, h) < 0.$$

Combining the facts that  $r(w(1)) > 0$  and  $r(w_{\min}) < 0$  are both decreasing in  $\beta$ , there is a unique value of  $\beta = \beta_{opt}$  such that  $r(w(1)) + r(w_{\min}) = 0$ ; this  $\beta_{opt}$  will minimize  $\rho(L, h, \beta_{opt})$  under the assumption that  $W(\beta) = \{w(1), w_{\min}\}$ .

Step 5: We give an asymptotic formula for the unique  $\beta_{opt}$  solving the equioscillation problem  $W(\beta_{opt}) = \{w(1), w_{\min}\}$ . To this end, we make the ansatz  $\beta = 1 - c(h/L)^{2/3}$ , and we find that

$$r(w(1)) = 1 - \frac{2\pi}{c} \left(\frac{h}{L}\right)^{1/3} + O(h^{2/3}) \quad \text{and} \quad (10.13)$$

$$r(w_{\min}) = -1 + 4\sqrt{c} \left(\frac{h}{L}\right)^{1/3} + O(h^{2/3}). \quad (10.14)$$

Hence, the equioscillation occurs when  $c = (\pi/2)^{2/3}$ .

Step 6: We now show that the equioscillation  $W(\beta_{opt}) = \{w(1), w_{\min}\}$  occurs when  $w_{\min} \in (w(1), w(n))$ . Let  $\beta_{opt} = 1 - c(h/L)^{2/3} + O(h)$ . Then, from (10.12) and (10.3),

$$w_{\min} = \frac{\sqrt{c}}{L^{1/3}h^{2/3}} + O(h^{2/3}) < w(n) = \frac{\operatorname{arccosh}(3)}{h} + O(h),$$

provided that  $h$  is sufficiently small.

Step 7: If  $\beta < \beta_{opt}$ , then  $\rho(\beta) > \rho(\beta_{opt})$ . Indeed, we see from (10.11) that  $\frac{\partial r}{\partial \beta}(w_1, h, \beta) < 0$ . Hence, if  $\beta < \beta_{opt}$ , then  $\rho(\beta) \geq r(w_1, h, \beta) > r(w_1, h, \beta_{opt}) = \rho(\beta_{opt})$ . A similar argument shows that if  $\beta > \beta_{opt}$ , then  $\rho(\beta) > \rho(\beta_{opt})$ .

We therefore conclude that the  $\beta_{opt}$  minimizing  $\rho(\beta)$  is the unique solution to the equioscillation problem  $W(\beta_{opt}) = \{w(1), w_{\min}\}$ , and its asymptotic expansion is given by (10.8). We compute a series expansion of  $\rho(L, h, 1 - c(h/L)^{2/3})$  to obtain (10.9).  $\square$

This shows that the O0 method converges at a rate similar to the OO0 method. In a practical problem where the domain is not a strip, or the partial differential equation is not the Laplacian, if one wants to obtain the best possible convergence factor, then one should solve the nonlinear optimization problem (4.20). We now consider the convergence factor obtained when instead the linear minimization problem (4.21) is solved, or even if a minimization of the norm of (4.17) is obtained.

LEMMA 10.3. *For our model problem, the solution of the optimization problem*

$$\beta_0 = \arg \min_{\beta} \|(-\beta B_{12} + B_{13})\|$$

is

$$\beta_0 = \frac{3}{2} \frac{\left((1 + \sqrt{2})^{2/3} - 1\right) s}{\sqrt[3]{1 + \sqrt{2}}}, \quad (10.15)$$

where

$$s^{-1} = 2 - \cos\left(\frac{\pi h}{L}\right) + \sqrt{3 - \cos\left(\frac{\pi h}{L}\right)} \sqrt{1 - \cos\left(\frac{\pi h}{L}\right)}. \quad (10.16)$$

The resulting asymptotics are

$$\beta_0 = 0.894\dots - 2.8089\dots(h/L) + O(h^2) \quad \text{and} \quad \rho_0 = 1 - 62.477\dots(h/L) + O(h^2). \quad (10.17)$$

We mention that the classical Schwarz iteration as implemented, e.g., using RAS, is obtained with  $\beta = 0$ , yielding the asymptotic convergence factor

$$1 - 9.42\dots(h/L) + O(h^2).$$

In other words, our algorithm is asymptotically  $62.477/9.42 \approx 6.6$  times faster than a classical Schwarz iteration, in the sense that it will take about 6.6 iterations of a classical Schwarz method to equal one of our O0 method, with the parameter  $\beta = \beta_0$ , if  $h$  is small.

*Proof of Lemma 10.3* By proceeding as in the proof of Lemma 10.1, we find that

$$\|(-\beta B_{12} + B_{13})\| = \max_{w \in \{w(1), \dots, w(m)\}} |(\beta - e^{-wh})e^{-2wh}|.$$

We thus set

$$r_0(w) = r_0(w, \beta, h) = (\beta - e^{-wh})e^{-2wh}.$$

The function  $r_0(w)$  has a single extremum at  $w^* = w^*(h, \beta) = (1/h) \ln(3/(2\beta))$ . We further find that

$$r_0(w^*) = \frac{4}{27}\beta^3,$$

independently of  $h$ . We look for an equioscillation by setting

$$r_0(w(1, L, h), \beta_0, h) = r_0(w^*(\beta_0, h), \beta_0, h);$$

that is,

$$\frac{4}{27}\beta_0^3 + s^{-2}\beta_0 + s^{-3} = 0.$$

Solving for the unknown  $\beta$  yields (10.3) and (10.15). Substituting  $\beta = \beta_0$  and  $w = w(1)$  into (10.2) and taking a series expansion in  $h$  gives (10.17).  $\square$

### 11. Scaling experiments for a rectangular region with an approximate Schur complement.

In this last section, we present an experimental study showing the effectiveness of ILU when used to approximate the solution of systems with the atomic blocks, as in (8.2). To that end, we consider a simpler PDE, namely the Laplacian on the rectangle  $(-1, 1) \times (0, 1)$ ; i.e.,  $a = 1$ ,  $b = 0$ ,  $\eta = 0$  in (9.2). We use two overlapping subdomains, which then correspond to two overlapping blocks in the band matrix. We consider several systems of equations of increasing order, by decreasing the value of the mesh parameter  $h$ . We show that for this problem, the rate of convergence of our method when ILU is used, stays very close to that obtained with the exact solution of the systems with the atomic blocks. Furthermore, in the one-parameter approximations to the transmission matrices, the value of this parameter  $\beta$  computed with ILU is also very close to that obtained with the exact solutions. See Figure 11.1.

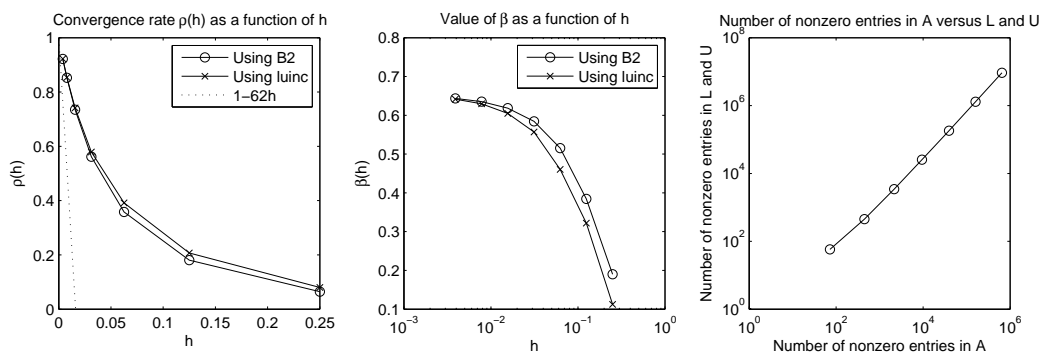


FIG. 11.1. Convergence factor of the new iterative method, using an approximate calculation for  $\beta$ . Left: the convergence factor as a function of  $h$ , for the value of  $\beta$  obtained from (8.3), as well as the value obtained by using the ILU approximation. We also plot the line  $1 - 62h$ , cf. (10.9). Middle: the two  $\beta$  parameters, as a function of  $h$ . Right: the number of nonzero entries of the  $L$  and  $U$  factors, compared to the number of nonzero entries of  $A$ .

We mesh the rectangle  $(-1, 1) \times (0, 1)$ , with a regular mesh, with the mesh interval  $h$ . The vertices of the first subdomain are all the vertices in  $(-1, h/2] \times (0, 1)$ , and the vertices of the second subdomain are all the vertices in  $[-h/2, 1) \times (0, 1)$ . We choose  $h$  in such a way that there are no vertices on the line  $x = 0$ , but instead the interfaces are at  $x = \pm h/2$ . By ordering the vertices such that all the vertices in  $x < -h/2$

occur first, then all the vertices with  $x = -h/2$  occur second, then all the vertices with  $x = h/2$  occur third, and finally all the vertices with  $x > h/2$  occur fourth, we obtain a stiffness matrix of the form (2.1), with additionally  $A_{13} = A_{42} = O$ . More specifically, the matrix  $A$  is a finite matrix of the form (10.1). We use  $D_i = \beta I$ , and use the optimized parameter given by (8.3).

As with the other experiments, we computed using the matrices  $B_{ij}$ . Since  $B_{12}, B_{13}$  are difficult to compute, we also used an Incomplete LU decomposition of  $A_2$  to obtain approximations  $B_{12}^{(ILU)}$  and  $B_{13}^{(ILU)}$ , which we then plug into (8.3). To obtain a good value of  $\beta$ , we used a drop tolerance of  $1/(n_2 + n_3 + n_4)$ , where  $n_2 + n_3 + n_4$  is the dimension of  $A_2$ . Using this drop tolerance, we found that the  $L$  and  $U$  factors have approximately ten times as many nonzero entries as the matrix  $A$ . Since the two subdomains are symmetrical, the value of  $\beta$  computed using (8.3), is the same for each subdomain.

**12. Concluding remarks.** Inspired by the optimized Schwarz methods for the solution (and preconditioning) of partial differential equations on simple domains, we have presented an algebraic view of these methods. These new methods can be applied to banded and block-banded matrices, again as iterative methods, and as preconditioners. The new method can be seen as the application of several local Schur complements. When these Schur complements are computed, the method, and preconditioner is guaranteed to converge in two steps. The new formulation presents these Schur complements as solutions of nonlinear least squares problems. We can approximate the solution of these problems by solving closely related linear least squares problems. Further approximations are obtained by restricting the solution of these linear minimizations to matrices of certain sparsity pattern. Experiments show that these approximations, as well as the use of ILU to approximate the computation of the inverses in the Schur complement, produce very fast methods and preconditioners.

**Acknowledgements.** Part of this research was performed during visits of the third author to the Université de Genève, which was supported by the Fond National Suisse under grant FNS 200020-121561/1. The second and third authors were supported in part by the U.S. Department of Energy under grant DE-FG02-05ER25672.

#### REFERENCES

- [1] M. Benzi. Preconditioning techniques for large linear systems: A survey. *Journal of Computational Physics*, 182:418–477, 2002.
- [2] M. Benzi, A. Frommer, R. Nabben, and D. B. Szyld. Algebraic theory of multiplicative Schwarz methods. *Numerische Mathematik*, 89:605–639, 2001.
- [3] X.-C. Cai and M. Sarkis. A restricted additive Schwarz preconditioner for general sparse linear systems. *SIAM Journal on Scientific Computing*, 21:792–797, 1999.
- [4] P. Chevalier and F. Nataf. Symmetrized method with optimized second-order conditions for the Helmholtz equation. *Contemporary Mathematics*, 218:400–407, 1998.
- [5] J. Côté, M. J. Gander, L. Laayouni, and S. Loisel. Comparison of the Dirichlet-Neumann and Optimal Schwarz Method on the Sphere. In R. Kornhuber, R. Hoppe, J. Périaux, O. Pironneau, O. B. Widlund, and J. Xu (eds.), *Domain Decomposition Methods in Science and Engineering*, Lecture Notes in Computational Science and Engineering, vol. 40, Springer, Berlin, Heidelberg, 2004, pp. 235–242.
- [6] V. Dolean, S. Lanteri, and F. Nataf. Optimized interface conditions for domain decomposition methods in fluid dynamics. *International Journal on Numerical Methods in Fluids*, 40:1539–1550, 2002.
- [7] E. Efstathiou and M. J. Gander. Why Restricted Additive Schwarz Converges Faster than Additive Schwarz, *BIT Numerical Mathematics*, 43:945–959, 2003.
- [8] A. Frommer and D. B. Szyld. Weighted max norms, splittings, and overlapping additive Schwarz iterations. *Numerische Mathematik*, 83:259–278, 1999.
- [9] A. Frommer and D. B. Szyld. An algebraic convergence theory for restricted additive Schwarz methods using weighted max norms. *SIAM Journal on Numerical Analysis*, 39:463–479, 2001.
- [10] M. J. Gander. Optimized Schwarz Methods. *SIAM Journal on Numerical Analysis*, 44:699–731, 2006.
- [11] M. J. Gander, Schwarz Methods in the Course of Time, *Electronic Transactions on Numerical Analysis*, 31:228–255, 2008.
- [12] M. J. Gander, L. Halpern, and F. Nataf. Optimized Schwarz Methods. In T. Chan, T. Kako, H. Kawarada, O. Pironneau (eds.), *Proceedings of the Twelfth International Conference on Domain Decomposition*, DDM press, 2001, pp. 15–27.
- [13] A. George and J. W. Liu. *Computer Solution of Large Sparse Positive Definite Systems*. Prentice-Hall, Englewood Cliffs, New Jersey, 1981.
- [14] M. Griebel and P. Oswald. On the abstract theory of additive and multiplicative Schwarz algorithms. *Numerische Mathematik*, 70:163–180, 1995.

- [15] G.H. Golub and C.F. Van Loan. *Matrix Computations*. Third Edition, The Johns Hopkins University Press, Baltimore, Maryland 1996.
- [16] S. Loisel. Optimal and optimized domain decomposition methods on the sphere. Ph.D. Thesis, Department of Mathematics, McGill University, Montreal, 2005.
- [17] S. Loisel, J. Côté, M. J. Gander, L. Laayouni, A. Qaddouri. Optimized Domain Decomposition Methods for the Spherical Laplacian. Preprint, 2009.
- [18] F. Magoulès, F.-X. Roux, and S. Salmon. Optimal discrete transmission conditions for a non-overlapping domain decomposition method for the Helmholtz equation. *SIAM Journal on Scientific Computing*, 25:1497–1515, 2004.
- [19] F. Magoulès, F.-X. Roux, and L. Series. Algebraic way to derive absorbing boundary conditions for the Helmholtz equation. *Journal of Computational Acoustics*, 13:433–454, 2005.
- [20] F. Magoulès, F.-X. Roux, and L. Series. Algebraic approximation of Dirichlet-to-Neumann maps for the equations of linear elasticity. *Computer Methods in Applied Mechanics and Engineering*, 195:3742–3759, 2006.
- [21] T. P. A. Mathew. *Domain Decomposition Methods for the Numerical Solution of Partial Differential Equations*, Lecture Notes in Computational Science and Engineering, vol. 61, Springer, Berlin, Heidelberg, 2008.
- [22] R. Nabben and D. B. Szyld. Convergence theory of restricted multiplicative Schwarz methods. *SIAM Journal on Numerical Analysis*, 40:2318–2336, 2003.
- [23] R. Nabben and D. B. Szyld. Schwarz iterations for symmetric positive semidefinite problems. *SIAM Journal on Matrix Analysis and Applications*, 29:98–116, 2006.
- [24] A. Quarteroni and A. Valli. *Domain Decomposition Methods for Partial Differential Equations*. Oxford Science Publication, Clarendon Press, Oxford, 1999.
- [25] Y. Saad. *Iterative Methods for Sparse Linear Systems*. The PWS Publishing Company, Boston, 1996. Second edition, SIAM, Philadelphia, 2003.
- [26] V. Simoncini and D. B. Szyld. Recent Computational Developments in Krylov Subspace Methods for Linear Systems. *Numerical Linear Algebra with Applications*, 14:1–59, 2007.
- [27] B. F. Smith, P. E. Bjørstad, and W. D. Gropp. *Domain Decomposition: Parallel Multilevel Methods for Elliptic Partial Differential Equations*. Cambridge University Press, Cambridge, New York, Melbourne, 1996.
- [28] A. St-Cyr, M.J. Gander and S.J. Thomas, Optimized multiplicative, additive and restricted additive Schwarz preconditioning. *SIAM Journal on Scientific Computing*, 29:2402–2425, 2007.
- [29] W.-P. Tang. Generalized Schwarz splittings. *SIAM Journal on Scientific and Statistical Computing*, 13:573–595, 1992.
- [30] A. Toselli and O. Widlund. *Domain Decomposition Methods – Algorithms and Theory*. Springer Series in Computational Mathematics 34, Springer, Berlin, Heidelberg, 2005.
- [31] R. S. Varga. *Matrix Iterative Analysis*. Prentice-Hall, Englewood Cliffs, New Jersey, 1962. Second Edition, Springer Series in Computational Mathematics 27, Springer, Berlin, Heidelberg, New York, 2000.
- [32] D. M. Young. *Iterative Solution of Large Linear Systems*. Academic Press, New York, 1971.