# Finite element block-factorized preconditioners

Erik Bängtsson[*]    Maya Neytcheva[†]

**Abstract**

In this work we consider block-factorized preconditioners for the iterative solution of systems of linear algebraic equations arising from finite element discretizations of scalar and vector partial differential equations of elliptic type.

For the construction of the preconditioners we utilize a general two-level standard finite element framework and the corresponding block two-by-two form of the system matrix, induced by a splitting of the finite element spaces, referred to as *fine* and *coarse*, namely,

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{matrix} fine, \\ coarse. \end{matrix}$$

The matrix $A$ admits the exact factorization

$$A = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S_A \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix},$$

where $S_A = A_{22} - A_{21}A_{11}^{-1}A_{12}$ and $I_1$, $I_2$ are identity matrices of corresponding size. The particular form of preconditioners we analyze here is

$$M_B = \begin{bmatrix} B_{11} & 0 \\ A_{21} & S \end{bmatrix} \begin{bmatrix} I_1 & Z_{12} \\ 0 & I_2 \end{bmatrix},$$

where $S$ is assumed to be some available good quality approximation of the Schur complement matrix $S_A$.

We propose two methods to construct an efficient, sparse and computationally cheap approximation $B_{11}^{-1}$ of the inverse of the pivot block $A_{11}^{-1}$, required when solving systems with the block factorized preconditioner $M_B$. Furthermore, we propose an approximation $Z_{12}$ of the off-diagonal matrix block product $A_{11}^{-1}A_{12}$, which further reduces the computational complexity of the preconditioning step. All three approximations are based on element-by-element manipulations of local finite element matrices.

The approach is applicable for both selfadjoint and non-selfadjoint problems, in two as well as in three dimensions. We analyze in detail the 2D case and provide extensive numerical evidence for the efficiency of the proposed matrix approximations.

[*]Uppsala University, Box 337, 751 05 Uppsala, Sweden, email `Erik.Bangtsson@it.uu.se`
[†]Uppsala University, Box 337, 751 05 Uppsala, Sweden, email `Maya.Neytcheva@it.uu.se`

# 1 Introduction

Consider a nonsingular matrix $A$ of size $n$ and the task to solve a linear system with it, of the form

$$A\mathbf{x} = \mathbf{b}. \tag{1}$$

Assume that we split the degrees of freedom $\mathbf{x} \in V \subset \mathbb{R}^n$ into two nonintersecting classes (subspaces) $V^{(1)}$ of size $n_1$ and $V^{(2)}$ of size $n_2$, $V^{(1)} \cap V^{(2)} = \emptyset$ ($n = n_1 + n_2$), which in the sequel will be referred correspondingly to as *fine* and *coarse*,

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \begin{matrix} \}V^{(1)} \\ \}V^{(2)} = V \backslash V^{(1)} \end{matrix} \quad \begin{matrix} (fine) \\ (coarse). \end{matrix} \tag{2}$$

The above splitting induces in a natural way a 2-by-2 block splitting of the matrix $A$,

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{matrix} fine, \\ coarse. \end{matrix}$$

Denote by $S_A = A_{22} - A_{21}A_{11}^{-1}A_{12}$ the exact Schur complement of $A$.

One very much exploited form of preconditioning for $A$ is the two-by-two block-factorization

$$M_B = \begin{bmatrix} B_{11} & 0 \\ A_{21} & S \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}, \tag{3}$$

where $B_{11}$ and $S$ are some approximations of $A_{11}$ and $S_A$.

There is a vast amount of literature and research, related to preconditioners of type (3). Block-factorized preconditioners are used in a two-level as well as in a multilevel setting. By recursion on $S$, the block-factorization in (3) is straightforwardly extendable to the multilevel case and such approach has already been used in the Algebraic MultiLevel Iteration (AMLI) framework (cf. [9], and [10]), in the ILU framework (cf. [15]), in the approximate inverse context (cf. [6], [16], [1]), to name a few typical preconditioning strategies, based on some block-factorized form of $A$. The preconditioner is also applicable for matrices of saddle-point form, surveyed for example in [13] and [8].

Related to the construction and the aimed properties of $M_B$, several issues have to be considered.

(A) *Computational cost of one preconditioning step*, namely the cost to solve one system of equations with $M_B$.

As is clear from the structure of the preconditioner, to solve $M_B \mathbf{y} = \mathbf{d}$, the following steps are required:

| | | | |
|---|---|---|---|
| (S1) | $B_{11}\mathbf{z_1} = \mathbf{d_1}$ | (S3) | $B_{11}\mathbf{w}_1 = A_{12}\mathbf{y}_2$ |
| (S2) | $S\mathbf{y_2} = \mathbf{d}_2 - A_{21}\mathbf{z}_1$ | (S4) | $\mathbf{y}_1 = \mathbf{z}_1 - \mathbf{w}_1$ |

The prevailing part of the computational effort lies in steps S1 to S3, where we have to solve two systems with $B_{11}$ and one system with $S$. (Clearly, if we construct an approximation $B_{11}^{-1}$ to $A_{11}^{-1}$ instead of approximating $A_{11}$, then, in steps S1 and S3, the solution of the linear system could be replaced by a matrix-vector multiplication.)

(B) *How to define the splitting (2)*, i.e., the subspaces $V^{(1)}$ and $V^{(2)}$ so that some desirable properties of the preconditioner can be achieved?

The answer to this question depends on the framework, in which the preconditioner is constructed.

One criterion could be to first reorder $A$ so that a (nearly) diagonal main pivot block is obtained. For $A_{11}$ being diagonal, the explicit computation of $S_A$ is cheap. The preconditioner can be extended to its multilevel version after recursively repeating the procedure to the so-obtained Schur complement. This framework is first described in [15] and later used and elaborated in [25], [24] and others.

Another criterion to measure the quality of the two-by-two splitting, applicable for symmetric positive definite matrices only, is to compute the so-called Cauchy-Bunyakowski-Schwarz (CBS) constant $\gamma$, associated with the splitting. The CBS constant can be defined in different ways, for instance as the square root of the following spectral radius

$$\gamma^2 = \rho \left( A_{22}^{-1} A_{21} A_{11}^{-1} A_{12} \right).$$

For $V^{(1)} \cap V^{(2)} = \emptyset$ the constant $\gamma$ is always strictly less than unity and measures the relative strength of the off-diagonal couplings between $A_{11}$ and $A_{22}$. Clearly, $\gamma = 0$ is obtained for $A_{12} = 0$, which is the best possible splitting, if of course achievable. On the other hand, if $A_{12}$ induces a strong off-diagonal coupling, then $\gamma$ can become arbitrarily close to 1. The latter is undesirable since the condition number estimates of $M_B^{-1} A$ involve $\gamma$ and $\gamma \approx 1$ indicates possible unbounded increase of the latter. We refer for example to [1], Chapter 9, as a rich source with more details on related condition number estimates.

A profitable framework to obtain efficient two-by-two block matrix splittings and values of $\gamma$ bounded away from unity is ensured within the two-level Finite Element (FE) framework and the usage of the Hierarchical Basis Functions (HBF) method. There, the fine and the course degrees of freedom are associated with two regular mesh refinement levels. In the latter context $\gamma$ is defined as

$$\gamma = \sup_{\substack{\mathbf{x}_1 \in V^{(1)} \\ \mathbf{x}_2 \in V^{(2)}}} \frac{\mathbf{x}_1^T A_{12} \mathbf{x}_2}{(\mathbf{x}_1^T A_{11} \mathbf{x}_1)^{1/2} (\mathbf{x}_2^T A_{22} \mathbf{x}_2)^{1/2}} \tag{4}$$

As shown in [7, 17], $\gamma$ can be estimated locally and is proven to be independent of mesh parameters, number of levels of refinement, jumps of problem coefficients (as long as they are aligned with the coarse mesh), geometry of the domain (for instance, reentrant corners), as well as mesh and problem anisotropies. The latter implies that as long as the condition number is expressed as a function of $\gamma$, it is independent of the influence of the same parameters as the CBS constant itself.

It is considered as a slight drawback of the HBF method that in order to profit from the very good condition number estimates, available in this case, one must use the corresponding

HBF matrix $\widehat{A}$, which is less sparse than $A$ (it is related to it via a congruence transformation of the type $\widehat{A} = J^T A J$).

We note that, unfortunately, there is no analog to $\gamma$ established in the case of nonsymmetric matrices, at least up to the knowledge of the authors.

(C) *How to approximate the Schur complement $S_A$ and the pivot block $A_{11}$?* A relevant related question is whether the approximations $S$ and $B_{11}$ have to be associated to each other in order to ensure good spectral properties of the preconditioned matrix $M_B^{-1} A$. In [22] it is shown that for certain approximate inverse techniques, used to construct $B_{11}^{-1}$ it is necessary to relate the approximation $S$ to that of the pivot block.

This work discusses in more detail issue (C) within the following setting. We do not make use of HBF, however we utilize the two-level standard basis functions FE framework. The question how to approximate the Schur complement falls out of the scope of this paper. We assume that we possess a good approximation $S$ of $S_A$, which inherits the properties of $S_A$, such as positive definiteness, symmetricity or nonsymmetricity. The particular choice of $S$ is stated in Section 3.

The paper is organized as follows. In Section 2 we consider a novel idea when constructing the full block-factorized preconditioner. Section 3 discusses two approximations of the pivot block, based of the FE framework. In section 4 we briefly discuss multilevel extensions of the proposed techniques. The properties of the suggested preconditioner are illustrated by numerical experiments, presented in Section 5.

## 2 More on full block-factorized preconditioners

The approximate block-factorization is a general framework suited for both selfadjoint and non-selfadjoint problems. Our interest is focused on the latter class and all considerations are done for nonsymmetric matrices mostly. For the purpose of having a convenient framework to illustrate the ideas, which by no means limits their generality, we adopt as a reference a second order elliptic problem of the form $\mathcal{L}u = f$ in $\Omega \subset \mathbb{R}^2$. The operator $\mathcal{L}$ can be of the form $\mathcal{L} = \nabla \cdot (K \nabla)$ or $\mathcal{L} = \nabla \cdot (\nabla) + \mathbf{b} \cdot \nabla$, where the coefficient matrix $K$ is diagonal but could vary in space, and leads to symmetric problems with anisotropy or discontinuous coefficients. The vector $\mathbf{b}$ describes some convection and the discrete analog of the second operator is in general nonsymmetric.

We assume that the problem is discretized on some triangular or quadrilateral mesh, referred to as *coarse*, which is then subjected to one regular refinement by subdividing the edges into $m$ equal parts. In this work the results are derived for $m = 2$ and some comments are included for $m > 2$.

Let the number of finite elements on the coarse mesh (referred to as macroelements) be $M$ and a fine-coarse splitting of the degrees of freedom be imposed on macroelement level also.

Consider two forms of full block-factorizations of the given matrix $A$,

$$A = L_B^A U_B^A = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S_A \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1} A_{12} \\ 0 & I_2 \end{bmatrix} \tag{5}$$

4

and

$$A = L_S^A D_S^A U_S^A = \begin{bmatrix} I_1 & 0 \\ A_{21}A_{11}^{-1} & I_2 \end{bmatrix} \begin{bmatrix} A_{11} & 0 \\ 0 & S_A \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}. \tag{6}$$

We start by observing that in the upper-triangular block factor $U_B^A$ in the factorization (3) we do not need to approximate $A_{11}^{-1}$ itself but rather the product $A_{11}^{-1}A_{12}$. Similar argument holds for the off-diagonal blocks in the factors $L_S^A$ and $U_S^A$ in (6). Therefore, we consider the following two full block-factorized preconditioners,

$$M_B = L_B U_B = \begin{bmatrix} B_{11} & 0 \\ A_{21} & S \end{bmatrix} \begin{bmatrix} I_1 & Z_{12} \\ 0 & I_2 \end{bmatrix} \tag{7}$$

and

$$M_S = L_S D_S U_S = \begin{bmatrix} I_1 & 0 \\ Z_{21} & I_2 \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ 0 & S \end{bmatrix} \begin{bmatrix} I_1 & Z_{12} \\ 0 & I_2 \end{bmatrix} \tag{8}$$

The prerequisite is to allow independent approximations for $A_{11}^{-1}$, $S_A$, and for $A_{11}^{-1}A_{12}$, respectively $A_{21}A_{11}^{-1}$, as whole blocks. We show in the sequel that the FE framework and the availability of two (consecutive) mesh refinements provide a possibility to construct an approximation $B_{11}^{-1}$ to $A_{11}^{-1}$ and $Z_{12}$, $Z_{21}$ to $A_{11}^{-1}A_{12}$ and $A_{21}A_{11}^{-1}$ respectively, on low computational cost, which also leads to efficient two-level preconditioners to $A$ of the form (7) and (8). In Section 4 we also indicate that the two-level framework can be extended to the multilevel case.

We analyze now the spectra of the preconditioned matrices $M_B^{-1}A$ and $M_S^{-1}A$. As is well known, the aim of preconditioning is to cluster the corresponding eigenvalues around unity (in general in the complex plain), which will ensure fast convergence of the iterative solution method. The convergence of a preconditioned iterative method can also be estimated when the spectrum of the preconditioned matrix has a positive real part and is contained in an ellipse or in two well-separated ellipse-like ovals, based on two disjoint real intervals. It is known (cf., for instance [1]) that in such cases the aim is to obtain narrow ellipses, which then ensure a small asymptotic convergence factor and thus, fast convergence of the iterative method.

**Case 1**: Consider $M_B^{-1}A$.

For any choice of the approximations $B_{11}$, $S$ and $Z_{12}$, this preconditioner is nonsymmetric even if $A$ is symmetric.

To obtain a general idea how the errors in the approximations of the three blocks ($B_{11}$, $S$ and $Z_{12}$) influence the quality of the preconditioner, we introduce an error matrix $\Delta_B$, such that $M_B = A + \Delta_B$ which is found to be as follows

$$\Delta_B = \begin{bmatrix} B_{11} - A_{11} & B_{11}Z_{12} - A_{12} \\ 0 & S + A_{21}Z_{12} - A_{22} \end{bmatrix}.$$

5

For the matrix $\widetilde{\Delta}_B = M_B^{-1}\Delta_B$ we obtain $\widetilde{\Delta}_B = \begin{bmatrix} \widetilde{\Delta}_{11} & \widetilde{\Delta}_{12} \\ \widetilde{\Delta}_{21} & \widetilde{\Delta}_{22} \end{bmatrix}$ with blocks

$$
\begin{aligned}
\widetilde{\Delta}_{11} &= (I_1 + Z_{12}S^{-1}A_{21})(I_1 - B_{11}^{-1}A_{11}) \\
\widetilde{\Delta}_{12} &= (I_1 + Z_{12}S^{-1}A_{21})(I_1 - B_{11}^{-1}A_{11})A_{11}A_{12} + A_{11}^{-1}A_{12}(S^{-1}S_A - I_2) + \\
&\quad (Z_{12} - A_{11}^{-1}A_{12})S^{-1}S_A \\
\widetilde{\Delta}_{21} &= S^{-1}A_{21}(B_{11}^{-1}A_{11} - I_1) \\
\widetilde{\Delta}_{22} &= (I_2 - S^{-1}S_A) + S^{-1}A_{21}(B_{11}^{-1}A_{11} - I_1)A_{11}^{-1}A_{12}
\end{aligned}
\tag{9}
$$

We see from (9) that the approximation $Z_{12}$ plays a role only in two of the blocks of $\widetilde{\Delta}$, while how well $B_{11}$ approximates $A_{11}$ (respectively how $B_{11}^{-1}$ approximates $A_{11}^{-1}$) influences all four blocks of $\widetilde{\Delta}$.

We estimate the eigenvalues of $A\mathbf{v} = \lambda M_B\mathbf{v}$. Using (5) and (7), and the general assumption that $A$ and $M_B$ are nonsingular, we see that we can analyze two equivalent eigenvalue problems instead:

$$
Q_1\mathbf{w} \equiv L_B^{-1}L_B^A U_B^A U_B^{-1}\mathbf{w} = \lambda\mathbf{w} \tag{10}
$$

$$
Q_2\mathbf{z} \equiv L_B^{A^{-1}}L_B U_B U_B^{A^{-1}}\mathbf{z} = \frac{1}{\lambda}\mathbf{z} \tag{11}
$$

where $\mathbf{z} = U_B^A\mathbf{v}$ and $\mathbf{w} = U_B\mathbf{v}$.

Computation reveals that

$$
Q_1 = \begin{bmatrix} I_1 - E_{11} & -(I_1 - E_{11})E_{12} \\ -(I_2 - E_{22})S_A^{-1}A_{21}E_{11} & (I_2 - E_{22})(I_2 - S_A^{-1}A_{11}E_{11}E_{12}) \end{bmatrix}
$$

$$
Q_2 = \begin{bmatrix} I_1 - F_{11} & (I_1 - F_{11})E_{12} \\ S_A^{-1}A_{21}F_{11} & I_2 - E_{22} + S_A^{-1}A_{12}F_{11}E_{12} \end{bmatrix}
$$

where $E_{11} = I_1 - B_{11}^{-1}A_{11}$, $F_{11} = I_1 - A_{11}^{-1}B_{11}$, $E_{22} = I_2 - S^{-1}S_A$, $F_{22} = I_2 - S_A^{-1}S$ and $E_{12} = Z_{12} - A_{11}^{-1}A_{12}$.

We find then that for appropriate splitting of the vectors $\mathbf{w}$ and $\mathbf{z}$,

$$
\begin{aligned}
\lambda\mathbf{w}^*\mathbf{w} &= \mathbf{w}_1^*\mathbf{w}_1 + \mathbf{w}_2^*\mathbf{w}_2 - \mathbf{w}_1^*E_{11}\mathbf{w}_1 + \mathbf{w}_2^*E_{22}\mathbf{w}_2 \\
&\quad -\mathbf{w}_2^*(I_2 - E_{22})S_A^{-1}A_{21}E_{11}E_{12}\mathbf{w}_2 \\
&\quad -\mathbf{w}_1^*(I_1 - E_{11})E_{12}\mathbf{w}_2 - \mathbf{w}_2^*(I_2 - E_{22})S_A^{-1}A_{21}E_{11}\mathbf{w}_1
\end{aligned}
\tag{12}
$$

and

$$
\begin{aligned}
\frac{1}{\lambda}\mathbf{z}^*\mathbf{z} &= \mathbf{z}_1^*\mathbf{z}_1 + \mathbf{z}_2^*\mathbf{z}_2 - \mathbf{z}_1^*F_{11}\mathbf{z}_1 - \mathbf{z}_2^*F_{22}\mathbf{z}_2 + \mathbf{z}_2^*S_A^{-1}A_{21}F_{11}E_{12}\mathbf{z}_2 \\
&\quad +\mathbf{z}_1^*(I_1 - F_{11})E_{12}\mathbf{z}_2 + \mathbf{z}_2^*S_A^{-1}A_{21}F_{11}\mathbf{z}_1.
\end{aligned}
\tag{13}
$$

Above, $^*$ denotes a conjugate transpose vector. From (12), (13), and using the inequality $2|ab| \leq a^2 + b^2$, we obtain the following bound for the eigenvalues of $M_B^{-1}A$:

$$
\begin{aligned}
|\lambda| &\leq 1 + \max\left\{\|E_{11}\|, \|E_{22}\| + \|I_2 - E_{22}\| \|S_A^{-1}A_{21}E_{11}\| \|E_{12}\|\right\} \\
&\quad + \tfrac{1}{2}\left(\|I_1 - E_{11}\| \|E_{12}\| + \|I_2 - E_{22}\| \|S_A^{-1}A_{21}E_{11}\|\right) \\
|\lambda| &\geq 1/R_D
\end{aligned}
\tag{14}
$$

with

$$
\begin{aligned}
R_D &= 1 + \max\left\{\|F_{11}\|, \|F_{22}\| + \|S_A^{-1}A_{21}F_{11}\| \|E_{12}\|\right\} \\
&\quad + \tfrac{1}{2}\left(\|I_1 - F_{11}\| \|E_{12}\| + \|S_A^{-1}A_{21}F_{11}\|\right).
\end{aligned}
\tag{15}
$$

The narrow ellipse will be guaranteed by ensuring that the skew-symmetric part $1/2(M_B^{-1}A - A^T M_B^{-T})$ of the preconditioned matrix has small eigenvalues. We do not analyze this issue any further. Instead, we provide numerical evidence that the above preconditioner, constructed using the framework described in Section 3 does cluster the spectrum in a narrow ellipse. We also illustrate that for the considered classes of problems, the eigenvalues are well clustered around unity with very few outliers, which, even though very large, can easily be eliminated by an iterative solution method.

**Case 2**: Consider $M_S^{-1}A$.

Utilizing now (6) and (8), we note that the spectrum of the generalized eigenvalue problem $A\mathbf{v} = \lambda M_S \mathbf{v}$ can be analyzed via the following two equivalent eigenproblems:

$$
T_1 \mathbf{w} \equiv D_S^{-1}L_S{}^{-1}L_S^A D_S^A U_S^A U_S{}^{-1}\mathbf{w} = \lambda \mathbf{w}
\tag{16}
$$

$$
T_2 \mathbf{z} \equiv D_S^{A^{-1}}L_S^{A^{-1}}L_S D_S U_S U_S^{A^{-1}}\mathbf{z} = \frac{1}{\lambda}\mathbf{z}
\tag{17}
$$

where $\mathbf{z} = U_S^A \mathbf{v}$ and $\mathbf{w} = U_S \mathbf{v}$. Computing $T_1$ and $T_2$ explicitly, we find

$$
T_1 = \begin{bmatrix} I_1 - E_{11} & -(I_1 - E_{11})E_{12} \\ (I_2 - E_{22})S_A^{-1}E_{21}A_{11} & (I_2 - E_{22})(I_2 + S_A^{-1}E_{21}A_{11}E_{12}) \end{bmatrix}
$$

and

$$
T_2 = \begin{bmatrix} I_1 - F_{11} & (I_1 - F_{11})E_{12} \\ S_A^{-1}E_{21}A_{11}(I_1 - F_{11}) & (I_2 - F_{22}) + S_A^{-1}E_{21}A_{11}(I_1 - F_{11})E_{12}, \end{bmatrix}
$$

where $E_{ij}$ and $F_{1j}$, $i, j = 1, 2$ are defined as in Case 1. For appropriate splitting of the vectors $\mathbf{z}$ and $\mathbf{w}$ we find that

$$
\begin{aligned}
\lambda \mathbf{w}^*\mathbf{w} &= \mathbf{w}_1^*\mathbf{w}_1 + \mathbf{w}_2^*\mathbf{w}_2 - \mathbf{w}_1^*E_{11}\mathbf{w}_1 + \mathbf{w}_2^*E_{22}\mathbf{w}_2 \\
&\quad - \mathbf{w}_2^*E_{22}S_A^{-1}E_{21}A_{11}E_{12}\mathbf{w}_2 \\
&\quad - \mathbf{w}_1^*(I_1 - E_{11})E_{12}\mathbf{w}_2 - \mathbf{w}_2^*(I_1 - E_{22})S_A^{-1}E_{21}A_{11}\mathbf{w}_1
\end{aligned}
\tag{18}
$$

and

$$
\begin{aligned}
\tfrac{1}{\lambda}\mathbf{z}^*\mathbf{z} &= \mathbf{z}_1^*\mathbf{z}_1 + \mathbf{z}_2^*\mathbf{z}_2 - \mathbf{z}_1^*F_{11}\mathbf{z}_1 + \mathbf{z}_2^*(I_2 - F_{22})\mathbf{z}_2 + \mathbf{z}_2^*S_A^{-1}E_{21}A_{11}(I_1 - F_{11})E_{12}\mathbf{z}_2 \\
&\quad + \mathbf{z}_1^*(I_1 - F_{11})E_{12}\mathbf{z}_2 + \mathbf{z}_2^*S_A^{-1}E_{21}A_{11}(I_1 - F_{11})\mathbf{z}_1.
\end{aligned}
\tag{19}
$$

Using (18), (19), the bounds for the eigenvalues of $M_S^{-1}A$ are found to be:

$$
\begin{aligned}
|\lambda| &\leq 1 + \max\left\{\|E_{11}\|, \|E_{22}\| + \|E_{22}\|\,\|S_A^{-1}E_{21}A_{11}\|\,\|E_{12}\|\right\} \\
&\quad + \tfrac{1}{2}\left(\|I_1 - E_{11}\|\,\|E_{12}\| + \|I_2 - E_{22}\|\,\|S_A^{-1}E_{21}A_{11}\|\right) \\
|\lambda| &\geq \frac{1}{R_S}
\end{aligned}
\tag{20}
$$

with

$$
\begin{aligned}
R_S &= 1 + \max\left\{\|F_{11}\|, \|F_{22}\| + \|S_A^{-1}E_{21}A_{11}\|\,\|I_1 - F_1\|\,\|E_{12}\|\right\} \\
&\quad + \tfrac{1}{2}\left(\|I_1 - F_{11}\|\,\|E_{12}\| + \|S_A^{-1}E_{21}A_{11}\|\,\|I_1 - F_1\|\right).
\end{aligned}
\tag{21}
$$

As expected, and as is also seen from (20) and (21), in this case the quality of the approximations of both $Z_{12}$ and $Z_{21}$ must be very good, in order to ensure the aimed clustering of the eigenvalues $\lambda$ around unity.

For this form of the preconditioner, analysis, not included here, reveals that if $A$ is spd, the spectrum bounds can be refined. However, since estimates (20) – (21) do not bring new insight in the considered matter, we do not provide numerical illustrations for the behavior of $M_S$. We note that this form, combined with a different (symmetric) scaling used when constructing $B_{11}^{-1}$ and $Z_{12}, Z_{21}$ may turn out to be the method of choice for selfadjoint problems.

# 3   Finite element-based approximations of the pivot block

We choose to approximate directly $A_{11}^{-1}$ and to this end consider the $k$th macroelement ($k = 1, \cdots, M$) together with the corresponding macroelement stiffness matrix, written in two-by-two block form

$$
A_k = \begin{bmatrix} A_{11,k} & A_{12,k} \\ A_{21,k} & A_{22,k} \end{bmatrix} \begin{array}{l} \}V^{(1)} \\ \}V^{(2)} = V \backslash V^{(1)} \end{array} \begin{array}{l} (fine) \\ (coarse) \end{array}.
\tag{22}
$$

The approximation $B_{11}^{-1}$ is constructed in an element-by-element (EBE) fashion, as follows

$$
B_{11}^{-1} = \sum_{k=1}^{M} R_k^T A_{11,k}^{-1} R_k,
\tag{23}
$$

where $R_k$ denotes the Boolean matrix representing the restriction to the degrees of freedom of the $k$th macroelement. That is, $B_{11}^{-1}$ is the assembled sum of all local exact inverses of the pivot blocks in the macroelement stiffness matrices. A similar idea is first used in [18] in the context of the construction of an AMLI method and is found unsatisfactory, compared to a standard incomplete factorization of $A_{11}$. In [5] it is shown that $B_{11}^{-1}$ and $A_{11}^{-1}$ are spectrally equivalent, namely that for some $0 < \alpha_1 < \alpha_2$ there holds

$$
\alpha_1 A_{11}^{-1} \leq B_{11}^{-1} \leq \alpha_2 A_{11}^{-1}.
\tag{24}
$$

There, it is also pointed out, that the spectral equivalence constants $\alpha_1$ and $\alpha_2$ depend on the ratio $\varkappa_1/\varkappa_2$, where

$$
0 < \varkappa_1 \leq \lambda_{min}(A_{11,k}) \leq \cdots \leq \lambda_{max}(A_{11,k}) \leq \varkappa_2, \text{ for all } k = 1, \cdots M.
$$

8

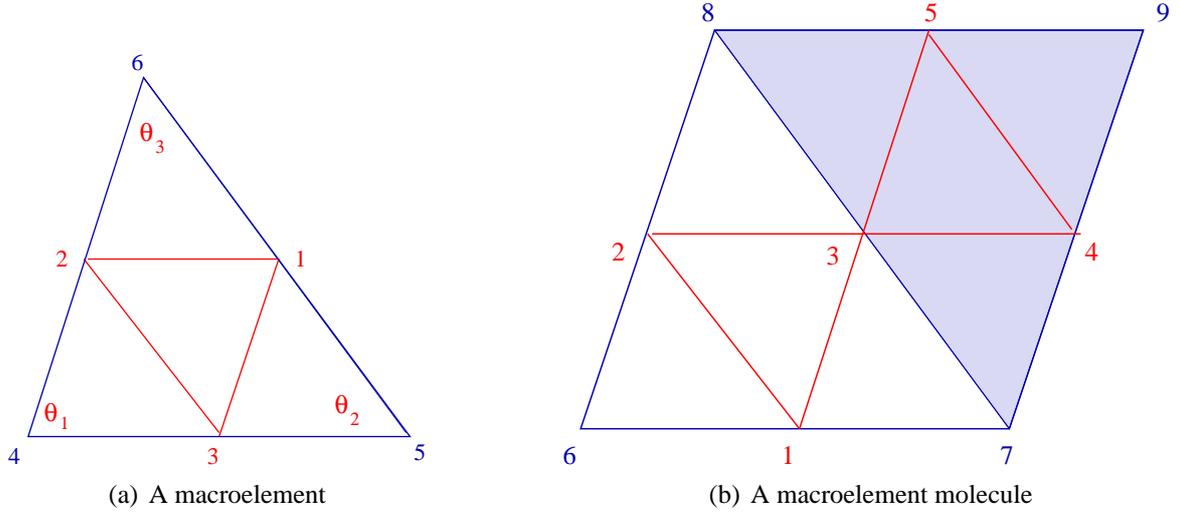(a) A macroelement          (b) A macroelement molecule

Figure 1: Macroelements

Thus, the spectral equivalence constants are independent on the mesh parameter $h$ but they are robust neither with respect to problem and mesh-anisotropies, nor to jumps in the problem coefficients. Furthermore, it is illustrated in [5] that the condition number increases faster than quadratically with $m$, which denotes the number of divisions that are performed on the coarse mesh edges to obtain the fine mesh. In Figure 1, $m = 2$.

Let us now take a closer look at $B_{11}^{-1}$. We observe that

$$
\begin{aligned}
A_{11} B_{11}^{-1} &= \sum_{l=1}^{M} R_l^T A_{11,l} R_l \sum_{k=1}^{M} R_k^T A_{11,k}^{-1} R_k \\
&= \sum_{k=1}^{M} R_k^T A_{11,k} A_{11,k}^{-1} R_k + \sum_{l=1}^{M} \sum_{\substack{k=1 \\ k \neq l}}^{M} R_l^T A_{11,l} R_l R_k^T A_{11,k}^{-1} R_k \\
&= D_{11} + \sum_{l=1}^{M} \sum_{\substack{k=1 \\ k \neq l}}^{M} R_l^T P_{l,k} R_k = D_{11} + W_{11}
\end{aligned}
\tag{25}
$$

where $D_{11}$ is a diagonal matrix with entries equal to either 1 or 2. The value 2 corresponds to nodes which are midpoints of interior edges, and thus, belong to two (and only two) macroelements. The latter hints that the element inverses $A_{11,k}^{-1}$ can be scaled properly in such a way that the corresponding sum $D_{11}$ becomes the identity matrix. From now on we assume that such a scaling is imposed on each individual $A_{11,k}^{-1}$. The exact form of the scaling is explained in Section 3.1.

Next, we note that for $l \neq k$ the product matrices $W_{l,k} = R_l^T A_{11,l} R_l R_k^T A_{11,k}^{-1} R_k$ are very

sparse (for example, for triangular meshes and $m = 2$ they have at most nine nonzero entries). In the sum $W_{11}$, $W_{l,k}$ are nonzero only for such pairs $(k,l)$ which denote neighboring macroelements, that is, where two matrices $R_l^T A_{11,l} R_l$ and $R_k^T A_{11,k}^{-1} R_k$ are intersecting only via the nodes which share a common edge in the set of all macroelements. This means that we are able to estimate the norm of $W_{11}$ based on local arguments. Furthermore, the estimate depends on problem parameters such as anisotropies and coefficient jumps, but not on the number of the macroelements, which again implies mesh-independence of the condition number of $A_{11} B_{11}^{-1}$.

In the analysis below we restrict ourselves to 2D, $m = 2$, arbitrary triangular mesh and piecewise linear FE basis functions. We denote by $q \ll M$ the number of nonzero terms $W_{l,k}$, and see that $q$ equals to the number of interior points belonging to the fine level that are shared by macroelement $l$ and its neighboring macroelements. Expression (25) shows that $A_{11} B_{11}^{-1}$ can be expressed as diagonal matrix, which, in the 2D case, is perturbed by a sum of $q$ rank-one matrices (in the literature referred to as *dyads*).

Let $rank(W_{11}) = r$. The general theory for such matrices states that $r \leq q$ and the eigenvalues of $I + W_{11}$ are exactly given by $1 + \lambda_j$, there $\lambda_j$ are generic eigenvalues of $W_{11}$. (cf., for instance, [14]). Therefore, $A_{11} B_{11}^{-1}$ has $n - r$ eigenvalues equal to one, where $r = rank(W_{11}) << n$ and $r$ eigenvalues of the above described form $1 + \lambda_j$, $j = 1, \cdots r$.

Let denote by $\mathcal{S}(W_{11})$ the nonzero spectrum of $W_{11}$ and consider the numerical range $\mathcal{W}(W_{11})$. We recall that the numerical range (or the field of values) of an arbitrary (real) matrix $A$ is defined as

$$\mathcal{W}(A) = \{\mathbf{x}^T A \mathbf{x}, \mathbf{x} \in \mathbb{R}^n, \mathbf{x}^T \mathbf{x} = 1\}.$$

As is well known, the numerical radius $r(A)$ of a general matrix $A$ is defined as $r(A) = \sup\{\mathbf{x}^T A \mathbf{x}, \mathbf{x} \in \mathbb{R}^n, \mathbf{x}^T \mathbf{x} = 1\}$ and the following properties hold true (see, for instance Theorem 4.5 in [1]):

1. $\mathcal{S}(A) \subset \mathcal{W}(A)$,

2. $\rho(A) = \|A\|_2 \leq r(A)$,

3. $r(A) \leq \|A\| \leq 2r(A)$, where $\|A\| = \sqrt{\lambda_{max}(A^T A)}$ is the spectral norm of $A$.

Then, for appropriate vectors $\mathbf{x}$ and for $k \neq l$ we have

$$
\begin{aligned}
\mathbf{x}^T \left( \sum_{k,l=1}^{q} W_{l,k} \right) \mathbf{x} &= \sum_{k,l=1}^{q} \left( \mathbf{x}^T R_l^T \right) P_{l,k} \left( R_k \mathbf{x} \right) \\
&\leq \sigma_{max}(P_{l,k}) \sum_{k,l=1}^{q} \left( \mathbf{x}^T R_l^T \right) \left( R_k \mathbf{x} \right) \leq 3\sigma_{max}(P_{l,k})
\end{aligned}
\tag{26}
$$

The factor $3$ reflects the fact that for 2D triangular elements the global index set for each macroelement pivot matrix intersects that of at most three other pivots from neighboring macroelements.

We analyze the properties of $W_{11}$ for the following simple model problem:

**Problem 1** Consider the scalar Poisson equation, discretized on an arbitrary triangular mesh. As pointed out, for example in [3], it is equivalent to consider either a scalar anisotropic Laplace

operator on isosceles triangles, or an isotropic Laplace operator on an arbitrary mesh. Here we consider the former case on meshes of the following three types:

$\mathcal{T}_1$: Right-angled isosceles triangles.

$\mathcal{T}_2$: Triangles with one large and two small angles. The triangle is "flat" and an example is shown in Figure 3(a).

$\mathcal{T}_3$: Triangles with two large and one small angle. The triangle is "sharp" as is depicted in Figure 3(b).

To analyze further the properties of the matrices $W_{l,k}$ also in the case of discontinuous coefficients, it suffices to consider two neighboring macroelements, as indicated in Figure 1(b). Assume that there is a discontinuity of size $s$ aligned with the two macroelements which we incorporate by multiplying the element coefficient matrix in the non-shaded area by $s$. Without loss of generality we can assume that $s > 1$ since the worst case scenario is either when $s = 1$ in macroelement 1 and $s \ll 1$ in macroelement 2 or when $s \gg 1$ in macroelement 1 and $s = 1$ in macroelement 2.

Now, let us introduce the element stiffness matrix for an arbitrary shaped triangle (cf. [3]),

$$A_k = \frac{1}{2} \begin{bmatrix} b+c & -c & -b \\ -c & a+c & -a \\ -b & -a & a+b \end{bmatrix},$$

where $a = \cot\theta_1$, $b = \cot\theta_2$, $c = \cot\theta_3$ and $\theta_1 \geq \theta_2 \geq \theta_3$ are the angles in the triangle, as shown in Figure 1(a). The assembled pivot matrix for one macroelement has the form

$$A_{11,k} = \begin{bmatrix} a+b+c & -c & -b \\ -c & a+b+c & -a \\ -b & -a & a+b+c \end{bmatrix},$$

and its exact inverse is then

$$A_{11,k}^{-1} = \frac{1}{2} \begin{bmatrix} \frac{b+2a+c}{(a+c)(a+b)} & (a+b)^{-1} & (a+c)^{-1} \\ (a+b)^{-1} & \frac{a+c+2b}{(a+b)(b+c)} & (b+c)^{-1} \\ (a+c)^{-1} & (b+c)^{-1} & \frac{a+2c+b}{(a+c)(b+c)} \end{bmatrix}.$$

For completeness we include the matrices $R_1$ and $R_2$, corresponding to the two macroelements in Figure 1(b)

$$R_1 = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

11

and

$$R_2 = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Then $R_1 R_2^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$. With the help of the above matrices, the product $P_{12} = s A_{11,1} R_1 R_2^T A_{11,2}^{-1}$ is found to be

$$P_{12} = \frac{s}{2} \begin{bmatrix} \frac{(a+b+c)(b+2\,a+c)}{(a+b)(a+c)} & \frac{(a+b+c)}{a+b} & \frac{(a+b+c)}{a+c} \\ -\frac{c\,(b+2\,a+c)}{(a+b)(a+c)} & -\frac{c}{a+b} & -\frac{c}{a+c} \\ -\frac{b\,(b+2\,a+c)}{(a+b)(a+c)} & -\frac{b}{a+b} & -\frac{b}{a+c} \end{bmatrix}. \tag{27}$$

Simple computations reveal that $P_{12}$ can be represented as a product of two vectors and is thus a rank-one matrix. That is,

$$P_{12} = \frac{s}{2} \mathbf{v} \mathbf{w}^T,$$

where

$$\mathbf{v} = \begin{bmatrix} a+b+c \\ -c \\ -b \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} (a+b)^{-1} + (a+c)^{-1} \\ (a+b)^{-1} \\ (a+c)^{-1} \end{bmatrix},$$

and the only singular value of $P_{12}$ is found to be equal to

$$\sigma = \frac{s\sqrt{2}}{2} \frac{\sqrt{f(a,b,c)g(a,b,c)}}{(a+c)(a+b)} \tag{28}$$

where

$$f(a,b,c) = 2(a+b)^2 + (b+c)^2 + 2ac$$

and

$$g(a,b,c) = (a+b)^2 + (b+c)^2 + (ab+ac+cb).$$

Relation (28) shows explicitly how the condition number of $A_{11}B_{11}^{-1}$ depends on the jump in the coefficients, and on the anisotropy in the problem induced by the stretching of the mesh.

The value of $\sigma$ in Equation (28) is evaluated using the symbolic computations software `Maple` ([19]) for $\mathcal{T}_i, i = 1, 2, 3$.

| Triangulation | $a, b, c$ | $\sigma$ |
|:---:|:---:|:---:|
| $\mathcal{T}_1$ | $a = 0$ <br> $b = 1$ <br> $c = 1$ | $1/2\sqrt{36}\,s = 3\,s$ |
| $\mathcal{T}_2$ | $a = -7.9158$ <br> $b = 15.8945$ <br> $c = 15.8945$ | $3.559123937\sqrt{2}\,s = 5.0334\,s$ |
| $\mathcal{T}_3$ | $a = -0.0315$ <br> $b = 0.0629$ <br> $c = 31.8520$ | $717.9941955\sqrt{2}\,s = 1015.397\,s$ |

Note, that so far no scaling on $A_{11,k}^{-1}$ has been imposed.

## 3.1   Scaled macroelement pivot inverses (EBES)

We now consider a scaled version of the EBE approximation for the pivot block, referred to as EBES, and denote it by $\widetilde{B}_{11}^{-1}$. The EBES approximation is constructed as in (23) with the only difference that $A_{11,k}^{-1}$ is scaled from the right by a diagonal scaling matrix $D_k$, namely,

$$\widetilde{B}_{11}^{-1} = \sum_{k=1}^{M} R_k^T A_{11,k}^{-1} D_k R_k. \tag{29}$$

The entries of $D_k$ are $1$ and $1/2$ (in 2D), where the entries $1/2$ correspond to the columns of $A_{11,k}^{-1}$ associated with nodes, shared by two neighboring macroelements. Then, for $A_{11}\widetilde{B}_{11}^{-1}$ we obtain ($k \neq l$)

$$
\begin{aligned}
A_{11}\widetilde{B}_{11}^{-1} &= I_1 + \sum_{l,k=1}^{q} R_l^T A_{11,l} R_l R_k^T A_{11,k}^{-1} D_k R_k \\
&= I_1 + \sum_{l,k=1}^{q} R_l^T \widetilde{P}_{l,k} R_k = I_1 + \sum_{l,k=1}^{q} \widetilde{W}_{l,k} = I_1 + \widetilde{W}_{11}.
\end{aligned}
\tag{30}
$$

We repeat the calculations using the scaled local inverses for the computational molecule in Figure 1(b), where the scaling matrices are of the form $D_1 = diag(1, 1, 1/2)$ and $D_2 = diag(1/2, 1, 1)$. For Problem 1, $\widetilde{P}_{12}$ can now be expressed as

$$\widetilde{P}_{12} = \frac{s}{2}\widetilde{\mathbf{v}}\widetilde{\mathbf{w}}^T, \tag{31}$$

where

$$
\widetilde{\mathbf{v}} = \begin{bmatrix} \frac{a+b+c}{2} \\ -c \\ -b \end{bmatrix} \qquad \widetilde{\mathbf{w}} = \frac{1}{2}\begin{bmatrix} (a+c)^{-1} \\ (a+b)^{-1} \\ (a+c)^{-1} \end{bmatrix}.
$$

13

For the EBES approximation, the expression for the singular value of $\widetilde{P}_{12}$ is

$$\sigma = \frac{s}{4} \frac{\sqrt{\widetilde{f}(a,b,c)\widetilde{g}(a,b,c)}}{(a+c)(a+b)}, \tag{32}$$

where

$$\widetilde{f}(a,b,c) = 5(a+b)^2 + 5(a+c)^2 + 2(ab+ac+bc)$$

and

$$\widetilde{g}(a,b,c) = 2(a+b)^2 + (b+c)^2 + 2ac.$$

Using `Maple` we evaluate Equation (32) and for the three different classes of triangles we have the following values of $\sigma$:

| Triangulation | $a, b, c$ | $\sigma$ |
|---|---|---|
| $\mathcal{T}_1$ | $a = 0$ <br> $b = 1$ <br> $c = 1$ | $2.1213s$ |
| $\mathcal{T}_2$ | $a = -7.9158$ <br> $b = 15.8945$ <br> $c = 15.8945$ | $3.5591s$ |
| $\mathcal{T}_3$ | $a = -0.0315$ <br> $b = 0.0629$ <br> $c = 31.8520$ | $802.505s$ |

We combine the latter values with (26) and compare with the results of the numerically computed spectral bounds, presented in Table 1. The theoretical prediction of the upper bound of the spectrum of $A_{11}B_{11}^{-1}$ is $1 + 3\sigma$, which for $\mathcal{T}_1$ gives $1 + 3 \cdot 2.1213 = 7.3639$, compared to the numerically computed values of $\sigma$ which are of order 4. For the degenerated case $\mathcal{T}_3$, the prediction is $2408.515$, compared with the numerical values of order $1012$.

The theoretical estimate seems to be somewhat larger than the numerically computed values. The reason for this is that the estimate (32) is based on element inverses, which are assumed to share only one node with a neighboring element. The latter holds for elements near the boundary of the problem domain. In practice, the majority of the element inverses are away from the boundary and the scaling matrices are of the form $diag(1/2, 1/2, 1/2)$, which diminishes the corresponding value of $\sigma$, as indicated by the expression (32), compared to the expression in (28), where no scaling is taken into account.

**Remark 3.1** *The above described scaling of the local inverses is one-sided and serves the purpose to obtain an expression for $A_{11}B_{11}^{-1}$, which consists of the exact identity matrix and some error term. The scaling has been shown to play a crucial role when using the element-by-element assembly of local inverses to construct and approximate inverse within the FE context. From the construction it is clear that the scaling is not unique. The scaling can be applied from the left and $B_{11}^{-1}A_{11}$ can be analyzed instead. The particular form of the approximation may depend on the context and the form of the preconditioner.*

*Numerical experiments, not included here, were performed, where the scaling was applied symmetrically. In this way, if $A_{11}$ is symmetric, the symmetricity of $B_{11}^{-1}$ is preserved. However, the tests were not found better than those with the one-sided scaling and further work is needed in order to provide convincing pros and cons with respect to how the scaling should be applied.*

## 3.2 Restricted scaled macroelement pivot inverses (EBERS)

As is seen from (28) and (32), and also from the numerical experiments in Table 1, the relative jump in the problem coefficients acts as a multiplicative factor in the expression of $\sigma$ and entails very large interval containing the real part of the spectrum of $A_{11}B_{11}^{-1}$. This is illustrated in Figures 7 and 8.

We elaborate on the EBES pivot block approximation in order to reduce the dependence of jumps in the problem coefficients. To this end we propose the following idea. Instead of inverting the local macroelement pivot matrices $A_{11,k}$, we proceed as follows:

(i) restrict the assembled block $A_{11}$ to a macroelement $k$,

(ii) invert the so-obtained local matrix, scale it with $D_k$ (name it $\widehat{A}_{11,k}^{-1}$),

(iii) let

$$\widehat{B}_{11}^{-1} = \sum_{k=1}^{M} R_k^T \widehat{A}_{11,k}^{-1} R_k, \tag{33}$$

We refer this approach to as EBERS. In this case we obtain ($k \neq l$)

$$A_{11}\widehat{B}_{11}^{-1} = \widehat{D}_{11} + \sum_{k,l=1}^{q} \widehat{W}_{l,k} = \widehat{D}_{11} + \widehat{W}_{11}, \tag{34}$$

where $\widehat{D}_{11}$ is not identity anymore. It is even not diagonal since the local contributions in $\widehat{B}_{11}^{-1}$ are restrictions of $A_{11}$ to the elements, rather than (scaled) element matrices. In Figures 2(a) and 2(b), $\widehat{D}_{11}$ and $\widehat{W}_{11}$ are shown for a small-sized problem.

Up to the knowledge of the authors, there is no general theory how to estimate the spectrum of the sum of two matrices. Even if $A_{11}$ is spd and we impose a symmetric scaling to the local inverses, the matrix $\widehat{D}_{11}$ will be positive definite but not necessarily symmetric. The EBERS idea is based on intuitive arguments and its efficiency is confirmed by all the test problems in this work. It is however not theoretically justified.

For Problem 1 and the two macroelements in Figure 1, the corresponding pivot matrices $\widehat{A}_{11,1}$ and $\widehat{A}_{11,2}$ are as follows:

$$\widehat{A}_{11,1} = \begin{bmatrix} 2\,(a+b+c)\,(s+1) & -sc & -sb \\ -sc & 2\,s\,(a+b+c) & -sa \\ -sb & -sa & 2\,s\,(a+b+c) \end{bmatrix}$$

15

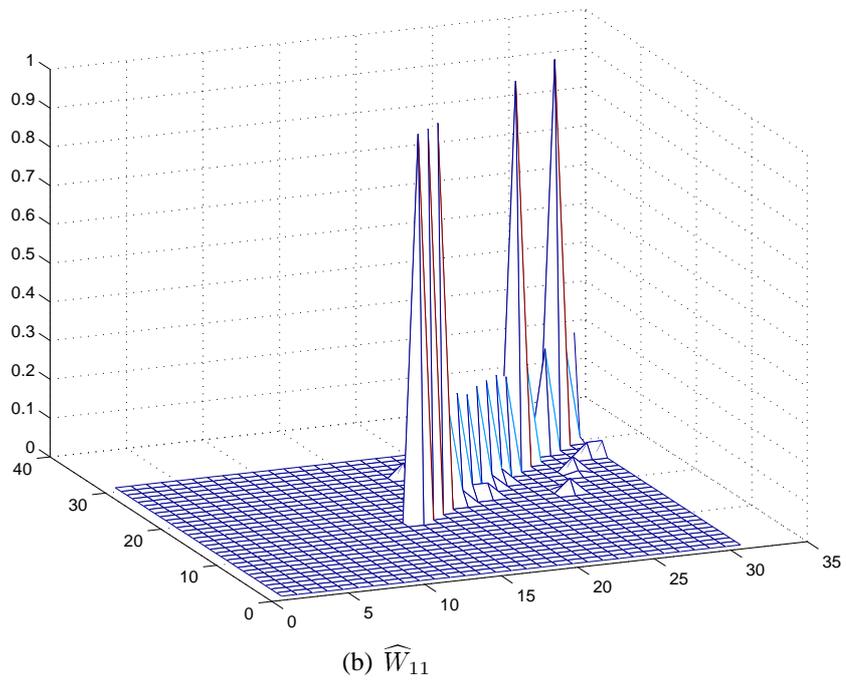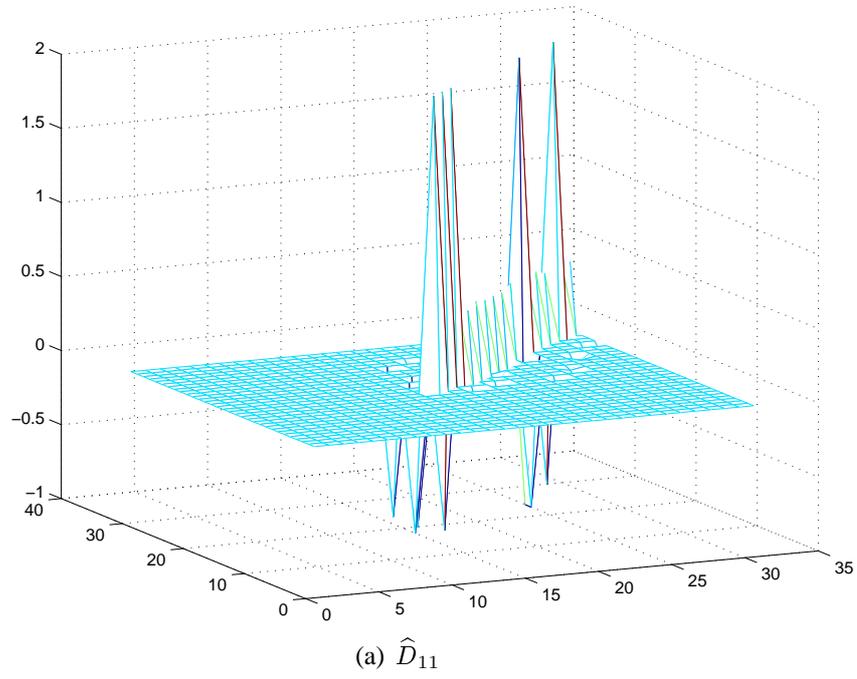(a) $\widehat{D}_{11}$



(b) $\widehat{W}_{11}$

Figure 2: EBERS, portraits of $\widehat{D}_{11}$ and $\widehat{W}_{11}$ for a small-sized problem

16

$$\widehat{A}_{11,2} = \begin{bmatrix} 2\,(a+b+c)\,(s+1) & -c & -b \\ -c & 2\,(a+b+c) & -a \\ -b & -a & 2\,(a+b+c) \end{bmatrix}.$$

Because the expressions become tedious, we do not present here the explicit forms of the matrices $A_{11,l}R_l R_k^T \widehat{A}_{11,k}^{-1} D_k$, their singular value $\sigma$ and the vectors $\widehat{\mathbf{v}}$ and $\widehat{\mathbf{w}}$ as in (31). Instead, we show the evaluation of the analytical expression of $\sigma$ for the different triangulations $\mathcal{T}_i, i = 1, 2, 3$ and three representative values of the jump $s$.

| Triangulation | $a, b, c$ | $\sigma\big|_{s=0.001}$ | $\sigma\big|_{s=1}$ | $\sigma\big|_{s=10000}$ |
|---|---|---|---|---|
| $\mathcal{T}_1$ | $a = 0$ $b = 1$ $c = 1$ | 3.908e-4 | 0.183 | 0.3423 |
| $\mathcal{T}_2$ | $a = -7.9158$ $b = 15.8945$ $c = 15.8945$ | 4.939e-4 | 0.221 | 0.4005 |
| $\mathcal{T}_3$ | $a = -0.0315$ $b = 0.0629$ $c = 31.8520$ | 4.939e-4 | 0.221 | 0.4005 |

As an illustration, we include below the form of the matrix $\widehat{W}_{12}$ for $\mathcal{T}_1$:

$$\begin{bmatrix} 16\,\frac{s}{32\,s+28} & 4\,\frac{s}{32\,s+28} & 4\,\frac{s}{32\,s+28} \\ 0 & 0 & 0 \\ -8\,\frac{s}{32\,s+28} & -2\,\frac{s}{32\,s+28} & -2\,\frac{s}{32\,s+28} \end{bmatrix}$$

The corresponding singular value $\sigma = 3/2\,\sqrt{10}\,\frac{s}{28+32\,s}$ i.e., the dependency on $s$ is much less pronounced.

In the numerical experiments, for EBERS the same right diagonal scaling is applied as for EBES

**Remark 3.2** *The question how to scale the macroelement pivot inverses in the EBERS case should be considered further so that the matrix $\widehat{D}_{11}$ becomes the identity matrix again.*

## 3.3 The construction of $Z_{12}$ and $S$

We apply the idea to approximate a block in the assembled global stiffness matrix by the sum of its corresponding local finite element counterparts to the construction of $Z_{12}$ as an approximation of $A_{11}^{-1}A_{12}$ and to that of the Schur complement approximation $S$.

The off-diagonal block $Z_{12}$ is constructed as

$$Z_{12} = \sum_{k=1}^{M} R_k^T \widetilde{A}_{11,k}^{-1} A_{12,k} R_k,$$

17

where $\widetilde{A}_{11,k}^{-1} = A_{11,k}^{-1} D_k$. If we consider $A_{11} Z_{12}$, then

$$
\begin{aligned}
A_{11} Z_{12} &= \left( \sum_{k=1}^{M} R_k^T A_{11,k} R_k \right) \left( \sum_{l=1}^{M} \widetilde{A}_{11,l}^{-1} A_{12,l} R_l \right) \\
&= \sum_{k=1}^{M} R_k^T A_{11,k} \widetilde{A}_{11,k}^{-1} A_{12,k} R_k + \sum_{k,l=1}^{M} R_k^T A_{11,k} R_k R_l^T \widetilde{A}_{11,l}^{-1} A_{12,l} R_l \\
&= \sum_{k=1}^{M} R_k^T A_{11,k} A_{11,k}^{-1} D_k A_{12,k} R_k + \sum_{k,l=1}^{M} R_k^T A_{11,k} R_k R_l^T A_{11,l}^{-1} D_l A_{12,l} R_l \\
&= \sum_{k=1}^{M} R_k^T \widetilde{A}_{12,k} R_k + \sum_{k,l=1}^{M} R_k^T P_{k,l} \widetilde{A}_{12,l} R_l = \widetilde{A}_{12} + \sum_{k,l=1}^{M} W_{12,kl}
\end{aligned}
$$

Evidently, the matrices $W_{12,kl}$ are of low rank and can be analyzed similarly to the corresponding blocks in (30). We note that the blocks $Z_{12}$ are sparse and cheap to compute explicitly. The results from the numerical experiments in Section 5 convincingly confirm the efficiency of the approach.

The Schur complement approximation $S$ to $S_A$, is constructed again using the same element-by-element approach, i.e., we assembly the local Schur complement matrices $S_k$, computed exactly on each macroelement $k$,

$$
S = \sum_{k=1}^{M} R_k^T (A_{22,k} - A_{21,k} (A_{11,k})^{-1} A_{12,k}) R_k, \tag{35}
$$

For symmetric positive definite (spd) matrices, it is shown in [5] that $S$ and $S_A$ are spectrally equivalent, namely

$$
(1 - \gamma^2) S_A \leq S \leq S_A \tag{36}
$$

where $\gamma$ is the CBS constant, related to the two-level FE splitting. Without having a rigorous proof for the nonselfadjoint case, we have used the same approximation of the Schur complement for nonsymmetric, as well as for symmetric and nonsymmetric saddle point matrices. All numerical results indicate that good spectral bounds, analogous to (36) hold also for more general classes of matrices. We note, that in addition, the above approximation is cheap to compute, sparse and the computational procedure possesses a high degree of parallelism across the macroelements. Furthermore, $S$ automatically inherits the symmetricity or nonsymmetricity of the original Schur complement.

**Remark 3.3** *We point out that, when the local, macroelement Schur complement matrices $S_k = A_{22,k} - A_{21,k} A_{11,k}^{-1} A_{12,k}$ are computed, the local pivot matrices $A_{11,k}^{-1}$ are* not *scaled.*

# 4 Multilevel extension of the two-level preconditioner

Let us assume that we can construct a sequence of FE triangulations $\tau_l$, where $l = L_0, \ldots, L_N$ with $L_0$ being the coarsest mesh. The elements of $\tau_l$ are obtained by $m$-fold regular refinement of the elements of $\tau_{l-1}$. Hence, the meshes are nested such that

$$
\tau_{l+1} \supset \tau_l \quad \text{for all } l = L_0, \ldots, L_{N-1}.
$$

The preconditioner $M_B^{(N)}$ is recursively defined from top to bottom as

$$M_B^{(l)} = \begin{bmatrix} B_{11}^{(l)} & 0 \\ A_{21}^{(l)} & M_B^{(l-1)} \end{bmatrix} \begin{bmatrix} I_1^{(l)} & Z_{12}^{(l)} \\ 0 & I_2^{(l)} \end{bmatrix}, \quad M_B^{(l-1)} = S^{(l)} \tag{37}$$

$$l = L_N, \ldots, L_0 + 1.$$

On each level $l$ the matrix $M_B^{(l)}$ is split into two-by-two block form, aligned with the fine-coarse splitting of the nodes on that level, and the corresponding blocks $B_{11}^{(l)-1}$, $Z_{12}^{(l)}$ and $S^{(l)}$ are constructed as described in Section 3 based on macroelement matrices. The block $A_{21}$ is computed by a standard assembly procedure. On the finest level $L_N$, the element matrices are the standard FE element stiffness matrices. On each coarser level $l - 1$, the element matrices are the local Schur complement matrices $S_k^{(l)}$, computed on the macroelements on level $l$.

Equation (37) describes a $V$-cycle multilevel preconditioner. As for the classical HBF and AMLI methods, the condition number of the preconditioned matrix $M_B^{(L)-1} A^{(L)}$ is known to deteriorate with a growing number of levels $L - L_0$. This growth can, for instance, be stabilized by a few iterations with an inner iterative solution method on some of the levels in the hierarchy. For further details we refer to [9], [10], [11], [23], and [2], [21], [26].

# 5   Numerical illustrations

We illustrate the performance of the proposed approximations $Z_{12}$, $\widetilde{B}_{11}^{-1}$ and $\widehat{B}_{11}^{-1}$, and the corresponding block-factorized preconditioners $\widetilde{M}_B$, and $\widehat{M}_B$ on the following set of test problems:

**Problem 2** Scalar isotropic Poisson equation, solved in a model domain, depicted in Figure 4(b).

**Problem 3** Scalar anisotropic Poisson equation, solved in a model domain, depicted in Figure 3(a) and 3(b). In this case we can control the mesh anisotropy by choosing the size of one of the angles, equal to $\pi/4$ in Figure 3(a) and $\pi/50$ in Figure 3(b).

**Problem 4** Scalar isotropic Poisson equation with discontinuous coefficients $a(\mathbf{x})$

$$-\nabla \cdot (a(\mathbf{x})\nabla u) = f. \tag{38}$$

We consider two different geometries of $\Omega$, referred to as (a) and (b). The former one with the corresponding initial (coarsest) triangulation are shown in Figure 4(a), where $\Omega_\varepsilon$ occupies the shaded region, $a = \varepsilon$ in $\Omega_\varepsilon$ and $a = 1$ elsewhere.

The second computational domain is the union of two subdomains, $\Omega = \Omega_1 \cup \Omega_2$, and the coefficient $a(\mathbf{x})$ is piecewise constant in $\Omega_1$ and $\Omega_2$, where

$$a(\mathbf{x}) = a_0 \qquad \forall \mathbf{x} \in \Omega_1 = \{\mathbf{x} : 0.0 \le x_2 < 0.5\}$$
$$a(\mathbf{x}) = \varepsilon\, a_0 \qquad \forall \mathbf{x} \in \Omega_2 = \{\mathbf{x} : 0.5 \le x_2 \le 1.0\},$$

with $\varepsilon = 0.001, 1$,and $10000$ and constant $a_0 = O(1)$.
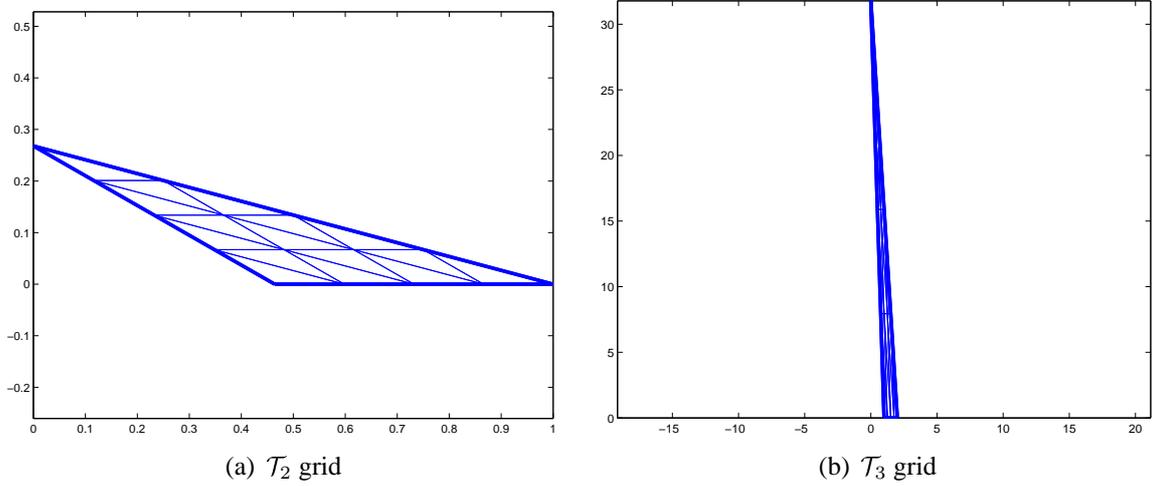
(a) $\mathcal{T}_2$ grid

(b) $\mathcal{T}_3$ grid

Figure 3: Anisotropic meshes of type $\mathcal{T}_2$ and $\mathcal{T}_3$. In Figure 3(a), the small angles equal $\pi/4$, whereas in Figure 3(b), the small angle is $\pi/50$
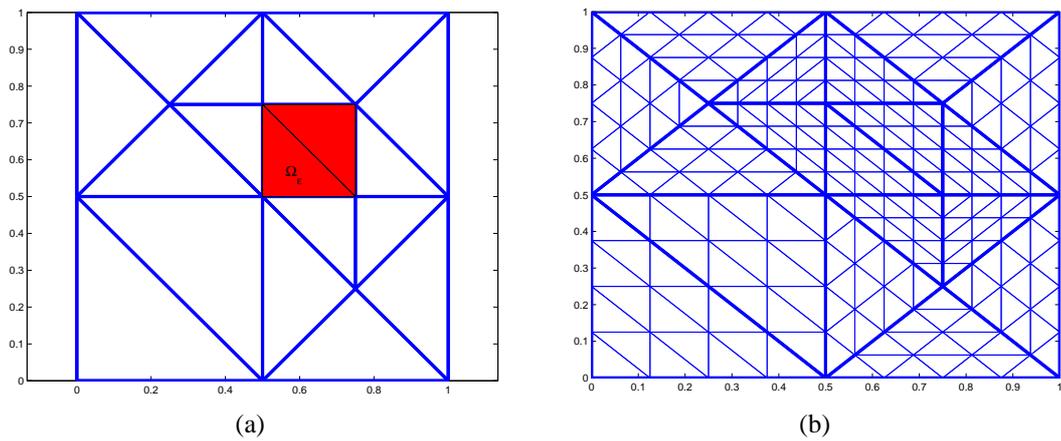


(a)

(b)

Figure 4: Meshes of $\mathcal{T}_1$ type. The geometry of Problem 4(a)

**Problem 5** Scalar convection-diffusion equation with constant convection on the unit square

$$-\nabla(a \cdot \nabla u) - \mathbf{b} \cdot \nabla u = f, \tag{39}$$

where $\mathbf{b} \in \mathbb{R}^2$. The convective wind has magnitude $R$, and its direction is determined by $p = 0, 1, 2, \ldots, 6$, that is, $b_1 = R\cos(p\pi/6)$ and $b_2 = R\sin(p\pi/6)$.

For Problems 2, 3, and 4(a) the PDEs are discretized on a triangular finite element mesh, and standard linear basis functions are used. Problem 4(b) and 5 are discretized on a uniform regular quadrilateral grid, where the elements are equipped with bilinear basis functions. The outer iterative solution method is chosen as GCG-MR, and it is iterated until the residual norm is decreased by six orders of magnitude compared to the initial residual. The Schur complement approximation $S$ is solved exactly by a direct solution method.

The code used for Problems 2, 3, and 4(a) is written in MATLAB, whereas for Problems 4(b) and 5 the implementation is made in C++. The latter code is based on the open source packages `deal.II` [12] and `PETSc` [20].

## 5.1   Results for the EBES approach

For all experiments in this section $\widetilde{B}_{11}^{-1}$ is constructed using EBES.

In Table 1 we present the numerically computed extreme values of the spectrum of $A_{11}\widetilde{B}_{11}^{-1}$ for Problems 2 - 4(a). We include also the iteration counts required for the outer solution method to converge up to the given tolerance.

The crucial step when computing the action of $B_M$ on a vector is the (inexact) solution of the pivot block $A_{11}$. In Table 1 we present three different approaches to do this: (i) a direct (exact) solve with the block $A_{11}$, (ii) replacing $A_{11}^{-1}\mathbf{b}_1$ by $\widetilde{B}_{11}^{-1}\mathbf{b}_1$, and (iii) using an inner iterative solution method with a weaker stopping criterion ($10^{-4}$) to solve the system, preconditioned by $B_{11}^{-1}$. Note that the preconditioner is multiplicative, the approximation is sparse by construction and thus, the inner iterations are cheap. As inner solver we recommend again GCG-MR, since although $A_{11}$ is symmetric in the test problems, $B_{11}^{-1}$ is nonsymmetric due to the applied one-sided scaling.

As is expected, the spectrum of $A_{11}\widetilde{B}_{11}^{-1}$ is mesh-independent. We also see that while the smallest eigenvalue is bounded away from zero, the largest is proportional to the amount of mesh anisotropy or the jump of the coefficients in the differential equation.

In the case of not very strong anisotropy, we see from Figure 5 that the spectrum is not clustered but is contained in a narrow ellipse, which is also favorable for the convergence of an iterative method. The seemingly most difficult case is that with a strong anisotropy (almost degenerate triangles). From Figure 6 we see that the spectrum is almost real but it is not clustered and spans a very large interval. Still, solving the systems in step (S1) by a preconditioned GCG-MR leads to outer iteration counts, which are competitive to the case when systems with $A_{11}$ are solved exactly.

We also see from Figures 7 and 8, that for discontinuous coefficients the eigenvalues are well-clustered and with small imaginary parts apart from a few very large outliers, which could be easily eliminated by an inner iterative solution method.

| Problem size | $\lambda(A_{11}\widetilde{B}_{11}^{-1})$ | | Iterations (block-factorized prec. GCG) | | |
|---|---|---|---|---|---|
| | $\lambda_{min}$ | $\lambda_{max}$ | exact solve | mult. by $B_{11}^{-1}$ | inner (gcg) |
| Problem 2 | | | | | |
| 153 | 1 | 3.8355 | 8 | 22 | 8(6) |
| 561 | 1 | 3.9551 | 8 | 26 | 8(6) |
| 2145 | 1 | 3.9884 | 8 | 34 | 9(7) |
| Problem 3: Mesh $\mathcal{T}_2$, small angles $\pi/12$ | | | | | |
| 153 | 1 | 5.63+$i$0.01 | 11 | nc | 13(11) |
| 561 | 1 | 5.64+$i$0.002 | 12 | nc | 13(11) |
| 2145 | 1 | 5.71+$i$0.0002 | 12 | nc | 13(11) |
| Problem 3: Mesh $\mathcal{T}_2$, small angles $\pi/50$ | | | | | |
| 153 | 1 | 5.63+$i$0.01 | 13 | nc | 13(11) |
| 561 | 1 | 5.89+$i$0.002 | 13 | nc | 14(11) |
| 2145 | 1 | 5.96+$i$0.0002 | 13 | nc | 15(11) |
| Problem 3: Mesh $\mathcal{T}_3$, small angle $\pi/50$ | | | | | |
| 153 | 1 | 976.0 | 5 | nc | 6(38) |
| 561 | 1 | 1004.8 | 6 | nc | 6(85) |
| 2145 | 1 | 1012.1 | 6 | nc | 7(106) |
| Problem 4(a), $\varepsilon = 1$ | | | | | |
| 161 | 0.5 | 4.2133 | 7 | 21 | 7(9) |
| 609 | 0.5 | 4.2258 | 8 | 21 | 8(9) |
| 2369 | 0.5 | 4.2263 | 8 | 21 | 8(9) |
| Problem 4(a), $\varepsilon = 10^{-3}$ | | | | | |
| 161 | 0.5 | 1001.69 | 8 | 29 | 8(16) |
| 609 | 0.5 | 1001.69 | 8 | 37 | 8(18) |
| 2369 | 0.5 | 1001.69 | 8 | 53 | 8(19) |
| Problem 4(a), $\varepsilon = 10^4$ | | | | | |
| 161 | 0.5 | 7501.95 | 7 | nc | 8(8) |
| 609 | 0.5 | 7501.95 | 8 | nc | 9(8) |
| 2369 | 0.5 | 7501.95 | 8 | nc | 9(12) |

Table 1: Problems 2, 3 and 4(a): EBES approach. The extreme eigenvalues of $A_{11}\widetilde{B}_{11}^{-1}$ and the number of iterations for the outer GCG-MR solution method. The numbers in the parentheses in the right-most column are the inner iterations required for convergence.
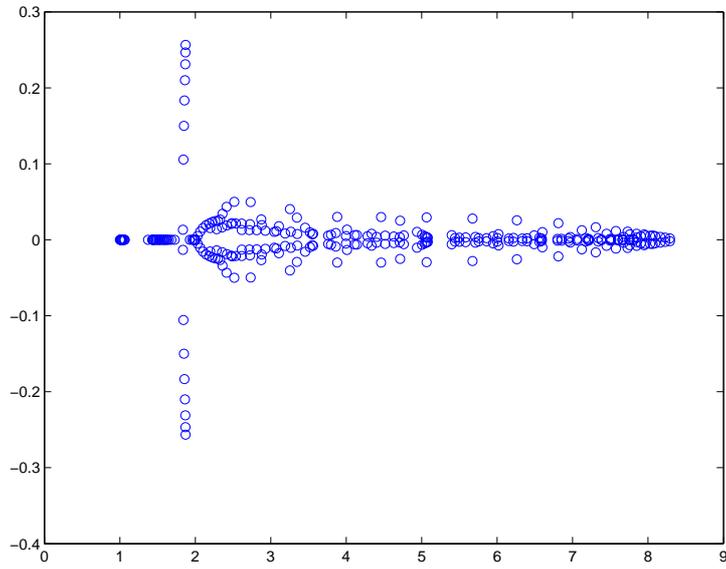
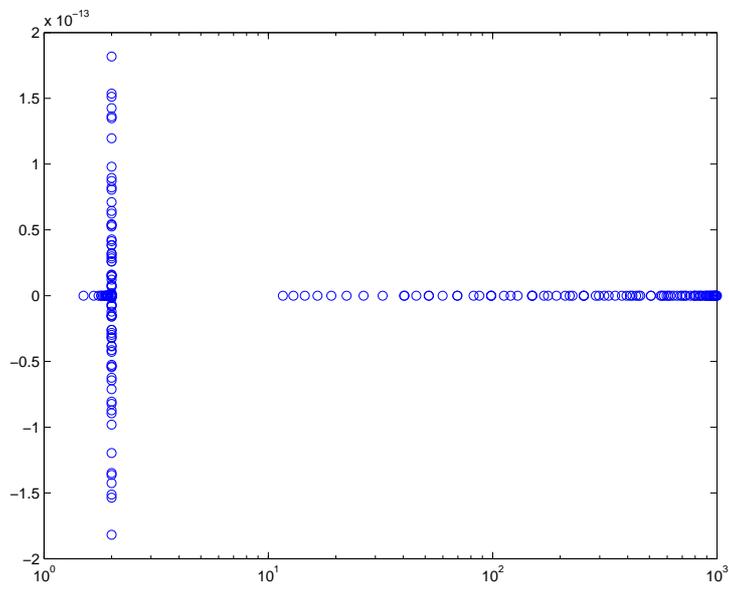Figure 5: Problem 3: Mesh $\mathcal{T}_3$, small angle $\pi/4$: Spectrum of $A_{11}\widetilde{B}_{11}^{-1}$, (`size(A)=561`)



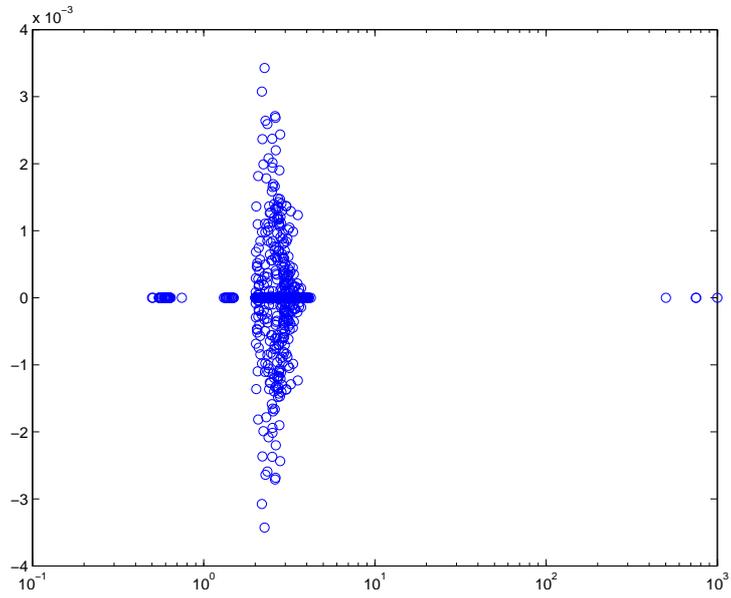Figure 6: Problem 3: Mesh $\mathcal{T}_3$, small angle $\pi/50$: Spectrum of $A_{11}\widetilde{B}_{11}^{-1}$, (`size(A)=561`).

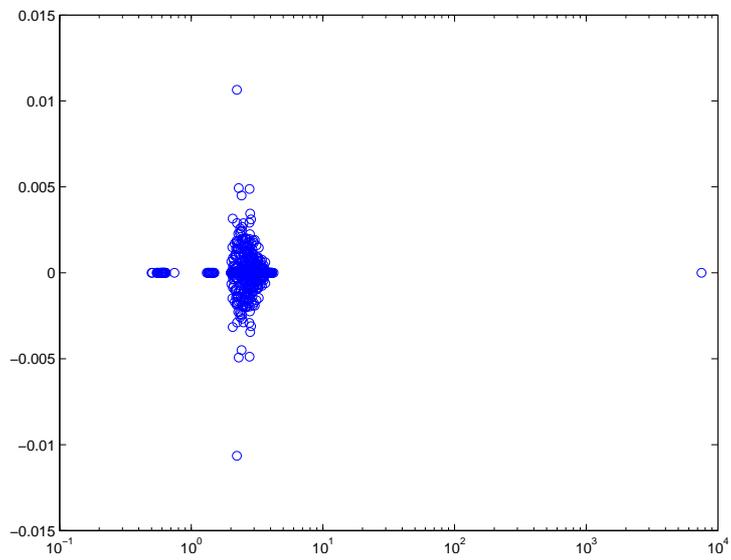Figure 7: Problem 4(a), $\varepsilon = 10^{-3}$: Spectrum of $A_{11}\widetilde{B}_{11}^{-1}$, (`size(A)=2369`)



Figure 8: Problem 4(a), $\varepsilon = 10^4$: Spectrum of $A_{11}\widetilde{B}_{11}^{-1}$, (`size(A)=2369`)

| Problem size | $\lambda(A_{11}\widehat{B}_{11}^{-1})$ | | Iterations (block-factorized prec. GCG) | | |
|---|---|---|---|---|---|
| | $\lambda_{min}$ | $\lambda_{max}$ | exact solve | mult. by $B_{11}^{-1}$ | inner (gcg) |
| $\varepsilon = 1$ | | | | | |
| 161 | 0.48438 | 1.2743 | 8 | 25 | 8(5) |
| 609 | 0.46155 | 1.2804 | 9 | 52 | 9(6) |
| 2369 | 0.45529 | 1.2838 | 9 | >100 | 11(6) |
| $\varepsilon = 10^{-3}$ | | | | | |
| 161 | 0.495 | 1.2793 | 8 | 24 | 8(6) |
| 609 | 0.465 | 1.2832 | 9 | 52 | 9(6) |
| 2369 | 0.457 | 1.2859 | 9 | >100 | 12(6) |
| $\varepsilon = 10^4$ | | | | | |
| 161 | 0.47648 | 1.302 | 7 | nc | 8(5) |
| 609 | 0.46396 | 1.298 | 8 | nc | 8(7) |
| 2369 | 0.45783 | 1.297 | 9 | nc | 10(7) |

Table 2: Problem 4(a): EBERS approach. The extreme eigenvalues of $A_{11}\widehat{B}_{11}^{-1}$ and the number of iterations for the outer GCG-MR solution method. The numbers in the parentheses in the right-most column are the inner iterations required for convergence.

## 5.2 Results for the EBERS approach

With the next series of experiments we demonstrate the performance of the EBERS approximation $\widehat{B}_{11}^{-1}$ of the pivot block on the five test problems.

In Table 2 we present data regarding the extreme eigenvalues of $A_{11}\widehat{B}_{11}^{-1}$ for Problem 4(a). Also included are the iteration counts required for the outer iterative solution method to meet the convergence criterion ($10^{-6}$). As is expected from the theoretical results, the eigenvalues of $A_{11}\widehat{B}_{11}^{-1}$ are independent of the size of the mesh, and jumps in the coefficients. Furthermore, they are bounded away from zero.

Figures 9 and 10 show the spectra of $A_{11}\widehat{B}_{11}^{-1}$ for Problem 4(a), for two different values of the discontinuity parameter $\varepsilon$. It is $10^{-3}$ in the former case, and $10^4$ in the latter. The spectra are contained in a narrow ellipse, and are insensitive to the choice of $\varepsilon$. Figure 11 show the spectrum of the whole preconditioned matrix $\widehat{M}_B^{-1}A$, for Problem 4(b) with $\varepsilon$ equal to $10^4$. Apart from two outliers, the spectrum of $\widehat{M}_B^{-1}A$ is clustered at unity and contained in a narrow ellipse.

**Remark 5.1** *In Problem 4(b) and 5, $M_B$ is implemented as a "left preconditioner", which is in contrast to the rest of the report, where it is constructed as a "right preconditioner".*

In Figures 12 and 13 respectively, the spectrum of the full preconditioned matrix $\widehat{M}_B^{-1}A$, and the preconditioned pivot block $\widehat{B}_{11}^{-1}A_{11}$, for Problem 5 are shown. The wind has unit magnitude ($R = 1$), and it is aligned with the $x_1$-axis ($p = 0$). Both spectra are clustered around unity and bounded away from zero.

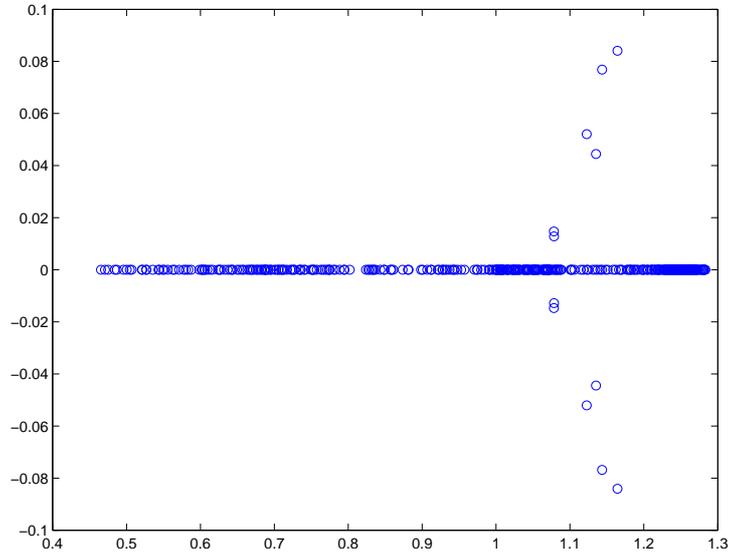Figure 9: Problem 4(a),$\varepsilon = 10^{-3}$: Spectrum of $A_{11}\widehat{B}_{11}^{-1}$, (`size(A)=609`)



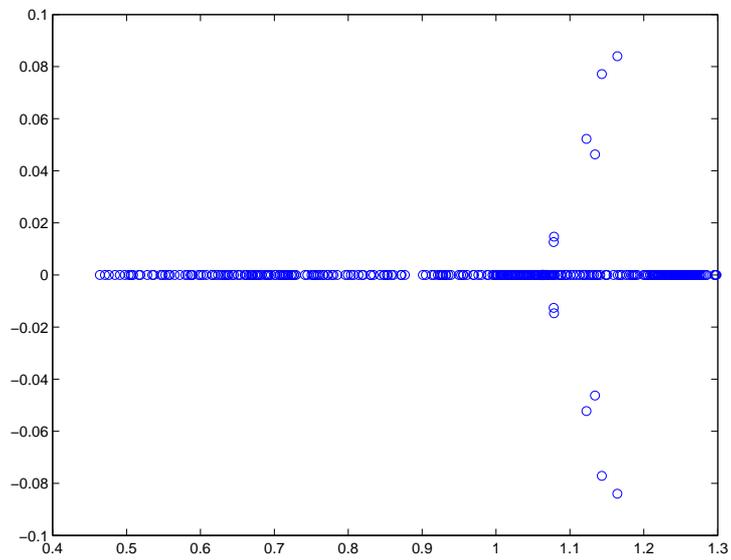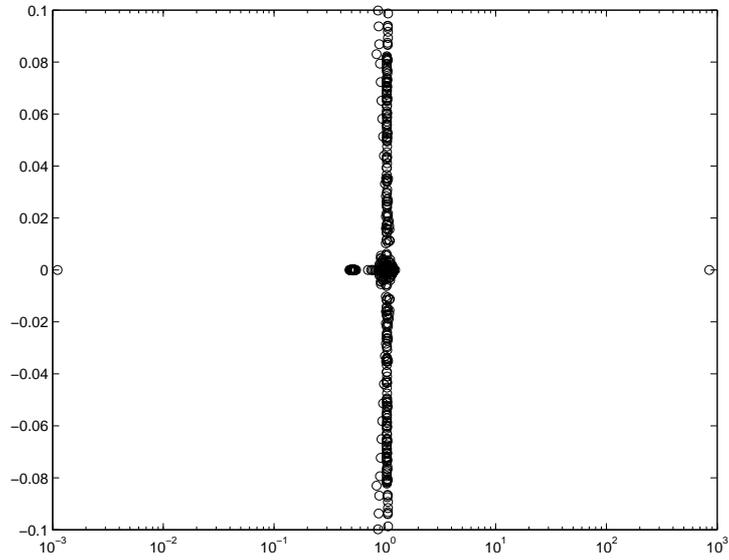Figure 10: Problem 4(a), $\varepsilon = 10^4$: Spectrum of $A_{11}\widehat{B}_{11}^{-1}$, (`size(A)=609`)

Figure 11: Problem 4(b), $\varepsilon = 10^4$: Spectrum of $\widehat{B}_M^{-1}A$, (`size(A)=1089`)

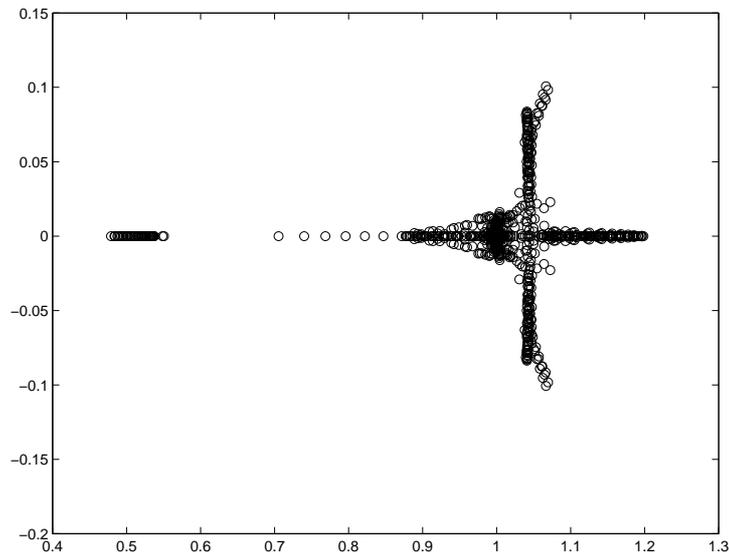

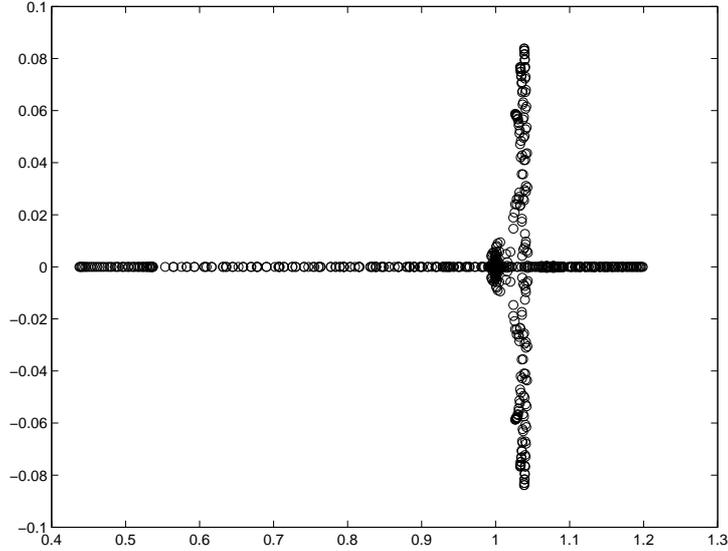Figure 12: Problem 5, $R = 1$, $p = 0$: Spectrum of $\widehat{B}_M^{-1}A$, (`size(A)=1089`)

27

Figure 13: Problem 5, $R = 1$, $p = 0$: Spectrum of $\widehat{B}_{11}^{-1} A_{11}$, (`size(A)=1089`)

In the experiments accounted for in Tables 3, 4, and 5, we solve Problems 4(b) and 5. Here, the inner iterative solution method for the pivot block performs a fixed number of iterations, rather than meets a given tolerance. The number of inner iterations in Tables 3, 4, and 5 is denoted by "Inner it.".

Table 3 shows the number of outer iterations required to solve Problem 4(b) for different problem sizes, number of inner iterations to solve with $A_{11}$, and values of the jump coefficient $\varepsilon$. The stars ($*$) in the table indicate stagnation of the outer iterative solution method. The preconditioner is sensitive to the discontinuity in the coefficients, and for $\varepsilon = 10000$ the convergence of the outer iterative solver is destroyed when the pivot block is not solved accurate enough.

Table 4 shows iteration counts for Problem 5 for different problem sizes and number of inner iterations for the pivot block. The convection is of unit magnitude ($R = 1$) and it is aligned with the $x_1$-axis. In Table 5 we present outer iteration counts for Problem 5 for different magnitudes and directions of the convective field. The size of the problem is $N = 4225$. The results show that the preconditioner is robust with respect to the considered range of the problem parameters, and that the convergence of the outer iterative solution method does not depend on the size of the problem.

# 6 Conclusions

In this report, we consider block-factorized preconditioners to a matrix $A$ given in a two-by-two block form based on a splitting of the unknowns as fine and coarse. The matrix arises from a two-level finite element discretization of an elliptic PDE. We propose and analyze two novel strategies of element-by-element type to construct a sparse approximation of the inverse of the

| Inner it. | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| | $\varepsilon = 0.001$ | | | | |
| $N = 1089$ | 11 | 8 | 7 | 7 | 7 |
| $N = 4225$ | 11 | 8 | 7 | 7 | 7 |
| $N = 16641$ | 11 | 8 | 7 | 7 | 7 |
| $N = 66049$ | 11 | 8 | 7 | 7 | 7 |
| | $\varepsilon = 1$ | | | | |
| $N = 1089$ | 9 | 7 | 7 | 7 | 7 |
| $N = 4225$ | 9 | 7 | 7 | 7 | 7 |
| $N = 16641$ | 9 | 7 | 7 | 7 | 7 |
| $N = 66049$ | 8 | 7 | 7 | 7 | 7 |
| | $\varepsilon = 10000$ | | | | |
| $N = 1089$ | * | 12 | 10 | 10 | 9 |
| $N = 4225$ | * | * | * | 9 | 9 |
| $N = 16641$ | * | * | 10 | 8 | 8 |
| $N = 66049$ | * | * | 10 | 9 | 9 |

Table 3: Problem 4(b). Iteration counts for different problem sizes and different values of $\varepsilon$. The star (∗) indicates stagnation of the iterative solution method.

| Inner it. | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $N = 1089$ | 9 | 7 | 7 | 7 | 7 |
| $N = 4225$ | 9 | 7 | 7 | 7 | 7 |
| $N = 16641$ | 9 | 7 | 7 | 7 | 7 |
| $N = 66049$ | 9 | 7 | 7 | 7 | 7 |

Table 4: Problem 5. Iteration count for different problem sizes. The convection is parallel with the $x_1$-axis and $R = 1$.

| Inner it. | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|---|---|---|
| | | $R = 1$ | | | | | $R = 2$ | | | |
| $p = 0$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 1$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 2$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 3$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 4$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 5$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 6$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| | | $R = 3$ | | | | | $R = 4$ | | | |
| $p = 0$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 1$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 2$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 3$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 4$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 5$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |
| $p = 6$ | 9 | 7 | 7 | 7 | 7 | 9 | 7 | 7 | 7 | 7 |

Table 5: Problem 5. Iteration counts for different magnitude and directions on the convection. The problem size is $N = 4225$.

pivot block $A_{11}^{-1}$. The approximation is based on assembly of scaled, locally inverted matrices, computed exactly on macroelements in the finite element mesh. In the first approach (EBES), the local matrix is the pivot block $A_{11,k}$ in the macroelement matrix $A_k$, while in the second strategy (EBERS), the local matrix is a restriction of the global, already assembled pivot matrix $A_{11}$ to the macroelement. The two approaches are analyzed theoretically for a symmetric test problem in two space dimensions. The analysis reveals that both strategies result in an approximation that is independent of the size of the problem, and that the spectrum of the rank-one matrices $\widehat{W}_{l,k}$ in the EBERS approach is independent of discontinuities in the coefficients of the underlying PDE.

Beside the strategies to approximate the inverse of the pivot block, we also propose a method to construct a sparse approximation $Z_{12}$ of the off-diagonal block $A_{11}^{-1}A_{12}$ in $M_B$. The matrix $Z_{12}$ if assembled from local, elementwise, exactly computed, scaled contributions of the form $Z_{12,k} = A_{11,k}^{-1}A_{12,k}$.

Extensive numerical experiments on a series of symmetric and nonsymmetric test problems confirm that (i) the analytical results for the EBER and EBERS strategies, (ii) the EBERS approach is robust with respects to discontinuities in the coefficients of the underlying PDE, (iii) the efficiency of the $Z_{12}$ approximation, and (iv) that the proposed block-factorized two-level preconditioner is robust also for the considered nonsymmetric convection-diffusion problem.

The theoretical analysis for the pivot block approximation is provided in two space dimensions, and for two-fold refinement of the coarse mesh. It is straightforward to extend the theory to 3D and for $m > 2$. However in these cases the analysis is harder to perform because the local products $A_{11,k}R_kR_l^TA_{11,l}^{-1}$ will still be of low rank but in general higher than one. On the other

hand, these can be represented as a sum of rank-one matrices, converting back to the framework used in this paper.

The strategy to refine the meshes using values of $m$ larger than two needs more attention. To apply it straightforwardly is shown to be not efficient, however for the case when no scaling is used. One question to study further is what will be the effect of applying a proper scaling on the properties of the approximate inverse in this case. Further, this strategy will make the coarsening of the two- or multilevel method more aggressive, resulting in a smaller size of the coarse mesh problem. However, for larger $m$, the size of the local matrices grows, and so does the arithmetic cost to invert the local pivot blocks. Therefore, the choice of $m$ must balance fast coarsening with the cost to compute the inverse of the local pivot block. In addition, with increasing $m$, the bound on the CBS-constant $\gamma$ in Equation (36) approaches one (cf. [4]).

As a final remark, we would like to emphasize that the proposed techniques to construct approximations of matrix expressions are suitable for implementation on parallel computers, since all computations in the construction phase are local and fully decoupled.

# Acknowledgements

# References

[1] O. Axelsson. *Iterative solution methods*. Cambridge University Press, 1994.

[2] O. Axelsson. Stabilization of algebraic multilevel iteration methods; additive methods. *Numerical Algorithms*, 21:23 – 47, 1999.

[3] O. Axelsson and V.A. Barker. *Finite Element Solution of Boundary Value Problems. Theory and Computation*. Academic Press, Inc, 1984.

[4] O. Axelsson and R. Blaheta. Two simple derivations of universal bounds for the C.B.S inequality constant. *Applications of Mathematics*, 49(1):57 – 72, 2001.

[5] O. Axelsson, R. Blaheta, and M. Neytcheva. Preconditioning of boundary value problems using elementwise Schur complements. Technical Report 2006-048, Department of Information Technology, Uppsala University, November 2006.

[6] O. Axelsson and V. Eijkhout. The nested recursive two-level factorization method for nine-point difference matrices. *SIAM Journal on Statistical and Scientific Computing*, 12(6):1373 – 1400, 1991.

[7] O. Axelsson and I. Gustafsson. Preconditioning and two-level mulitgrid methods of arbitrary degree of approximation. *Mathematics of Computation*, 40(161):219 – 242, 1983.

[8] O. Axelsson and M. Neytcheva. Preconditioning methods for linear systems arising in constrained optimization problems. *Numerical Linear Algebra with Applications*, 10:3–31, 2003.

[9] O. Axelsson and P.S. Vassilevski. Algebraic multilevel preconditioning methods I. *Numerische Mathematik*, 56(2-3):157–177, 1989.

[10] O. Axelsson and P.S. Vassilevski. Algebraic multilevel preconditioning methods II. *SIAM Journal on Numerical Analysis*, 27(6):1569 – 1590, 1990.

[11] O. Axelsson and P.S. Vassilevski. Variable-step multilevel preconditioning methods, I: Self-adjoint and positive definite elliptic problems. *Numerical Linear Algebra with Applications*, 1(1):75 – 101, 1994.

[12] W. Bangerth, R. Hartmann, and G. Kanschat. `deal.II` *Differential Equations Analysis Library, Technical Reference*. IWR. `http://www.dealii.org`.

[13] M. Benzi, G.H. Golub, and J. Liesen. Numerical solution of saddle point problems. *Acta Mathematica*, 14:1–137, 2005.

[14] R. Bhatia. *Matrix analysis*. Springer-Verlag, New York, 1997.

[15] E.F.F. Botta and F.W. Wubs. Matrix renumbering ILU: an effective algebraic multilevel ILU preconditioner for sparse matrices. *SIAM Journal on Matrix Analysis and Applications*, 20(4):1007–1026, 1999.

[16] E. Chow and Y. Saad. Approximate inverse techniques for block-partitioned matrices. *SIAM Journal on Scientific Computing*, 18(6):1657 – 1675, 1997.

[17] V. Eijkhout and P.S. Vassilevski. The role of the strenghened Cauchy-Buniakowskii-Schwarz inequality in multilevel methods. *SIAM Review*, 33(3):405 – 419, 1991.

[18] J.K. Kraus. Algebraic multilevel preconditioning of finite element matrices using local Schur complements. *Numerical Linear Algebra with Applications*, 13(1):49–70, 2006.

[19] Maplesoft, a division of Waterloo Maple Inc. *Maple 9 - Learning Guide*, 2003.

[20] Mathematics and Computer Science Division, Argonne National Laboratory. *Portable, Extensible Toolkit for Scientific computation (PETSc) suite*. `www-unix.mcs.anl.gov/petsc/`.

[21] Y. Notay. Optimal V-cycle algebraic multilevel preconditioning. *Numerical Linear Algebra with Applications*, 5:441 – 459, 1998.

[22] Y. Notay. Using approximate inverses in algebraic multilevel methods. *Numerische Mathematik*, 80(3):397–417, 1998.

[23] Y. Notay. Robust parameter-free algebraic multilevel preconditioning. *Numerical Linear Algebra with Applications*, 9:409 – 428, 2002.

[24] Y. Saad. Multilevel ILU with reorderings for diagonal dominance. *SIAM Journal on Scientific Computing*, 27(3):1032 – 1057, 2005.

[25] Y. Saad and B. Suchomel. ARMS: an algebraic recursive multilevel solver for general sparse linear systems. *Numerical Linear Algebra with Applications*, 9:359–378, 2002.

[26] P.S. Vassilevski. On two ways of stabilizing the hierarchical basis multilevel methods. *SIAM Review*, 39(1):18–53, March 1997.