

# Preconditioning of boundary value problems using elementwise Schur complements

Owe Axelsson<sup>1</sup>, Radim Blaheta<sup>2</sup>, Maya Neytcheva<sup>3</sup>

November 2, 2006

## Abstract

Based on a particular node ordering and corresponding block decomposition of the matrix we analyse an efficient, algebraic multilevel preconditioner for the iterative solution of finite element discretizations of elliptic boundary value problems. Thereby an analysis of a new version of block-factorization preconditioning methods is presented. The approximate factorization requires an approximation of the arising Schur complement matrix. In this paper we consider such approximations derived by the assembly of the local macro-element Schur complements. The method can be applied also for non-selfadjoint problems but for the derivation of condition number bounds we assume that the corresponding differential operator is selfadjoint and positive definite.

## 1 Introduction

Let  $\mathcal{L}$  be a second order elliptic operator on a plane domain  $\Omega$  with proper boundary conditions and consider the finite element solution of  $\mathcal{L}u = f$  on a triangular or quadrilateral mesh. The initial mesh is assumed to be adjusted to the geometry of the domain and to discontinuities of the coefficients in the differential operator. Each element of the mesh is subdivided into  $m^2$ ,  $m \geq 2$ , congruent elements using a uniform refinement. Each such element forms then a macro-element. It is also possible to use locally regular refinements as will be briefly mentioned later. Except when mentioned otherwise, we assume that the coefficients in the differential operator are constant on each macro-element. Arbitrary jumps can, however, occur between macro-elements. By ordering the edge nodes and the interior nodes first, and the vertex nodes last, and partitioning the finite element matrix  $A$  correspondingly, we obtain

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}. \quad (1)$$

---

<sup>1</sup>Department of Information Technology, Uppsala University, Sweden & Institute of Geonics AS CR, Ostrava, The Czech Republic, email owea@it.uu.se

<sup>2</sup>Institute of Geonics AS CR, Ostrava, The Czech Republic, email blaheta@ugn.cas.cz

<sup>3</sup>Department of Information Technology, Uppsala University, Sweden, email maya@it.uu.se

The block matrix factorization of  $A$  can be written as

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

Here  $I_1, I_2$  are unit matrices,  $A_{11}$  is the pivot block and

$$S_A = A_{22} - A_{21}A_{11}^{-1}A_{12}$$

is the corresponding Schur complement matrix.

To obtain an approximate factorization  $B$  of  $A$  we let  $B_{11}$  be an approximation of  $A_{11}$ ,  $\tilde{S}_A$  be some approximation of  $S_A$  and let

$$B = \begin{bmatrix} B_{11} & 0 \\ A_{21} & I_2 \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1}A_{12} \\ 0 & \tilde{S}_A \end{bmatrix}. \quad (2)$$

Note that  $B$  is nonsingular if and only if  $B_{11}$  and  $\tilde{S}_A$  are nonsingular and  $B$  is symmetric and positive definite if and only if  $B_{11}$  and  $\tilde{S}_A$  are symmetric positive definite. Various approximations  $B_{11}$  of  $A_{11}$  can be found in earlier publications, see e.g [8]. In the framework of block-factorized preconditioners of type (2), a common choice for the approximations of  $S_A$  has been the coarse mesh finite element matrix  $A_H$ ,  $H = mh$  where  $H$  and  $h$  denote the coarse and fine meshsizes, respectively. It has been shown that

$$(1 - \gamma^2)A_H \leq S_A \leq A_H, \quad (3)$$

where  $\gamma^2 = \rho \left( A_H^{-1} \hat{A}_{21} A_{11}^{-1} \hat{A}_{12} \right)$ . Here  $\rho$  denotes the spectral radius and  $\gamma$  is the so-called Cauchy-Bunyakowski-Schwarz (CBS) constant for the corresponding hierarchical basis functions (HBF) matrix  $\hat{A} = \begin{bmatrix} A_{11} & \hat{A}_{12} \\ \hat{A}_{21} & A_H \end{bmatrix}$ , where the vertex node basis functions have been replaced by the basis functions for the vertex nodes for the coarse mesh elements, see [6, 4]. A special Schur complement approximation for Stieltjes matrices has been considered in [9].

We note that the matrix  $\hat{A}$  can be written as an assembly of local finite macro-element matrices, which have the same partitioning as in (1).

In this paper we consider instead approximations of  $S_A$  based on the assembly of the local Schur complements on each macro-element. Such an approximation has been considered earlier in [17] for the case  $m = 2$  and rectangular elements. Here we extend these results to general values of  $m$ , consider triangular elements and derive more general eigenvalue bounds and in a more direct way than in [17]. Our analysis is based on the CBS constant.

The paper is organized as follows. In Section 2, the ideas presented in this introduction are formulated more thoroughly and in a multilevel setting. Sections 3 and 4 are devoted to the analysis of two-level preconditioners and possible improvement of the hierarchical

method by congruence transformation. The elementwise Schur complement approximation is introduced and investigated in Section 5. Section 6 describes a possible further improvement which can be obtained by combining the elementwise Schur complements with the coarse mesh finite element matrices. Section 7 discusses an approximation of the pivot block using again the elementwise concept. Section 8 deals with the complexity and optimality, as well as some implementational aspects including a possible local refinement variant of the multilevel method. The results are illustrated by numerical tests presented in Section 9.

## 2 AMLI preconditioners

Consider a sequence of nested meshes  $\mathcal{T}_k$ ,  $k = k_0, \dots, \ell$ . Partitioning the finite element mesh on each level in two subsets, containing edge and interior nodes, and vertex nodes, respectively, induces a block  $2 \times 2$  structure of the corresponding finite element matrix,

$$A_k = \begin{bmatrix} \widehat{A}_{k,11} & A_{k,12} \\ A_{k,21} & A_{k,22} \end{bmatrix},$$

where  $k$  denotes the current level number. For the solution of a linear system with  $A_k$  one can use the original Algebraic Multilevel Iteration (AMLI) method, see [11]. It is based indirectly on the related hierarchical basis representation of the matrix,

$$\widehat{A}_k = \begin{bmatrix} A_{k,11} & \widehat{A}_{k,12} \\ \widehat{A}_{k,21} & A_{k-1} \end{bmatrix},$$

where  $A_{k-1}$  stands for the finite element matrix on the coarse (vertex set) level. The relation between the two representations is  $\widehat{A}_k = J_k^T A_k J_k$ , where  $J_k = \begin{bmatrix} I_{k,1} & J_{k,12} \\ 0 & I_{k-1} \end{bmatrix}$ . Here  $I_{k,1}$ ,  $I_{k-1}$  denote identity matrices and  $J_{k,12}$  is an interpolation matrix from coarse to fine element mesh points basis functions. Hence,  $\widehat{A}_{k,12} = A_{k,12} + A_{k,11} J_{k,12}$ . Since the lower-right block of  $\widehat{A}_k$  equals  $A_{k-1}$ , it is called a two-level the HBF matrix. Another observation is that the two Schur complements  $S_{\widehat{A}_k} = A_{k-1} - \widehat{A}_{k,21} A_{k,11}^{-1} \widehat{A}_{k,12}$  and  $S_{A_k} = A_{k,22} - A_{k,21} A_{k,11}^{-1} A_{k,12}$  are identical, which is revealed by an elementary computation.

The following constant  $\gamma_k$ , which can be derived from the strengthened CBS inequality, plays a fundamental role in the AMLI method. It measures the strength of the off-diagonal blocks in  $\widehat{A}_k$  in relation to the main diagonal blocks and can be defined simply as

$$\gamma_k^2 = \rho(A_{k-1}^{-1} \widehat{A}_{k,21} A_{k,11}^{-1} \widehat{A}_{k,12}).$$

Next we introduce the exact factorization of  $A_k$ , namely,

$$A_k = \begin{bmatrix} A_{k,11} & 0 \\ A_{k,21} & I_{k-1} \end{bmatrix} \begin{bmatrix} I_{k,1} & A_{k,11}^{-1} A_{k,12} \\ 0 & S_{A_k} \end{bmatrix}. \quad (4)$$

To construct an efficient preconditioner to  $A_k$  we must first construct accurate approximations  $B_{k,11}$  to  $A_{k,11}$  and  $\widetilde{S}_{A_k}$  to  $S_{A_k}$ , the latter being a full matrix in general. Furthermore, the approximation of  $S_{A_k}$  must be done recursively. We recall that in the classical AMLI method, [11], this is done as follows.

Let  $P_\nu(t)$  be a given polynomial of degree  $\nu$ , such that  $0 \leq P_\nu(t) < 1$ ,  $0 < t \leq 1$  and is normalized,  $P_\nu(0) = 1$ . Further, let  $B_{k,11}$  be an approximation of  $A_{k,11}$ , which is spectrally related to it in a form to be shown later.

#### The AMLI method

Let  $M_0 = A_0$

For  $k = 1, 2, \dots$ , let

$$M_k = \begin{bmatrix} B_{k,11} & 0 \\ \widetilde{A}_{k,21} & I_{k,1} \end{bmatrix} \begin{bmatrix} I_{k,1} & B_{k,11}^{-1} \widetilde{A}_{k,12} \\ 0 & \widetilde{M}_{k-1} \end{bmatrix}$$

and  $\widetilde{M}_{k-1} = (I_{k-1} - P_\nu(M_{k-1}^{-1} A_{k-1})) A_{k-1}^{-1}$ .

It is seen that the preconditioner  $M_k$  to  $A_k$  is only implicitly (recursively) defined.

There are several possibilities how to define the off diagonal blocks. Besides  $\widetilde{A}_{k,12} = A_{k,12}$  and  $\widetilde{A}_{k,21} = A_{k,21}$ , which is used in the numerical experiments in Section 9, we can choose  $\widetilde{A}_{k,12} = \widehat{A}_{k,12}$  and  $\widetilde{A}_{k,21} = \widehat{A}_{k,21}$  or, alternatively,

$$\begin{aligned} \widetilde{A}_{k,12} &= A_{k,12} + (A_{k,11} - B_{k,11}) J_{k,12} \\ \widetilde{A}_{k,21} &= A_{k,21} + J_{k,12}^T (A_{k,11} - B_{k,11}) \end{aligned} \quad (5)$$

(see also Remark 2.1 below).

The reason for the last choice of perturbing the off-diagonal block-matrices, as is done in (5) is that in this way the HBF counterpart of  $M_k$ ,  $\widehat{M}_k = J_k^T M_k J_k$  takes the form

$$\widehat{M}_k = \begin{bmatrix} B_{k,11} & \widehat{A}_{k,12} \\ \widehat{A}_{k,21} & \widetilde{M}_{k-1} + \widehat{A}_{k,21} B_{k,11}^{-1} \widehat{A}_{k,12} \end{bmatrix}$$

which follows from an elementary computation. Hence,  $\widehat{M}_k$  can be considered as a preconditioner to  $\widehat{A}_k$  and since

$$\sup_{\mathbf{v}} \frac{\mathbf{v}^T A_k \mathbf{v}}{\mathbf{v}^T M_k \mathbf{v}} = \sup_{\widehat{\mathbf{v}}} \frac{\widehat{\mathbf{v}}^T \widehat{A}_k \widehat{\mathbf{v}}}{\widehat{\mathbf{v}}^T \widehat{M}_k \widehat{\mathbf{v}}} \quad \text{and} \quad \inf_{\mathbf{v}} \frac{\mathbf{v}^T A_k \mathbf{v}}{\mathbf{v}^T M_k \mathbf{v}} = \inf_{\widehat{\mathbf{v}}} \frac{\widehat{\mathbf{v}}^T \widehat{A}_k \widehat{\mathbf{v}}}{\widehat{\mathbf{v}}^T \widehat{M}_k \widehat{\mathbf{v}}} \quad (6)$$

the extreme eigenvalues of  $M_k^{-1} A_k$  equal those of  $\widehat{M}_k^{-1} \widehat{A}_k$ . Since the off-diagonal blocks in  $\widehat{M}_k$  equal those in  $\widehat{A}_k$ , the estimate of the extreme eigenvalues of  $\widehat{M}_k^{-1} \widehat{A}_k$  can be more readily done (see e.g. [12]).

The polynomial  $P_\nu$  is constructed as a shifted and scaled Chebyshev polynomial,

$$P_\nu(t) = \left( T_\nu \left( \frac{1 + \alpha - 2t}{1 - \alpha} \right) + 1 \right) / \left( T_\nu \left( \frac{1 + \alpha}{1 - \alpha} \right) + 1 \right),$$

where  $T_\nu(x) = \frac{1}{2}[(x + \sqrt{x^2 - 1})^\nu + (x - \sqrt{x^2 - 1})^\nu]$  is the Chebyshev polynomial of first kind and  $\alpha > 0$  is a lower bound of the eigenvalues of  $M_{k-1}^{-1}A_{k-1}$ . As is seen, the upper bound of these eigenvalues is bounded by the unit number.

**Remark 2.1** *The AMLI method can alternatively be defined from the HBF matrix. The matrices  $\widehat{A}_{k,12}$  and  $\widehat{A}_{k,21}$  are less sparse than  $A_{k,12}$  and  $A_{k,21}$ . Therefore, when an action of the HBF matrix on a vector is required, in practice one uses the transformation  $\widehat{A}_k = J_k^T A_k J_k$ . The formulation in (5) is simpler, however. When  $B_{k,11}$  is a sufficiently accurate approximation to  $A_{k,11}$  one may neglect the perturbation terms in (5). In Section 4 we extend the above results to more general congruence transformations of the form  $Z^T A Z$ .*

One can construct an AMLI method for general positive definite matrices, i.e., without assuming any underlying hierarchy of meshes and thus avoiding any (implicit or direct) transformation to a corresponding HBF block structure of the matrices. However, in order to construct an optimal order preconditioner for a sequence of approximations to elliptic boundary value problems, the approximation  $B_{k,11}$  to  $A_{k,11}$  must then be related to the Schur complement matrix  $S_{A_k}$  or its approximation  $S_{B_k} = A_{k,22} - A_{k,21}B_{k,11}^{-1}A_{k,12}$  in a certain way. As the multilevel extension can be done as shown above, it suffices to consider the two-level form of the method. Therefore, in the sequel, we omit the subscript  $k$ .

### 3 Condition number bounds for the two-level method

We consider now approximate factorizations of the same form as in (4). To derive the condition number of the corresponding two-level method, we make use of the following Lemma (see [12], [13]).

**Lemma 3.1** *Let  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$  and  $E = \begin{bmatrix} E_{11} & 0 \\ 0 & E_{22} \end{bmatrix}$  be symmetric matrices, where  $A_{11}$  is positive definite and  $E_{11}$  is positive semidefinite. Consider a transformation matrix  $K = \begin{bmatrix} I_1 & -A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}$ . Then the transformed matrices  $K^T A K$  and  $K^T E K$  take the following form:*

$$(i) \quad K^T A K = \begin{bmatrix} A_{11} & 0 \\ 0 & S_A \end{bmatrix}$$

$$(ii) \quad K^T E K = \begin{bmatrix} E_{11} & -E_{11}A_{11}^{-1}A_{12} \\ -A_{21}A_{11}^{-1}E_{11} & E_{22} + A_{21}A_{11}^{-1}E_{11}A_{11}^{-1}A_{12} \end{bmatrix}$$

The corresponding quadratic forms satisfy

$$(iii) \quad [\mathbf{v}_1, \mathbf{v}_2]^T K^T A K \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix} = \mathbf{v}_1^T A_{11} \mathbf{v}_1 + \mathbf{v}_2^T S_A \mathbf{v}_2$$

$$(iv) \quad [\mathbf{v}_1, \mathbf{v}_2]^T K^T E K \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix} \leq (1 + \xi) \mathbf{v}_1^T E_{11} \mathbf{v}_1 + (1 + \xi^{-1}) \mathbf{v}_2^T A_{21} A_{11}^{-1} E_{11} A_{11}^{-1} A_{12} \mathbf{v}_2 + \mathbf{v}_2^T E_{22} \mathbf{v}_2$$

for any  $\xi > 0$  and for all  $\mathbf{v}_1, \mathbf{v}_2$  of dimensions consistent with the block matrix partitioning of  $A$ .

*Proof* Parts (i)-(iii) follow by straightforward computations. Part (iv) follows from

$$\begin{aligned} & |\mathbf{v}_1^T E_{11} A_{11}^{-1} A_{12} \mathbf{v}_2 + \mathbf{v}_2^T A_{21} A_{11}^{-1} E_{11} \mathbf{v}_1| = 2 |\mathbf{v}_1^T E_{11} A_{11}^{-1} A_{12} \mathbf{v}_2| \\ & = 2 |(\mathbf{v}_1^{1/2})^T E_{11}^{1/2} A_{11}^{-1} A_{12} \mathbf{v}_2| \leq 2 \left\{ \mathbf{v}_1^T E_{11} \mathbf{v}_1 \mathbf{v}_2^T A_{21} A_{11}^{-1} E_{11} A_{11}^{-1} A_{12} \mathbf{v}_2 \right\}^{1/2} \\ & \leq \xi \mathbf{v}_1^T E_{11} \mathbf{v}_1 + \xi^{-1} \mathbf{v}_2^T A_{21} A_{11}^{-1} E_{11} A_{11}^{-1} A_{12} \mathbf{v}_2. \end{aligned}$$

■

Consider now a preconditioner to  $A$  in block-matrix factored form

$$B = \begin{bmatrix} B_{11} & 0 \\ A_{21} & I_2 \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1} A_{12} \\ 0 & S_B \end{bmatrix} \quad (7)$$

where  $B_{11}$  and  $S_B$  are approximations to  $A_{11}$  and  $S_A$ , respectively. Note that  $B = \begin{bmatrix} B_{11} & A_{12} \\ A_{21} & B_{22} \end{bmatrix}$ , where  $B_{22} = S_B + A_{21} B_{11}^{-1} A_{12}$ , and that  $S_B$  is the Schur complement of  $B$ .

Consider further the auxiliary matrix  $\tilde{A} = \begin{bmatrix} B_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ . We assume that  $B_{11}$  and  $S_B$  are spectrally equivalent to  $A_{11}$  and  $S_A$ , that  $\tilde{A}$  is positive semi-definite, and that the following inequalities hold:

$$\beta^{-1} \mathbf{v}_1^T A_{11} \mathbf{v}_1 \leq \mathbf{v}_1^T B_{11} \mathbf{v}_1 \leq \mathbf{v}_1^T A_{11} \mathbf{v}_1 \text{ for all } \mathbf{v}_1, \quad (8)$$

$$\eta^{-1} \mathbf{v}_2^T S_A \mathbf{v}_2 \leq \mathbf{v}_2^T S_B \mathbf{v}_2 \leq \mathbf{v}_2^T S_A \mathbf{v}_2 \text{ for all } \mathbf{v}_2, \quad (9)$$

where  $\beta \geq 1$  and  $\eta \geq 1$ .

**Theorem 3.1** *Let  $A$  be symmetric and positive definite. Assume that the approximations  $B_{11}$  and  $S_B$  are such that (8) and (9) hold. Then*

$$\kappa^{-1} \mathbf{v}^T A \mathbf{v} \leq \mathbf{v}^T B \mathbf{v} \leq 2 \mathbf{v}^T A \mathbf{v} \quad (10)$$

holds true for all  $\mathbf{v}$  with  $\kappa = \eta + \beta + \alpha - 1$  where

$$\alpha = \sup_{\mathbf{v}_2} \frac{\mathbf{v}_2^T A_{21} B_{11}^{-1} (A_{11} - B_{11}) B_{11}^{-1} A_{12} \mathbf{v}_2}{\mathbf{v}_2^T S_B \mathbf{v}_2}. \quad (11)$$

*Proof* Since  $B_{11} \leq A_{11}$ , it follows that  $A_{11}^{-1} \leq B_{11}^{-1}$ , so  $S_{\tilde{A}} = A_{22} - A_{21} B_{11}^{-1} A_{12} \leq A_{22} - A_{21} A_{11}^{-1} A_{12} = S_A$ . Further, it follows that

$$B = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} - \begin{bmatrix} A_{11} - B_{11} & 0 \\ 0 & S_{\tilde{A}} - S_B \end{bmatrix} \leq A + \begin{bmatrix} 0 & 0 \\ 0 & S_B - S_{\tilde{A}} \end{bmatrix}.$$

Since the lower diagonal block in  $A^{-1}$  equals  $S_A$ , it holds that

$$A^{-1/2} B A^{-1/2} \leq I + \begin{bmatrix} 0 & 0 \\ 0 & S_A^{-1/2} (S_B - S_{\tilde{A}}) S_A^{-1/2} \end{bmatrix},$$

that is,

$$\mathbf{v}^T A^{-1/2} B A^{-1/2} \mathbf{v} \leq \mathbf{v}^T \mathbf{v} + \mathbf{v}_2^T S_A^{-1/2} (S_B - S_{\tilde{A}}) S_A^{-1/2} \mathbf{v},$$

where  $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}$ . It holds that  $S_B - S_{\tilde{A}} = S_B - A_{22} + A_{21} B_{11}^{-1} A_{12} = S_B - S_A + A_{21} B_{11}^{-1} A_{12} - A_{21} A_{11}^{-1} A_{12}$ . Since, by assumptions made,  $S_B \leq S_A$  and  $A_{21} B_{11}^{-1} A_{12} \leq A_{22}$ , it follows that  $S_A - S_{\tilde{A}} \leq A_{22} - A_{21} A_{11}^{-1} A_{12} = S_A$ . Hence,  $A^{-1/2} B A^{-1/2} \leq \begin{bmatrix} I_1 & 0 \\ 0 & 2I_2 \end{bmatrix}$ , so  $\mathbf{v}^T B \mathbf{v} \leq 2\mathbf{v}^T A \mathbf{v}$ , which is the right-hand-side inequality in (10).

To prove the left-hand-side inequality, let  $K = \begin{bmatrix} I_1 & -B_{11}^{-1} A_{12} \\ 0 & I_2 \end{bmatrix}$ . It follows from Lemma 3.1 with  $E_{11} = A_{11} - B_{11}$ ,  $E_{22} = S_{\tilde{A}} - S_B$  that

$$\mathbf{v}^T K^T B K \mathbf{v} = \mathbf{v}_1^T B_{11} \mathbf{v}_1 + \mathbf{v}_2^T S_B \mathbf{v}_2$$

and

$$\begin{aligned} \mathbf{v}^T K^T (A - B) K \mathbf{v} &= \mathbf{v}_1^T (A_{11} - B_{11}) \mathbf{v}_1 + 2\mathbf{v}_1^T (I_1 - A_{11} B_{11}^{-1}) A_{12} \mathbf{v}_2 + \\ &\quad \mathbf{v}_2^T [(S_{\tilde{A}} - S_B) + A_{21} B_{11}^{-1} (A_{11} - B_{11}) B_{11}^{-1} A_{12}] \mathbf{v}_2. \end{aligned}$$

Here, for any  $\xi > 0$ ,

$$2|\mathbf{v}_1^T (I_1 - A_{11} B_{11}^{-1}) A_{12} \mathbf{v}_2| \leq \xi \mathbf{v}_1^T (A_{11} - B_{11}) \mathbf{v}_1 + \xi^{-1} \mathbf{v}_2^T A_{21} B_{11}^{-1} (A_{11} - B_{11}) B_{11}^{-1} A_{12} \mathbf{v}_2.$$

Further, we use  $S_{\tilde{A}} - S_B \leq S_A - S_B$ . Hence,

$$\begin{aligned} \frac{\mathbf{v}^T K^T (A - B) K \mathbf{v}}{\mathbf{v}^T K^T B K \mathbf{v}} &\leq \max \left\{ (1 + \xi) \sup_{\mathbf{v}_1} \frac{\mathbf{v}_1^T (A_{11} - B_{11}) \mathbf{v}_1}{\mathbf{v}_1^T B_{11} \mathbf{v}_1}, \right. \\ &\quad \left. \sup_{\mathbf{v}_2} \frac{\mathbf{v}_2^T (S_A - S_B) \mathbf{v}_2}{\mathbf{v}_2^T S_B \mathbf{v}_2} + (1 + \xi^{-1}) \sup_{\mathbf{v}_2} \frac{\mathbf{v}_2^T A_{21} B_{11}^{-1} (A_{11} - B_{11}) B_{11}^{-1} A_{12}}{\mathbf{v}_2^T S_B \mathbf{v}_2} \right\} \\ &= \max\{(1 + \xi)(\beta - 1), \eta - 1 + (1 + \xi^{-1})\alpha\}. \end{aligned}$$

To minimize the upper bound, let  $\xi$  satisfy  $(1 + \xi)(\beta - 1) = \eta - 1 + (1 + \xi^{-1})\alpha$ , i.e.,

$$\xi = \frac{\alpha + \eta - \beta}{2(\beta - 1)} \left[ 1 + \sqrt{1 + \frac{4\alpha(\beta - 1)}{(\alpha + \eta - \beta)^2}} \right]$$

Hence,

$$\begin{aligned} \frac{\mathbf{v}^T K (A - B) K \mathbf{v}}{\mathbf{v}^T K^T B K \mathbf{v}} &\leq (1 + \xi)(\beta - 1) = \frac{\alpha + \eta + \beta - 2}{2} \left[ 1 + \sqrt{1 - \left( \frac{4(\beta - 1)(\eta - 1)}{(\alpha + \beta + \eta - 2)^2} \right)} \right] \\ &\leq \alpha + \eta + \beta - 2. \end{aligned}$$

This implies

$$\frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{v}^T B \mathbf{v}} \leq \eta + \alpha + \beta - 1, \quad \text{for all } \mathbf{v},$$

where equality is taken when both  $\eta = 1$  and  $\beta = 1$ . ■

In our framework we can control the values of  $\eta$  and  $\beta$ . However, we see from the definition of  $\alpha$  (expression (11)) that it can take large values, in particular for vectors near the eigenvectors for eigenvalues of the order of  $\lambda_{\min}(S_B)$  when the quantity  $\mathbf{v}_2^T S_B \mathbf{v}_2 / \mathbf{v}_2^T \mathbf{v}_2$  is small. This occurs typically for "smooth" vectors, i.e., corresponding to the first harmonics of the matrix  $A$ .

An upper bound for  $\alpha$  can be derived using the CBS constant  $\tilde{\gamma}$  for  $B$ . Then one finds that  $\alpha \leq \frac{(\beta-1)\tilde{\gamma}^2}{1-\tilde{\gamma}^2}$ , see [13]. Since the matrices  $A_{12}$ ,  $A_{21}$  in  $B$  are derived from standard basis functions,  $\tilde{\gamma}$  can take values arbitrary close to unity, which also shows why  $\alpha$  can take big values.

To avoid the latter we must therefore somehow relate the approximation  $B_{11}$  of  $A_{11}$  to these vectors. This can be done by constructing  $B_{11}$  in such a way that the denominator in the expression (11) takes the value zero or very small values for such vectors. Such or similar analysis can be found in [18, 19] and [13]. In practice, this seems to work satisfactory. For the case where the opposite inequalities to (8), (9) hold, see [12].

A similar approach has actually appeared earlier in [9], but based on Stieltjes matrices. There, one lets  $\tilde{B}_{11} A_{11} \mathbf{v}_1 = \mathbf{v}_1$  for a positive vector  $\mathbf{v}_1$  such that  $A\mathbf{v} > \mathbf{0}$ ,  $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}$ , where  $\tilde{B}_{11}$  is a nonnegative approximation of  $A_{11}^{-1}$ . The corresponding approximate Schur complement is here taken as  $S_{\tilde{A}} = A_{22} - A_{21} \tilde{B}_{11} A_{12}$ , where  $\tilde{A} = \begin{bmatrix} \tilde{B}_{11}^{-1} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ . The approach in [18] is similar to this. Assuming that  $A_{11}$  is an M-matrix, in [13] two constructions of  $B_{11}$  are considered, namely as a modified incomplete factorization, where  $B_{11}$  preserves the action of  $A_{11}$  on a particular vector  $\mathbf{v}$  ( $B_{11}\mathbf{v} = A_{11}\mathbf{v}$ ) and as a compensated explicit inverse, where  $B_{11}^{-1}(A_{11}\mathbf{v}) = A_{11}^{-1}(A_{11}\mathbf{v})$ .

## 4 Improvements of the condition number by congruence transformations

Consider the following form of a congruence transformation of  $A$ ,  $\tilde{A} = Z^T A Z$ , where  $Z = \begin{bmatrix} I_1 & Z_{12} \\ 0 & I_2 \end{bmatrix}$ . Let  $Z_{21} = Z_{12}^T$  and  $\tilde{A} = \begin{bmatrix} \tilde{A}_{11} & \tilde{A}_{12} \\ \tilde{A}_{21} & \tilde{A}_{22} \end{bmatrix}$ . An elementary computation shows that  $\tilde{A}_{11} = A_{11}$ ,  $\tilde{A}_{12} = A_{12} + A_{11} Z_{12}$ ,  $\tilde{A}_{21} = \tilde{A}_{12}^T$  and  $S_{\tilde{A}} = \tilde{A}_{22} - \tilde{A}_{21} \tilde{A}_{11}^{-1} \tilde{A}_{12} = S_A$ .

Our aim is to choose  $Z_{12}$  such that the off-diagonal blocks of  $\tilde{A}$  have a weaker influence on the condition number in (10) than the corresponding blocks in  $A$  have, i.e., to ensure a smaller value of  $\kappa$  by decreasing the value of the coefficient  $\alpha$  in (11). To achieve this the matrix  $Z_{12}$  must be chosen properly, while still preserving the sparsity of the matrices involved.

Matrix  $\tilde{A}$  can then be preconditioned as in (7), i.e., with

$$\tilde{B} = \begin{bmatrix} I_1 & 0 \\ \tilde{A}_{21} B_{11}^{-1} & I_2 \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ 0 & S_B \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1} \tilde{A}_{12} \\ 0 & I_2 \end{bmatrix},$$

where  $B_{11}$  is an approximation of  $\tilde{A}_{11}$  and  $S_B$  an approximation of  $S_{\tilde{A}}$  which latter, as we have seen, equal  $A_{11}$  and  $S_A$ , respectively.

However, instead of performing an iterative method, such as the conjugate gradient method, with  $\tilde{B}$  as a preconditioner to  $\tilde{A}$ , we show now that it can be more efficient to use the preconditioner  $Z^{-T}\tilde{B}Z^{-1}$  to  $A$ .

We note first that the extreme values of

$$\inf_{\mathbf{y}} / \sup_{\mathbf{y}} \frac{\mathbf{y}^T \tilde{A} \mathbf{y}}{\mathbf{y}^T \tilde{B} \mathbf{y}} \quad \text{and} \quad \inf_{\mathbf{x}} / \sup_{\mathbf{x}} \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T Z^{-T} \tilde{B} Z^{-1} \mathbf{x}}$$

are equal, which follows from a transformation  $\mathbf{x} = Z\mathbf{y}$ . Therefore, the extreme eigenvalues of  $\tilde{B}^{-1}\tilde{A}$  and  $(Z^{-T}\tilde{B}Z^{-1})^{-1}A$  are equal and, hence, the upper bound for the rate of convergence of the conjugate gradient method is unchanged. Next we simplify the matrix  $Z^{-T}\tilde{B}Z^{-1}$ . It holds

$$Z^{-T}\tilde{B}Z^{-1} = Z^{-T} \begin{bmatrix} I_1 & 0 \\ \tilde{A}_{21}B_{11}^{-1} & I_2 \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ 0 & S_B \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1}\tilde{A}_{12} \\ 0 & I_2 \end{bmatrix} Z^{-1}$$

where

$$\begin{bmatrix} I_1 & B_{11}^{-1}\tilde{A}_{12} \\ 0 & I_2 \end{bmatrix} Z^{-1} = \begin{bmatrix} I_1 & B_{11}^{-1}\tilde{A}_{12} \\ 0 & I_2 \end{bmatrix} \begin{bmatrix} I_1 & -Z_{12} \\ 0 & I_2 \end{bmatrix} = \begin{bmatrix} I_1 & B_{11}^{-1}(A_{12} + A_{11}Z_{12}) - Z_{12} \\ 0 & I_2 \end{bmatrix} = \begin{bmatrix} I_1 & B_{11}^{-1}[A_{12} + (A_{11} - B_{11})Z_{12}] \\ 0 & I_2 \end{bmatrix}.$$

Therefore, it holds that

$$Z^{-T}\tilde{B}Z^{-1} = \begin{bmatrix} I_1 & 0 \\ \hat{A}_{21}B_{11}^{-1} & I_2 \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ 0 & S_B \end{bmatrix} \begin{bmatrix} I_1 & B_{11}^{-1}\hat{A}_{12} \\ 0 & I_2 \end{bmatrix},$$

where  $\hat{A}_{12} = A_{12} + (A_{11} - B_{11})Z_{12}$  and  $\hat{A}_{21} = \hat{A}_{12}^T$ .

We consider now two choices of  $Z_{12}$  to improve the value of  $\alpha$  in (11).

- (i)  $Z_{12} = J_{12}$ , the interpolation matrix between hierarchical and standard basis functions in a finite element method. This is the choice considered in [11]. It can improve the value of  $\alpha$  dramatically, for instance from  $O(h^{-1})$  to  $1/(1 - \gamma^2) < 4$ , which latter holds for piece-wise linear basis functions and a uniform subdivision of triangular elements in four congruent triangles.
- (ii)  $Z_{12} = -H_1 A_{12}$ , where  $H_1$  is an approximation of  $A_{11}^{-1}$ . This makes  $\tilde{A}_{12} = A_{12} + A_{11}Z_{12}$  small.

If  $H_1 = A_{11}^{-1}$ , then  $\tilde{A}_{12} = 0$ . This choice is infeasible, however, as  $A_{11}^{-1}$  is a full matrix. Even very accurate approximations may involve too much loss of sparsity and increase of computational complexity. Instead we aim at choosing  $H_1$  such that  $\alpha$  is significantly reduced while matrix-vector multiplications with  $H_1$  can still be

performed cheaply. One possible choice is to let  $H_1$  be an extension of the inverse of the matrix  $A_{11,H}$  corresponding to a coarser mesh. The smallest eigenvalues of the latter matrix are good approximations of those of  $A_{11} = A_{11,h}$ , and the value of  $\alpha$  in (11) is therefore much reduced. Furthermore, by using the matrix graph of  $A$ , this method can often be extended also to other applications than those arising from finite elements.

**Remark 4.1** *If  $H_1$  is such that the action of  $A_{21}$  on "smooth" vectors is small, it may suffice to use a block-diagonal preconditioner to  $A$ , where the diagonal blocks approximate  $A_{11}$  and  $S_A$ , respectively. It is then assumed that an action of  $H_1$  is relatively cheap but the approximation  $B_{11}$  of  $A_{11}$  can be allowed to require more computational effort, as it is now applied only once at each iteration step.*

## 5 Element-by-element Schur complement approximation

While in the previous sections we introduced a general AMLI framework, we proceed now with a particular approximation to  $S_A$ . This approximation uses (macro)elementwise approach and the resulting AMLI variant will be denoted as AMLI-ES in the sequel. Such an approach has been previously considered in [17], but in a more limited context and with less general proofs of the resulting condition numbers.

Let us reconsider a sequence of nested meshes  $\mathcal{T}_k$  and an approximation to the Schur complement  $S_A = S_{A_k}$ . Let  $E \in \mathcal{T}_{k-1}$  be a macroelement unifying several elements  $e \in \mathcal{E}_E \subset \mathcal{T}_k$ . Then a macroelement matrix  $A_E = A_{k,E}$  arises from assembling the element matrices  $A_e$ ,  $e \in \mathcal{E}_E$ . This matrix can be partitioned and transformed to the hierarchical basis (as in Section 2),

$$A_E = \begin{bmatrix} A_{E,11} & A_{E,12} \\ A_{E,21} & A_{E,22} \end{bmatrix}, \quad \widehat{A}_E = \begin{bmatrix} A_{E,11} & \widehat{A}_{E,12} \\ \widehat{A}_{E,21} & \widehat{A}_{E,22} \end{bmatrix}$$

where  $\widehat{A}_{E,22}$  equals the element matrix  $A_{k-1,E}$ .

Then let

$$S_E = A_{E,22} - A_{E,21}A_{E,11}^{-1}A_{E,12}$$

be the local Schur complement for  $E$  and  $\widetilde{S}_A$  be the element-Schur approximation to  $S_A$ , namely,

$$\widetilde{S}_A = \sum_{E \in \mathcal{T}_{k-1}} R_E^T S_E R_E,$$

where  $R_E$  denotes a Boolean matrix representing the restriction to the macroelement degrees of freedom.

Let  $A_E = A_{k,E}$  denote the macroelement matrix. Then the following inequality holds for the local Schur complements  $S_E$  on each element  $E \in \mathcal{T}_{k-1}$ ,

$$(1 - \gamma_E^2)A_E \leq S_E, \tag{12}$$

where  $\gamma_E = \rho(A_E^{-1}\widehat{A}_{E,21}A_{E,11}^{-1}\widehat{A}_{E,12})^{1/2}$  is the local value of the CBS constant, see e.g. [11]. This inequality follows directly from the representation in the local hierarchical basis functions and the fact that the Schur complements  $S_E$  and  $\widehat{S}_E = A_{k-1,E} - \widehat{A}_{E,21}A_{E,11}^{-1}\widehat{A}_{E,12}$  for the standard and hierarchical basis function matrices are identical,  $S^{(\ell)} = \widehat{S}^{(\ell)}$ , see [11].

The following spectral bounds hold.

**Theorem 5.1** *Let  $\widetilde{S}_A$  be the assembly of the local Schur complements  $S_E$  and let  $S_A = A_{22} - A_{21}A_{11}^{-1}A_{12}$  be the global Schur complement matrix. Then*

$$(1 - \gamma^2)S_A \leq \widetilde{S}_A \leq S_A, \quad (13)$$

where  $\gamma = \max_E \gamma_E$ .

*Proof* Let  $A_H$  be the matrix corresponding to the finite element stiffness matrix on the discrete mesh  $\mathcal{T}_H$ . Since  $A_H$  and  $\widetilde{S}_A$  are the assembly of  $A_E$  and  $S_E$ , respectively, it follows from (12) that  $(1 - \gamma^2)A_H \leq \widetilde{S}_A$  which, combined with the bound  $A_H \geq S_A$ , is the left-hand-side inequality in (13).

The right-hand-side inequality in (13) follows directly from a general property which holds for Schur complements (see e.g. [1] (Theorem 3.8) and [17]). Let  $\mathbf{v}_E = \begin{bmatrix} \mathbf{v}_{E,1} \\ \mathbf{v}_{E,2} \end{bmatrix}$ ,  $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}$  be the corresponding partitioning of local and global finite element vectors. Then

$$\mathbf{v}_{E,2}^T S_E \mathbf{v}_{E,2} = \min_{\mathbf{v}_{E,1}} \mathbf{v}_E^T A_E \mathbf{v}_E.$$

Further, because  $S_E$  are local matrices,

$$\mathbf{v}^T \widetilde{S}_A \mathbf{v} = \mathbf{v}^T \sum_E R_E^T S_E R_E \mathbf{v} = \sum_E \mathbf{v}_{E,2}^T S_E \mathbf{v}_{E,2}$$

Finally,

$$\sum_E \min_{\mathbf{v}_{E,1}} \mathbf{v}_E^T A_E \mathbf{v}_E \leq \min_{\mathbf{v}_1} \mathbf{v}^T A \mathbf{v} = \mathbf{v}_2^T S_A \mathbf{v}_2,$$

where the inequality follows from the fact that the minimum on the left-hand-side is attained over a larger set of degrees of freedom than on the right-hand-side. Collecting these results we obtain that  $\widetilde{S}_A \leq S_A$ .  $\blacksquare$

We note that  $\widetilde{S}_A$  replaces  $S_B$  when applying Theorem 3.1. It is seen from Theorem 3.1 and (15) that the two-level spectral condition number of the preconditioning of  $S_A$  by  $\widetilde{S}_A$  satisfies

$$\text{cond}(\widetilde{S}_A^{-1} S_A) \leq 1/(1 - \gamma_H^2) \leq m^2 \quad (14)$$

**Remark 5.1** *The proof of the lower bound in Theorem 5.1 is more straightforward and more general than the corresponding bound in [17]. Furthermore, it shows that universal bounds hold for arbitrary coefficients in the differential operators (if they are piecewise constant) and even for degenerate elements.*

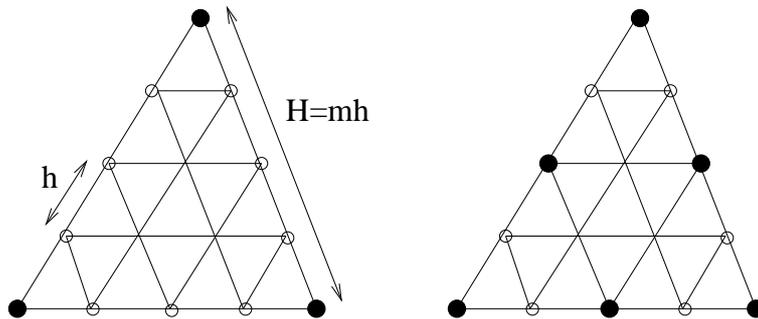


Figure 1: Two distributions of vertex and edge points (normal and composed macroelement)

**Remark 5.2** Clearly the method above is applicable both for triangular and for quadrilateral meshes. In the former case,  $\gamma^2 < 3/4$  for all triangles, even for degenerate ones (when one or two angles equal zero). The constant for an anisotropic operator for bilinear elements on a rectangular mesh is also  $\gamma^2 < 3/4$ , see [20] for a derivation.

**Remark 5.3** The preconditioner  $B_{11}$  to the block matrix  $A_{11}$  in (2) can be constructed in various ways. For the Laplacian operator and a nondegenerate mesh the matrix can be approximated by a diagonal matrix, see [1, 6, 2]. For more general operators such as anisotropic Laplacian and/or nearly degenerate meshes special, more robust techniques must be applied, see e.g. [1, 8]. For systems of partial differential operators one can frequently apply preconditioners based on a separate displacement ordering, see e.g. [16, 3].

The universal bound of  $\gamma^2$  for the  $m$ -partitioning of a general triangle and a general selfadjoint operator  $\mathcal{L}$  are

$$\gamma^2 < (m^2 - 1)/m^2, \quad m = 2, 3, \dots \quad (15)$$

see [5]. As is seen from (15), the value of  $\gamma^2$  approaches 1 when  $m$  increases, which is less favorable for the spectral bound in (3).

There are two possible choices how to handle the edge points in an  $m$ -times uniformly refined mesh, illustrated in Figure 1. In the first case (left), each node point on an edge is treated as an edge-point and in the second case (right), each second node-point is treated as an edge-point and the remaining ones – as vertex points. The advantage with the latter choice is that the value of  $\gamma^2$  will attain the universal bound for  $m = 2$ , namely,  $\gamma^2 < 3/4$ .

The two-level method considered above can be applied recursively as in the AMLI methods, see Section 4 and [11]. Thereby we first assemble the current local Schur complements and partition the assembled matrix in  $2 \times 2$  block form. The corresponding Schur complement matrix is again approximated by the assembly of the new local Schur complement matrices, now corresponding to a new (coarse mesh) partitioning of the mesh, which are constructed by first partitioning them in corresponding  $2 \times 2$  block forms.

The same types of approximations of the new  $B_{11}$  and  $S_A$  matrices are used. However, the new value of the CBS constant is not available because the coarse mesh matrix  $A_H$  corresponding to the new, coarser mesh is not available. One can say that the local Schur complement matrices correspond to a finite element methods for new, but unknown, basis functions. Hence, the progress of the method in a multilevel framework can only be shown experimentally.

***Remark 5.4*** In [15] the element-by-element Schur complement approximation technique was successfully applied to a nonsymmetric problem of saddle point form. It was first used to approximate the Schur complement of the saddle point matrix and second, to construct an AMLI preconditioner to the main pivot block and to the Schur complement itself.

## 6 A further improvement of the Schur complement preconditioner

The spectral bounds in Theorem 5.1 show that the eigenvalues of  $\tilde{S}_A^{-1}S$  are located in the interval  $[1, 1/(1 - \gamma_H^2)]$ . It is known that the conjugate gradient method converges in general faster for eigenvalues clustered on both sides of the origin, see e.g. [1]. Such a preconditioner can readily be constructed by taking a linear combination of the coarse mesh matrix  $A_H$  and  $\tilde{S}_A$ . Hence, let

$$C = \xi A_H + (1 - \xi)\tilde{S}_A, \quad 0 \leq \xi \leq 1$$

be a preconditioner to  $S_A$ . We derive below bounds for the eigenvalues of  $C^{-1}S$ .

Note first that on each element,

$$\begin{aligned} C^{(\ell)} &= \xi A_H^{(\ell)} + (1 - \xi)(A_H^{(\ell)} - \hat{A}_{21}^{(\ell)} A_{11}^{(\ell)-1} \hat{A}_{12}) \\ &\geq (\xi + (1 - \xi)(1 - \gamma_H^2))A_H^{(\ell)} \end{aligned}$$

Here the notation " $\hat{\phantom{x}}$ " indicates that the corresponding matrix blocks are in a HBF form. Hence

$$\begin{aligned} C &\geq (\xi + (1 - \xi)(1 - \gamma_H^2))A_H \\ &\geq (\xi + (1 - \xi)(1 - \gamma_H^2))S_A, \quad 0 \leq \xi \leq 1. \end{aligned}$$

Furthermore,

$$\begin{aligned} C &= \xi A_H + (1 - \xi)\tilde{S}_A \leq \left( \frac{\xi}{1 - \gamma_H^2} + 1 - \xi \right) S_A \\ &= \frac{\xi + (1 - \xi)(1 - \gamma_H^2)}{1 - \gamma_H^2} S_A. \end{aligned}$$

The eigenvalues of  $C^{-1}S_A$  are therefore contained in the interval

$$\left[ \frac{1}{1 - \xi + \frac{\xi}{1 - \gamma_H^2}}, \frac{1}{\xi + (1 - \xi)(1 - \gamma_H^2)} \right].$$

The upper bound for the condition number is still  $1/(1 - \gamma_H^2)$ , but for  $0 < \xi < 1$  the eigenvalues are now located on both sides of the unity.

Numerical illustrations show that the number of conjugate gradient iterations now can decrease somewhat for such values of  $\xi$  compared to the value  $\xi = 0$  when  $C = \tilde{S}_A$ .

## 7 Element-by-element pivot block approximation

One possible choice for the pivot block preconditioner is the following

$$B_{k,11}^{-1} = \sum_{E \in \mathcal{T}_{k-1}} R_E^T A_{E,11}^{-1} R_E,$$

where  $R_E$  denote a Boolean matrix representing the restriction to the macroelement degrees of freedom. Clearly, this type of preconditioners fits well in the concept of elementwise approximation of the Schur complement matrix.

Some basic spectral information about both  $A_{k,11}$  and  $B_{k,11}$  can be easily obtained from the spectral analysis of the macroelement matrices  $A_{E,11}$ . Let

$$0 < \lambda_0 \leq \lambda_{\min}(A_{E,11}) \text{ and } \lambda_{\max}(A_{E,11}) \leq \lambda_1 \text{ for all } E \in \mathcal{T}_{k-1}.$$

Then, for all appropriate vectors  $\mathbf{x}$ ,

$$\langle A_{k,11}x, x \rangle = \sum_E \langle A_{E,11}R_E x, R_E x \rangle \leq \lambda_1 \sum_E \langle R_E x, R_E x \rangle \leq 2\lambda_1 \langle x, x \rangle,$$

$$\langle A_{k,11}x, x \rangle = \sum_E \langle A_{E,11}R_E x, R_E x \rangle \geq \lambda_0 \sum_E \langle R_E x, R_E x \rangle \leq \lambda_0 \langle x, x \rangle.$$

The factor two in the upper estimate comes from the fact that the ‘‘fine’’ degrees of freedom can be shared by (at most) two macroelements.

Similarly, we get bounds for the preconditioner,

$$\langle B_{k,11}^{-1}x, x \rangle \leq \frac{2}{\lambda_0} \langle x, x \rangle \text{ and } \langle B_{k,11}^{-1}x, x \rangle \geq \frac{1}{\lambda_1} \langle x, x \rangle,$$

$$\frac{\lambda_0}{2} \langle x, x \rangle \leq \langle B_{k,11}x, x \rangle \leq \lambda_1 \langle x, x \rangle.$$

Combined, it gives

$$\varkappa^{-1} \langle B_{k,11}x, x \rangle \leq \langle A_{k,11}x, x \rangle \leq 4\varkappa \langle B_{k,11}x, x \rangle$$

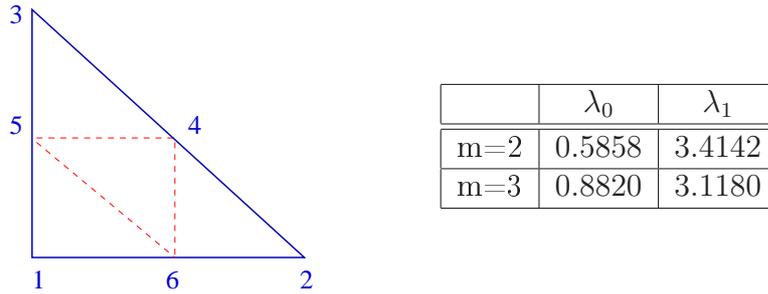


Figure 2: A macroelement

for  $\varkappa = \lambda_1/\lambda_0$ . However, this estimate is very pessimistic (it is worse than what we already got for the pivot block approximation by the identity matrix). Thus some other analysis, as e.g. the domain decomposition point of view, is necessary for better understanding and possible improvement of the preconditioner.

We comment now briefly on the values of  $\lambda_0$  and  $\lambda_1$ . They are easily computable in the case of Laplace differential operator and isosceles rectangular triangles, see Figure 2. In this case,

$$A_{E,11} = \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix}$$

and a computation of the eigenvalues shows that  $\lambda_0 = 2 - \sqrt{2} = 0.5858$  and  $\lambda_1 = 2 + \sqrt{2} = 3.4142$ .

The dependence of  $\lambda_0$  and  $\lambda_1$  on the parameter  $m$  is also illustrated in Figure 2. The improvement of the condition number with increasing  $m$  is in accordance with the domain decomposition point of view, as the condition number decreases with decreasing number of macroelements (subdomains).

A numerical computation reveals that  $\varkappa = \lambda_1/\lambda_0$  deteriorates with anisotropy in the case of a generalized anisotropic Laplacian or e.g. increasing Poisson ratio in the case of elasticity.

## 8 Implementation and computational complexity

We consider below triangular meshes. Dividing each element in  $m^2$  subelements can result in substantial savings of computer time when  $m$  is big, both on sequential and parallel computers, since much computational work can be done locally. Larger values of  $m$  also reduce the required number of levels, which is advantageous as the matrices  $A_k$  must be assembled in advance and stored.

For the outer iteration method we note however, that the value of  $\gamma$  increases with the number of vertex points  $m$ . The numerical results indicate, however, that the condition number bound for the method using local Schur complement approximations may be more accurate than the upper bound in (14) would imply.

We now show how to optimize the value of  $m$  in order to minimize the computational complexity for the iterative method. Within the AMLI methods, described in Section 2, a solution of a system with the preconditioner on level  $k$ , i.e.,  $M_k \mathbf{g} = \mathbf{r}$  can be implemented in the following steps:

$$\begin{aligned} \text{Step 1: } \mathbf{g}_1 &\leftarrow B_{k,11}^{-1} \mathbf{r}_1 \\ \text{Step 2: } \mathbf{g}_2 &\leftarrow \widetilde{M}_{k-1}^{-1} (\mathbf{r}_2 - \widetilde{A}_{k,21} \mathbf{g}_1) \\ \text{Step 3: } \mathbf{g}_1 &\leftarrow B_{k,11}^{-1} (\mathbf{r}_1 - \widetilde{A}_{k,12} \mathbf{g}_2) \end{aligned}$$

More generally, Step 2 solves the system  $A_{k-1} \mathbf{g}_2 = \mathbf{r}_2 - \widetilde{A}_{k,21} \mathbf{g}_1$  by  $\nu$  inner iterations with a proper preconditioner. If  $\nu$  is sufficiently big, typically  $\nu \geq (1 - \gamma^2)^{-1/2}$ , then the corresponding preconditioner is spectrally equivalent to the given matrix, see [13]. Since  $\gamma^2 < (m^2 - 1)/m^2$  (see [5]), it must then hold that  $m \leq \nu$ ,  $m = 2, 3, \dots$ . In practice we choose  $\nu$  as small as possible, i.e.,  $\nu = m$ .

Now, we are interested in the complexity  $W_k = w(M_k)$ , which after a simplification fulfills the following recurrence relation,

$$W_k = 2w(B_{k,11}^{-1}) + \nu W_{k-1}, \quad (16)$$

where  $w(B_{k,11}^{-1})$  is the complexity of the pivot block preconditioner. Note that in (16), we neglect the multiplications with the off-diagonal blocks  $\widetilde{A}_{k,21}$  and  $\widetilde{A}_{k,12}$  and neglect the work in the inner iterations, which is additional to the  $M_{k-1}$  preconditioner. The complexity can be further discussed from several points of view:

- (A) The first point of view is usually used in the multigrid theory. One assumes then that the coarsest grid  $\mathcal{T}_0$  is very coarse and fixed. Then one considers an arbitrary grid  $\mathcal{T}_k$ , which is a refinement of the coarsest grid. If  $n_k$  is the order of  $A_k$  (or, equivalently, number of the degrees of freedom of that level) then we can claim that the complexity of  $M_k$  is of optimal order if there is a  $k$ -independent constant  $C$  such that  $W_k \leq Cn_k$ . Let for  $k = 1, 2, \dots$ , it holds that  $n_k \leq \rho n_{k-1}$ ,  $\nu\rho < 1$  and there is a constant  $C$  such that  $w(B_{k,11}^{-1}) \leq Cn_k$  and  $W_0 = w(M_0) = w(A_0^{-1}) \leq Cn_0$ . Then (16) induces directly the optimal complexity of the multilevel preconditioners.

For 2D problems and macroelement-wise approximation to both the Schur complement and the pivot block (see Sections 5 and 7) all the above requirements for optimal complexity are fulfilled,  $\nu \sim m$  and  $\rho = m^{-2}$  for  $m \geq 2$ . In this respect, we can even afford bigger values of  $m_k = h_k/h_{k-1}$  which are advantageous as was already mentioned at the beginning of this section.

- (B) The second point of view, which is taken in this paper, respects the fact that a finer coarsest level is advantageous as well. It permits a triangulation which can better model more complex domains and a construction of the coarsest mesh such that no discontinuities occur within a macroelement, which latter is also favourable for the condition number of the preconditioned matrix.

**Remark 8.1** As has been shown in [7], an alternative technique to avoid carrying out the recursion in a multilevel method to a coarse mesh with a fixed number of degrees of freedom can be based on permutations of the local finite element matrices by addition of zero-order terms.

Let the coarsest level correspond to a finite element mesh with a characteristic size  $h_0$  and  $n_0 = O(h_0^{-2})$ . If the system corresponding to the coarsest level matrix is solved by a solution method, such as the unpreconditioned CG or a nested dissection method, then it suffices with  $O(h_0^{-3})$  operations. Thus, there is no asymptotic complexity increase if this  $O(h_0^{-3})$  complexity is comparable with  $O(h^{-2})$ , where  $h$  is the finest grid size.

To analyse the computational complexity of the method, let us start with such a finer and more flexible coarsest grid. Then it is a question how to develop the refinement process to get both good efficiency and complexity for the AMLI preconditioners. The answer to this question needs a more precise formulation and therefore, we assume that a macroelement-wise approximation is used for both  $A_{k,11}^{-1}$  and  $\tilde{S}_{A_k}$ . Then there is an initial cost in forming the Schur complements and in the computation of  $A_{E,11}^{-1}$  for each macroelement. On the  $k$ th level, this initial cost is  $O(h_{k-1}^{-2}m_k^6)$ , where  $m_k$  is the mesh size refinement factor,  $h_k = h_{k-1}/m_k$  and the number of elements for the  $k$ -th level is  $O(h_k^{-2})$ .

At each iteration step we perform matrix-times-vector multiplications with the exact local inverses and the major cost there is  $w(A_{E,11}^{-1}) = O(m_k^4)$ , arising from the matrices of order  $m_k^2 \times m_k^2$ .

From (14) it follows that  $O(m_k)$  inner iterations should take place at level  $k$ . Ignoring the constants involved, the computational complexity of the  $k$ -th level preconditioner is then as follows,

$$W_{k+1} = h_k^{-2}m_{k+1}^4 + m_{k+1}W_k, \quad k = 0, 1, \dots,$$

where  $W_0 = h_0^{-3}$  and  $h_k^{-2}$  is the number of macroelements on level  $k$ .

Let  $W_k = h_k^{-2-\alpha_k}$ ,  $0 < \alpha_k \leq 1$ ,  $\alpha_0 = 1$ , where  $\alpha_k$  measures the discrepancy in the order of computational complexity between the method and a method of optimal order at that level. Then

$$W_{k+1} = h_k^{-2} (m_{k+1}^4 + m_{k+1}h_k^{-\alpha_k}). \quad (17)$$

To balance the two terms above, we let

$$m_{k+1} = h_k^{-\alpha_k/3}. \quad (18)$$

Then

$$W_{k+1} = h_k^{-2-\frac{4}{3}\alpha_k} = h_{k+1}^{-2-\alpha_{k+1}} = (h_k/m_{k+1})^{-2-\alpha_{k+1}} = h_k^{(1+\frac{1}{3}\alpha_k)(-2-\alpha_{k+1})}.$$

A comparison of the exponents of  $h_k$  gives immediately that

$$\alpha_{k+1} = \frac{2}{3} \frac{\alpha_k}{1 + \alpha_k/3}. \quad (19)$$

Thus, it is seen that  $\alpha_k$  decreases monotonically to zero, and geometrically with a rate  $2/3$ . Accordingly, the complexity of multilevel preconditioner approaches the optimal order.

$h_0 =$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$
$m_1$	2	4	9	21
$m_2$	2	2	4	7
$m_3$	2	2	2	3
$m_4$	2	2	2	2

$h_0 =$	$10^{-1}$	$10^{-2}$	$10^{-3}$	$10^{-4}$
$m_1$	4	21	99	464
$m_2$	2	7	21	59
$m_3$	2	3	7	15
$m_4$	2	2	3	6

Table 1: The refinement for  $\alpha_0 = 1$  (left) and  $\alpha_0 = 2$  (right).

Simultaneously, the refinement factor  $m_k$  fulfills the relation

$$\begin{aligned}
m_{k+1} &= h_k^{-\alpha_k/3} = (h_{k-1}/m_k)^{-\alpha_k/3} = h_{k-1}^{-(1+\frac{\alpha_{k-1}}{3})\frac{\alpha_k}{3}} \\
&= m_k^{(1+\frac{\alpha_{k-1}}{3})\frac{\alpha_k}{\alpha_{k-1}}} = m_k^{2/3}.
\end{aligned} \tag{20}$$

If  $\alpha_0 = 1$ , then all  $\alpha_k \leq 1$  and

$$m_{k+1} \leq m_k^{2/3} \leq m_k.$$

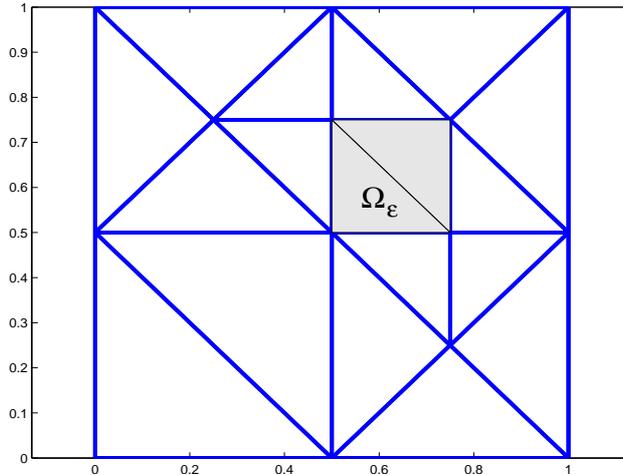
Thus  $m_k$  does not grow and therefore the condition number in (14) does not deteriorate. For  $\alpha_0 > 1$  the same happens after some steps. Thus with an increasing number of levels  $k \rightarrow \infty$ , we again approach an optimal order method.

Starting with  $\alpha_0 = 1$ , it follows from (19), that  $\alpha_1 = 0.5$ ,  $\alpha_2 = 0.2857$ ,  $\alpha_3 = 0.1739$  etc. The corresponding refining factors depend on the mesh size of the coarsest mesh. The denser the coarsest mesh is, the larger is the value of  $m$  that we can afford, see e.g. the values in Table 1.

The construction allows to compute an approximate solution at each such level and use it as an initial solution for the next refinement step. In this way one avoids the  $\log(h^{-1})$  factor in the iteration count which otherwise appears as the required relative accuracy for stopping the iterations would be  $O(h^2)$  (see e.g. [1]). A further advantage of the method is that if  $m$  is fixed one can extrapolate the approximate solutions to get a higher order of the approximation error than the standard  $O(h^2)$ . This requires however sufficient regularity of the exact solution of the boundary value problem.

For some of the levels, it is also possible to do the refinement locally and refine only some selected macroelements. To preserve conformity of the arising finite element mesh, the values in the arising hanging nodes can be e.g. interpolated. The arising composite grid finite element problem can then be solved iteratively. The convergence of whole iterative process depends again on the strengthened CBS constant. The subproblems corresponding to a coarser global grid and refined local grid can be solved iteratively with the use of the presented AMLI preconditioners.

Furthermore, the method offers a large amount of efficient parallel computation at each level. If the number of processors ( $p$ ) is  $p \leq h_0^{-2}$ , then each processor will be active at any computational step, except possibly when solving the final, coarsest level matrix by a direct solution method or by a preconditioned or unpreconditioned conjugate gradient



method. The efficiency of the latter step depends on the parallel computer platform, which is utilized.

Note that the presented complexity considerations are fully justified for standard AMLI preconditioners. They are also successfully applicable to AMLI-ES (as we can see in Section 9), however then we don't have sharp theoretical estimates for the constant  $\gamma$  on the coarser levels (see also Section 5).

## 9 Numerical illustrations

The iterative solvers used for the numerical tests are the Generalized Conjugate Gradient - Minimal Residual (GCG) method, described, for instance in [1] and the standard preconditioned Conjugate Gradient (PCG) method.

The performance of the preconditioning technique is illustrated on the following three test problems.

**Problem 9.1** We consider a Poisson problem with discontinuous coefficients  $k$

$$\nabla(k\nabla u) = f \quad \text{in } \Omega \subset \mathbb{R}^2,$$

where  $u = u(x, y)$ ,  $k = \varepsilon$  in  $\Omega_\varepsilon$  and  $k = 1$  elsewhere. The geometry of  $\Omega$  and the initial (coarsest) triangulation are shown in Figure 9, where  $\Omega_\varepsilon$  occupies the shaded region. The parameter  $m$  which relates  $H$  and  $h$  is varying as  $m = 2^s$ ,  $s = 1, 2, 3$ .

The results are shown in Table 2. For these experiments the coarsest level for the multilevel preconditioner is always the coarsest possible. The positions in the table, marked by "-", correspond to problem sizes, which cannot be obtained with that particular choice of the parameter  $m$ . We can see that the convergence of the method is not affected by the coefficient jumps, which in this case are aligned with the coarse mesh. It is also seen that the iteration count improves significantly for larger values of  $m$ .

Problem size	$\varepsilon_1 = 1$			$\varepsilon_1 = 10^{-3}$			$\varepsilon_1 = 10^3$		
	$m = 2$	$m = 4$	$m = 8$	$m = 2$	$m = 4$	$m = 8$	$m = 2$	$m = 4$	$m = 8$
161	8	6	-	8	6	-	8	6	-
609	11	-	6	11	-	6	11	-	6
2369	14	11	-	14	11	-	14	11	-
9345	18	-	-	18	-	-	18	-	-
37121	23	16	13	23	17	13	23	16	13

Table 2: Problem 9.1: Iteration counts for varying  $m$

**Problem 9.2** To illustrate the results in Section 6 we consider a two-dimensional anisotropic problem

$$-\varepsilon_1 u_{xx} - \varepsilon_2 u_{yy} = f \quad \text{in } \Omega \subset \mathbb{R}^2,$$

where  $u = u(x, y)$  and  $\Omega = [0, 1]^2$ .

The results are shown in Table 3. There we show runs for a sequence of problem sizes. The total possible number of levels is indicated in the row below the problem size. For  $\alpha = 0, 0.1, 0.5, 0.9, 1$  and  $\varepsilon = 10^{-3}, 1, 10^3$  we show the iteration counts of the standard conjugate gradient (CG) method, preconditioned by the  $V$ -cycle AMLI preconditioner, constructed as in (7). On each level the blocks  $A_{11}$  are solved exactly. The approximation  $S_A$  is obtained by assembling local Schur complement matrices on each level. We show the number of iterations for the full-length  $V$ -cycle preconditioner ('Coarsest=1'), as well as for the corresponding two-level method ('Coarsest=3,4,5,6,7').

From the numerical experiments we see that for strong anisotropy  $\alpha = 0$  is the best choice, while  $\alpha = 0.1$  leads to a slight improvement when  $\varepsilon = 1$ . In the two-level case the differences are minor and introducing  $\alpha$  is hardly justified. The results indicate that the Schur complement approximation  $\tilde{S}_2$  is more robust than  $A_H$ .

**Problem 9.3** We consider a simple linear elasticity problem in two dimensions, where a homogeneous body occupies a domain  $\Omega = [0, 1]^2$ . The material parameters (Young modulus  $E$  and Poisson ratio  $\nu$ ) are chosen as  $E = 1$  and  $\nu = 0.2$ . The discretization is done either by right-angled triangles and linear basis functions or by square mesh and bilinear basis functions. The block  $A_{11}$  is solved exactly. The stopping criterion is set to decrease the relative residual norm below  $10^{-6}$ .

The results are presented in Tables 4 and 5. In all cases  $m = 2$  and full-length recurrence is used. We show results for the pure  $V$ -cycle and for a stabilized preconditioner. In the latter case, a stabilization is done on some of the intermediate levels by solving the corresponding system with the same preconditioned method to a lower accuracy  $10^{-3}$ , which requires 3-4 inner iterations. An analogous way of stabilizing an AMLI-type solver of additive form is used in [10].

The parameter  $\mu$  in Tables 4 and 5 indicates how often the stabilization has been performed. The values  $\mu = 2$  and  $\mu = 3$  mean that we have done it on each second

Problem size		81		289		1089		4225		16049	
Total no.levels		4		5		6		7		8	
$\alpha$	$\varepsilon_2$	Full length	Two level								
0	$10^{-3}$	3	3	3	3	5	4	7	5	11	7
	1	6	6	8	6	11	7	14	7	18	7
	$10^3$	3	3	5	4	7	6	10	7	15	8
0.1	$10^{-3}$	4	3	6	4	9	5	13	6	19	7
	1	6	6	7	7	9	6	10	7	11	6
	$10^3$	5	4	8	5	12	6	18	7	23	8
0.5	$10^{-3}$	5	3	7	4	12	5	16	6	22	7
	1	7	6	8	7	9	7	10	7	11	7
	$10^3$	6	4	10	5	16	6	19	7	29	8
0.9	$10^{-3}$	4	3	7	4	12	5	17	6	24	7
	1	7	6	9	7	10	7	10	7	11	7
	$10^3$	6	4	11	5	16	6	23	7	29	8
1	$10^{-3}$	5	3	7	4	12	5	17	6	23	7
	1	7	6	9	7	10	7	10	7	12	7
	$10^3$	6	4	11	5	16	6	22	7	29	8

Table 3: Problem 9.2: Iteration counts for the full-length V-cycle and for the two-level preconditioners

Size	Total no levels	V-cycle		Stabilized GCG		Unprec. CG
		PCG	GCG	$\mu = 2$	$\mu = 3$	
2x81	4	8	11	10	10	71
2x289	5	12	16	11	11	150
2x1089	6	16	22	11	11	291
2x4225	7	21	30	12	12	461
2x16641	8	30	41	12	12	-

Table 4: Problem 9.3 (triangles)

Size	Total no levels	V-cycle		Stabilized GCG		Unprec. CG
		PCG	GCG	$\mu = 2$	$\mu = 3$	
2x81	4	7	8	6	6	32
2x289	5	10	11	8	8	60
2x1089	6	14	16	11	11	117
2x4225	7	19	22	12	12	-
2x16641	8	25	30	12	12	-

Table 5: Problem 9.3 (quadrilaterals)

level, respectively, on each third level. Stabilization on the second highest level is always required. Numerical tests, not included here, indicate that if we release the latter for  $\mu > 2$ , the method behaves as in the unstabilized V-cycle preconditioned case. An observation to mention here is that the method behaves better on quadrilateral meshes than on triangular meshes. This is even more pronounced in the case of a scalar equation.

In both Tables 4 and 5, under the V-cycle, there are two columns containing iteration counts for the PCG and for the GCG methods, correspondingly. Since in these tests the V-cycle preconditioner is fixed (due to the exact solve with the top-left blocks), we can use the standard PCG method. The experiments indicate that it is slightly better than the GCG method, which is expected for symmetric problems. In the stabilized case the preconditioner is varying, and therefore only the GCG method is used.

## 10 Concluding remarks

It has been shown how one can construct and estimate condition numbers for a general approach to approximate the arising global Schur complement matrices by use of local, macroelementwise computed Schur complement matrices.

Thereby the influence of the number of subdivisions of each macroelement to achieve an optimal order method has been shown.

The elementwise approach is applicable also for problems in three space dimensions. It is tested in the context of nonsymmetric saddle point problems and numerical experiments in 2D and 3D can be found in [14].

It has further been shown how one can use the strong properties of the HBF element matrices without actually having to implement them, thus avoiding their disadvantage of being less sparse.

Finally, the influence of approximations of the pivot block matrix  $A_{11}$  have been discussed in some detail.

## References

- [1] O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, New York, 1994.

- [2] O. Axelsson, A survey of algebraic multilevel iteration (AMLI) methods. *BIT Numerical Mathematics*, 43 (2003), 863-879.
- [3] O. Axelsson, On iterative solvers in structural mechanics; separate displacement orderings and mixed variable methods. *Mathematics and Computers in Simulation* 50 (1999), 11-30.
- [4] O. Axelsson and A. Barker, *Finite Element Solution of Boundary Value Problems*, Classics in Appl. Math., SIAM, Philadelphia, 2001.
- [5] O. Axelsson and R. Blaheta, Two simple derivations of universal bounds for the CBS inequality constant, *Applications of Mathematics* 49 (2004), 57-72.
- [6] O. Axelsson and I. Gustafsson, Preconditioned and two-level multigrid methods of arbitrary degree of approximations, *Mathematics of Computation* 40 (1983), 219-242.
- [7] O. Axelsson, Y. R. Hakopian and Y. A. Kuznetsov, Multilevel preconditioning for perturbed finite element matrices, *IMA Journal of Numerical Analysis*, 17 (1997), 125-149.
- [8] O. Axelsson and S. Margenov, On multilevel preconditioners which are optimal with respect to both problem and discretization parameters. *Computational Methods in Applied Mathematics* 3 (2003), 6-22.
- [9] O. Axelsson and M. Neytcheva, Algebraic multilevel iteration methods for Stieltjes matrices, *Numerical Linear Algebra with Applications*, 1 (1994), 213-236.
- [10] O. Axelsson and A. Padiy, On the additive version of the algebraic multilevel iteration method for anisotropic elliptic problems. *SIAM Journal on Scientific Computing* 20, 1807-1830 (1999).
- [11] O. Axelsson and P. S. Vassilevski, Algebraic multilevel preconditioning methods II, *SIAM Journal on Numerical Analysis* 27 (1990), 1569-1590.
- [12] O. Axelsson and P. S. Vassilevski, The AMLI method: an algebraic multilevel method for positive definite sparse matrices. Catholic University of Nijmegen, Report no. 0036, 1999.
- [13] O. Axelsson, P.S. Vassilevski. A survey of a class of algebraic multilevel iteration methods for positive definite symmetric matrices. Neittaanmäki, Pekka (ed.) et al., ENUMATH 99. Numerical mathematics and advanced applications. *Proceedings of the 3rd European conference*, Jyväskylä, Finland, July 26-30, 1999. Singapore: World Scientific. 16-30 (2000).
- [14] E. Bängtsson, B. Lund, A comparison between to solution techniques to solve the equations of linear isostasy. Submitted for publication.

- [15] E. Bängtsson, M. Neytcheva. An agglomerate multilevel preconditioner for linear isostacy saddle point problems. *Lecture Notes in Computer Science, Proceedings of the 5th International Conference on Large-scale Scientific Computations 2005* (eds: Lirkov, I. and Margenov, S. and Wasniewski, J.), 3743 (2006), pp. 113–120.
- [16] R. Blaheta, Displacement decomposition - incomplete factorization preconditioning techniques for linear elasticity problems. *Numerical Linear Algebra with Applications*, 1 (1994), 107-128.
- [17] J. Kraus, Algebraic multilevel preconditioning of finite element matrices using local Schur complements, *Numerical Linear Algebra with Applications*, 13 (2006), 49-70.
- [18] Y. Notay, Using approximate inverses in algebraic multilevel methods. *Numer. Math.* 80 (1998), 397-417.
- [19] Y. Notay, Optimal order preconditioning of finite difference matrices. *SIAM Journal on Scientific Computing* 21 (2000), 1991-2007.
- [20] I. Pultarová, The strengthened CBS inequality constant for second order elliptic partial differential operator and for hierarchical bilinear finite element functions, *Applications of Mathematics*, 50 (2005), 323-329.