

## ALGEBRAIC ANALYSIS OF TWO-LEVEL MULTIGRID METHODS FOR EDGE ELEMENTS\*

ARTEM NAPOV<sup>†</sup> AND RONAN PERRUSSEL<sup>‡</sup>

**Abstract.** We present an algebraic analysis of two-level multigrid methods for the solution of linear systems arising from the discretization of the curl-curl boundary value problem with edge elements. The analysis is restricted to the singular compatible linear systems as obtained by setting to zero the contribution of the lowest order (mass) term in the associated partial differential equation. We use the analysis to show that for some discrete curl-curl problems, the convergence rate of some Reitzinger-Schöberl two-level multigrid variants is bounded independently of the mesh size and the problem peculiarities. This covers some discretizations on Cartesian grids, including problems with isotropic coefficients, anisotropic coefficients and/or stretched grids, and jumps in the coefficients, but also the discretizations on uniform unstructured simplex grids.

**Key words.** convergence analysis, multigrid, algebraic multigrid, two-level multigrid, Reitzinger-Schöberl multigrid, preconditioning, aggregation, edge elements

**AMS subject classifications.** 65N55, 65N12, 65N22, 35Q60

**1. Introduction.** We present an algebraic analysis of two-level multigrid methods for the solution of the  $n \times n$  symmetric positive semi-definite linear systems

$$(1.1) \quad A \mathbf{u} = \mathbf{b}$$

arising from the discretization of the boundary value problem

$$(1.2) \quad \begin{cases} \operatorname{curl}(\tilde{\mu}^{-1} \operatorname{curl} \mathbf{E}) + \beta \mathbf{E} = \mathbf{f} & \text{in } \Omega, \\ \mathbf{E} \times \mathbf{n} = \mathbf{g}_D & \text{on } \Gamma_D, \\ (\tilde{\mu}^{-1} \operatorname{curl} \mathbf{E}) \times \mathbf{n} = \mathbf{g}_N & \text{on } \Gamma_N = \partial\Omega \setminus \Gamma_D, \end{cases}$$

with edge elements. In this boundary value problem, the domain  $\Omega$  is a polygonal bounded simply connected region of  $\mathbb{R}^2$  or  $\mathbb{R}^3$ , the coefficient  $\tilde{\mu}$  is a scalar positive function in two dimensions (2D) and a  $3 \times 3$  diagonal matrix  $\tilde{\mu} = \operatorname{diag}(\mu_x, \mu_y, \mu_z)$  in three dimensions (3D) with the diagonal entries  $\mu_x > 0$ ,  $\mu_y > 0$ , and  $\mu_z > 0$  being piecewise constant functions on  $\Omega$ , and  $\beta \geq 0$  is a constant function on  $\Omega$ . For isotropic problems in three dimensions we denote with  $\mu := \mu_x = \mu_y = \mu_z$  the diagonal entries of  $\tilde{\mu}$ , whereas in two dimensions we set  $\mu = \tilde{\mu}$ . Further,  $\mathbf{E}$  is an unknown vector function on  $\Omega$ ,  $\mathbf{f}$  is a given vector function on  $\Omega$ ,  $\mathbf{g}_D$  and  $\mathbf{g}_N$  are given function on  $\Gamma_D \subset \partial\Omega$  and  $\Gamma_N \subset \partial\Omega$ , respectively, and  $\mathbf{n}$  is the unit outward normal vector to the boundary surface  $\partial\Omega = \Gamma_D \cup \Gamma_N$ . By edge elements we mean the lowest-order elements of the first family proposed by Nédélec on simplex and Cartesian grids [21]. For these elements the degrees of freedom are associated with the individual edges of the grid. We therefore consider a discretization with Cartesian or simplex grids and further assume that each boundary edge (2D) and each boundary face (3D) of the discretization grid belongs either to  $\Gamma_D$  or to  $\Gamma_N$  and is therefore not “shared” between these subsets of  $\partial\Omega$ .

We develop our analysis for a particular case of the above problem: a *singular compatible* linear system (1.1) arising from the discretization of the boundary value problem (1.2) with

\*Received June 16, 2018. Accepted August 27, 2019. Published online on November 13, 2019. Recommended by Stefan Vandewalle.

<sup>†</sup>Service de Métrologie Nucléaire, Université Libre de Bruxelles, (C.P. 165-84), 50 Av. F.D. Roosevelt, B-1050 Brussels, Belgium (anapov@ulb.ac.be).

<sup>‡</sup>LAPLACE, Université de Toulouse, CNRS, INPT, UPS, Toulouse, France (perrussel@laplace.univ-tlse.fr).

$\beta = 0$ . Setting  $\beta = 0$  implies that the boundary value problem (1.2) is singular, its null space being spanned by the gradients of smooth enough functions that satisfy the boundary conditions. Likewise, the corresponding linear system (1.1) is also singular with the null space of the system matrix  $A$  given by the columns of the discrete gradient operator. Moreover, we assume that the system (1.1) is compatible, that is, that  $\mathbf{b} \in \mathcal{R}(A)$  so that it has solutions and the convergence to a solution can be quantified.

Our focus on the singular compatible linear systems for the case  $\beta = 0$  is motivated by two reasons. First, this particular case is important in its own right. For instance, the compatible linear systems corresponding to  $\beta = 0$  arise in situations where the magnetostatic approximation is considered. Although the solution techniques then typically tend to eliminate the singularity [1, 16, 27], the effective condition number of the singular system is typically smaller than the one of the corresponding regularized variants (see, e.g., [14]), and keeping the singularity is therefore more attractive for iterative solution methods. Regarding compatibility, we note that, although the compatibility of the continuous problem does not imply that of the discrete counterpart, in the considered applications this latter is typically enforced during or after the discretization [28].

The analysis of two-level multigrid methods for singular compatible systems corresponding to  $\beta = 0$  is also helpful in understanding the convergence of the two-level multigrid methods for regular systems associated with  $\beta > 0$ . This is because, on one hand, the multigrid methods for singular systems with  $\beta = 0$  typically behave similarly to those for regular systems with  $\beta \rightarrow 0$  and, on the other hand, the multigrid convergence when  $\beta \rightarrow 0$  is representative for the multigrid methods with any  $\beta > 0$ . The first point holds if in the case where  $\beta \rightarrow 0$  the smoothing of the (near)kernel component of the correction is efficient enough compared to the overall efficiency of the multigrid method. This of course depends on the multigrid ingredients, and for the ingredients considered here the efficiency of the (near)kernel smoothing is indeed observed in practice. The second point holds because the discrete counterpart of the lowest order (mass) term  $\beta \mathbf{E}$  is typically well conditioned, and therefore its impact on the multigrid convergence is mostly benign (this phenomenon is quantified in, e.g., [3]). Moreover, the contribution of the discretized lowest order term to the system matrix is proportional to the square of the mesh size, and therefore its impact on the multigrid convergence decreases as the grid is refined. The similarity between the two-level convergence properties in the cases  $\beta = 0$  and  $\beta > 0$  is further illustrated below with the numerical experiments.

The second reason for considering the singular compatible systems with  $\beta = 0$  is the simplification of the underlying multigrid algorithms, which actually makes an algebraic analysis possible. More specifically, as discussed in Section 2.1 below, the multigrid smoothers typically used for the discrete boundary value problem (1.2) can then be replaced by, or even reduce to, a simpler variant.

The presented analysis is primarily intended for algebraic multigrid (AMG) methods. Such methods require little input from the user, the specificity of the system (1.1) being captured at every multigrid level by a prolongation matrix. Here we focus on the prolongations introduced by Reitzinger and Schöberl [27]; these prolongations amount to group the nodes of the discretization grid into problem-dependent aggregates and are therefore both simple and flexible. The analysis is, however, not limited to the Reitzinger and Schöberl prolongations; see the extended report [19] for details. Besides, the AMG framework imposes some other algorithmic peculiarities, including the use of multigrid as a preconditioner and the choice of a Galerkin coarse grid correction.

Note that the use of two-level analyses is typical for the AMG framework. Such analyses help with the automatic construction of a prolongation matrix at every level, representing an important element of the AMG design. A *two-level* character of the analyses comes with the fact that AMG methods typically build a given prolongation matrix only once the matrices on the previous levels are available, but also because the extra flexibility required for the analyses makes a multilevel extension challenging. The use of a two-level analysis in the design of an AMG method of course does not imply that this method has only two (or few) levels. It rather hints that it is possible to approach the convergence of a two-level method in a multilevel setting by carefully choosing the associated multigrid recursion, so that the analysis remains valuable in a multilevel setting. The above observations are well illustrated by the classical algebraic multigrid methods [18, 23, 29] for Poisson-like problems, which are designed based on the associated two-level analyses from [9, 17, 29].

The focus on the Reitzinger-Schöberl multigrid preconditioner is not solely motivated by its simplicity. The low memory requirements and a moderate cost per iteration are the other attractive features of this preconditioner. On the negative side, the convergence of the original method deteriorates with increasing number of levels [27]. This led to further improvements of the approach [5, 6, 13, 26], which, however, did not give full satisfaction. The Reitzinger-Schöberl multigrid has lost in popularity since the development of alternative approaches [4, 12, 15] based on an auxiliary space preconditioning, i.e., based on the application of the classical AMG methods for Poisson-like problems to an extended and transformed system. These latter methods are now considered as a reference. Although they typically exhibit a level-independent convergence [12], their cost per iteration and memory requirements are often higher than those of the Reitzinger-Schöberl preconditioner. Therefore, a proper redesign of the Reitzinger-Schöberl multigrid method, as recently proposed in [20], can lead to a competitive approach, which can further benefit from the present two-level analysis.

We now overview the main outcomes of the analysis. We begin with the discretizations of the boundary value problem (1.2) on Cartesian grids. In this setting, we first consider the problem with constant isotropic coefficients discretized on a grid with a square/cubic mesh, and show that a typical two-level Reitzinger-Schöberl multigrid method has a convergence rate bounded independently of the mesh size. We then consider problems in which the jumps in the coefficient  $\mu^{-1}$  are aligned with the grid lines and show that the convergence of the same Reitzinger-Schöberl two-level multigrid in two dimensions is also bounded independently of the jumps amplitude. On the other hand, if the coefficients are anisotropic, that is, if  $\mu_x^{-1}$ ,  $\mu_y^{-1}$ , and  $\mu_z^{-1}$  have a different magnitude, we show that the convergence is bounded independently of the mesh size and of the coefficient values as long as the aggregates in the Reitzinger-Schöberl method are aligned in the direction of the weakest coefficient; similar results hold for a grid-based anisotropy as induced by stretched grids. We then continue with unstructured uniform grids in two and three dimensions, proving that for the Reitzinger-Schöberl method with aggregates of bounded size, the convergence is bounded independently of the mesh size. Most of these *two-level* results corroborate the observations made in [20] based on the numerical experiments with the *multilevel* Reitzinger-Schöberl method.

Let us also mention some related works. Several approaches are known in the analysis of multigrid methods for the (non-transformed) discrete boundary value problem (1.2). Amongst these, the analyses based on multi-level decompositions are presented in [2, 11], whereas those using the Fourier modes—the so-called (local) Fourier analysis—are introduced in [7]. The first family of approaches covers multi-level multigrid methods obtained by the progressive refinement of an initial discretization on a simplex grid, whereas the second family applies to two-level multigrid methods used on a structured Cartesian grid with a structured Cartesian coarse grid. In both cases the structure of the coarse and the fine grids are strongly related,

which is typical for multigrid methods of geometric type. On the other hand, the analysis presented here allows for more freedom in choosing the coarse grid structure, making it suitable for multigrid methods of algebraic type.

The remainder of this paper is structured as follows. In Section 2 we present the core of the analysis and show how this analysis can be applied to a two-level multigrid method of Reitzinger-Schöberl type. In Section 3 and Section 4 we highlight some of the outcomes of the analysis for structured and unstructured grids, respectively, and further illustrate them with numerical experiments that compare the cases  $\beta = 0$  and  $\beta > 0$ . Concluding remarks are given in Section 5.

**2. Analysis.** In this section we present our analysis of two-level multigrid methods for the considered singular problems. This is done by first introducing the two-level multigrid methods covered by our analysis, followed by recalling a general convergence estimate that holds when these methods are applied to the singular compatible linear systems, then presenting the main steps of the analysis, and eventually showing how the analysis can be used in the case of the Reitzinger-Schöberl two-level multigrid.

**2.1. Two-level preconditioner.** The two-level multigrid preconditioner covered by the analysis is defined by

$$(2.1) \quad B_{\text{TG}} = M^{-1}(2M - A)M^{-1} + (I - M^{-1}A) P A_c^g P^T (I - A M^{-1}),$$

where  $M$  is an  $n \times n$  symmetric matrix called smoother,  $P$  is an  $n \times n_c$  prolongation matrix with  $n_c < n$ ,  $A_c^g$  is an  $n_c \times n_c$  matrix representing a proper generalized inverse of the coarse grid matrix, and  $I$  is the  $n \times n$  identity matrix. The structure of the preconditioner is best understood from the corresponding iteration matrix

$$I - B_{\text{TG}}A = (I - M^{-1}A) (I - P A_c^g P^T A) (I - M^{-1}A),$$

which is a product of three simpler iteration matrices corresponding to the pre-smoothing iteration, the coarse-grid correction, and the post-smoothing iteration.

In typical multigrid applications, the smoothing iteration is a simple one-level method, such as a weighted Jacobi or Gauss-Seidel iteration. This is, however, not the case for the smoothers [2, 11] used in multigrid methods for the boundary value problem (1.2), as such smoothers are required to effectively reduce some error components in the space of the discrete gradients. In particular, a Hiptmair smoother [11] achieves this by combining together an iteration of a one-level method (Jacobi or Gauss-Seidel) for the original system (1.1) and an iteration of a one-level method for an auxiliary system, this latter system being obtained from the system (1.1) by projection on a subspace of discrete gradients. The second iteration has, however, no effect when considering singular compatible systems with  $\beta = 0$ , and in this case the Hiptmair smoother is equivalent to a one-level method for the original system (1.1).

The present analysis relies on a simple weighted Jacobi smoother  $M = \omega_J^{-1}D$ , with  $D = \text{diag}(A)$ . The weighting parameter  $\omega_J$  is typically chosen<sup>1</sup> as  $\omega_J^{-1} \approx \lambda_{\max}(D^{-1}A)$ ; the typical value for  $\omega_J^{-1}$  is then around 2–6, including for the problems with anisotropy or jumps in the coefficients as considered below. Of course, in practice a Gauss-Seidel iteration is preferable as a smoother as it typically gives better results and does not require any weighting parameter. However, the resulting analysis then still provides a useful indication of the multigrid convergence.

<sup>1</sup>One can show that this choice yields the best convergence bounds below; however, it does not always yields the best actual convergence rate for a two-level method based on a weighted Jacobi smoother.

Regarding the coarse grid correction, we assume that the coarse grid matrix is given by the Galerkin formula  $A_c = P^T A P$ . In the case of a singular system matrix  $A$ , the coarse grid matrix  $A_c$  is generally also singular, prompting the use of a generalized inverse  $A_c^g$  of  $A_c$  in the coarse grid correction step. Another implication of the Galerkin formula is that the corresponding coarse grid correction step is a projection matrix which is entirely determined by the system matrix  $A$  and the prolongation  $P$ .

We note that in the case where  $\beta = 0$ , the resulting system is symmetric positive semi-definite and compatible. Therefore it is suitable for the preconditioned conjugate gradient method. Then, the associated preconditioner should be symmetric positive definite (SPD), and the two-level multigrid preconditioner in (2.1) satisfies this requirement provided that  $2M - A$  is also SPD—a condition fulfilled by common multigrid smoothers.

**2.2. Two-level estimate.** The following theorem is at the foundation of our analysis. It provides a practical convergence estimate for two-level multigrid methods when applied to singular compatible systems. More specifically, it shows that under some rather general assumptions, the rate of convergence of the conjugate gradient method used with the two-level multigrid preconditioner (2.1) can be bounded above with the help of the parameter  $\kappa_\pi$  defined by (2.4) below. The results gathered in the theorem are borrowed from [25].

**THEOREM 2.1.** *Let  $A$  be an  $n \times n$  symmetric positive semi-definite matrix,  $P$  be an  $n \times n_c$  full rank matrix for some  $n_c < n$ , and  $A_c = P^T A P$ . Let  $B_{\text{TG}}$  be defined by (2.1) with a symmetric  $n \times n$  matrix  $M$  such that  $\omega M - A$  is an SPD matrix for some  $\omega \in (0, 2)$  and with a matrix  $A_c^g$  satisfying  $A_c A_c^g A_c = A_c$ , that is, being a proper generalized inverse of  $A_c$ . Let  $\pi$  be an  $n \times n$  projector whose null space satisfies*

$$(2.2) \quad \mathcal{N}(\pi) = \text{span}(\mathcal{N}(A), \mathcal{R}(P)),$$

where  $\mathcal{R}(\cdot)$  denotes the range of a matrix and  $\mathcal{N}(\cdot)$  is its null space.

Then the approximation  $\mathbf{u}_k$  of a solution  $\mathbf{u}$  of the compatible system (1.1) produced at the iteration  $k$  of the conjugate gradient method preconditioned with  $B_{\text{TG}}$  satisfies

$$(2.3) \quad \|\mathbf{u} - \mathbf{u}_k\|_A \leq 2 \left( \frac{\sqrt{\kappa_\pi} - 1}{\sqrt{\kappa_\pi} + 1} \right)^k \|\mathbf{u} - \mathbf{u}_0\|_A,$$

where

$$(2.4) \quad \kappa_\pi = \frac{1}{2 - \omega} \sup_{\mathbf{v} \notin \mathcal{N}(A)} \frac{\mathbf{v}^T \pi^T M \pi \mathbf{v}}{\mathbf{v}^T A \mathbf{v}}.$$

Moreover,

$$(2.5) \quad \dim(\mathcal{N}(\pi)) = \dim(\mathcal{N}(A)) + \text{rank}(A_c).$$

*Proof.* The inequality (2.3) with  $\kappa_\pi$  defined as in (2.4) follows directly from the combination of Lemma 2.4, Theorem 3.4, and Theorems 3.5 in [25], together with the fact that for  $X = M(2M - A)^{-1}M$  there holds

$$\frac{\mathbf{v}^T X \mathbf{v}}{\mathbf{v}^T M \mathbf{v}} = \frac{\mathbf{v}^T M (2M - A)^{-1} M \mathbf{v}}{\mathbf{v}^T M \mathbf{v}} \leq \frac{\mathbf{v}^T M (2M - \omega M)^{-1} M \mathbf{v}}{\mathbf{v}^T M \mathbf{v}} = \frac{1}{2 - \omega}.$$

The equality (2.5) is equivalent to equality (3.6) in [25], which is satisfied in the considered setting.  $\square$

Results similar to the above are typically used as a first step in the analyses of algebraic multigrid methods for *regular* systems [9, 10, 17, 22]; see also [24] for a review. However,

the above theorem is for *singular compatible* systems and differs from such results in that, in addition to the range  $\mathcal{R}(P)$  of the prolongation, the null space of the projector  $\pi$  should also contain the null space  $\mathcal{N}(A)$  of the system matrix  $A$  as stated in condition (2.2). Condition (2.2) actually ensures that the parameter  $\kappa_\pi$  is not trivially infinite as the denominator in (2.4) can only become zero when the numerator is also zero<sup>2</sup>. Now, the need to account for the null space  $\mathcal{N}(A)$  in condition (2.2) is of little importance in the case (common for multigrid applications) where the dimension of  $\mathcal{N}(A)$  is small, as the prolongation is then typically designed to contain this null space in its range. However, for the problems considered here the dimension of the null space  $\mathcal{N}(A)$  as spanned by the discrete gradients is quite large, and the range of the prolongation generally contains only part of this null space. In such a case, the presence of  $\mathcal{N}(A)$  in condition (2.2) is important.

The combination of the observations on the weighted Jacobi smoothers  $M = \omega_J^{-1}D$  from Section 2.1 with the results from Theorem 2.1 implies that, if  $\omega_J^{-1} = \lambda_{\max}(D^{-1}A)$ , then  $\omega = 1$  and

$$(2.6) \quad \kappa_\pi \leq \omega_J^{-1} \tilde{\kappa}_\pi,$$

where

$$(2.7) \quad \tilde{\kappa}_\pi = \sup_{\mathbf{v} \notin \mathcal{N}(A)} \frac{\mathbf{v}^T \pi^T D \pi \mathbf{v}}{\mathbf{v}^T A \mathbf{v}}$$

and the projector  $\pi$  is as in Theorem 2.1. As a result, the convergence rate of a two-level multigrid method can be kept under control if the parameter  $\tilde{\kappa}_\pi$  is bounded above; we therefore consider this parameter in the remainder of this paper.

Let us briefly comment on the sharpness of the inequality (2.3). It corresponds to the classical convergence bound for the conjugate gradient method for symmetric positive semi-definite compatible systems if  $\kappa_\pi$  is replaced by the effective condition number  $\kappa_{\text{eff}}$  defined via (3.1) below. Such a bound is not sharp [30] if, as is the case of the considered applications, the eigenvalues of the preconditioned matrix are clustered.

**2.3. Main steps.** We now highlight the main steps of our analysis. The key idea behind these steps is to bound both the numerator and the denominator in the definition (2.7) of the parameter  $\tilde{\kappa}_\pi$  by a sum of nonnegative terms associated with the oriented faces of the grid. Such a decomposition relies on the structure of the range  $\mathcal{R}(A)$  of the system matrix for the considered problems and is therefore a natural option for the denominator of (2.7). Regarding the numerator, a similar decomposition can be obtained by properly choosing the projector matrix  $\pi$ .

To better highlight the structure of  $\mathcal{R}(A)$ , we briefly review the assembly procedure of the system matrix  $A$ . We assume that the system matrix is assembled using the standard finite element methodology, that is, there holds

$$(2.8) \quad A = \sum_{e=0}^{n^{(e)}} T_e A_e T_e^T,$$

where  $A_e$  is the element matrix of the  $e$ th element,  $T_e$  is the local-to-global index mapping corresponding to this element, and  $n^{(e)}$  is the number of elements. Unless stated otherwise,

---

<sup>2</sup>The condition  $\mathbf{v} \notin \mathcal{N}(A)$  under the supremum in (2.4) does prevent the actual division by zero but does not exclude the vectors  $\mathbf{v}$  which are arbitrarily close to  $\mathcal{N}(A)$ . Therefore, if there is a vector  $\mathbf{v} \in \mathcal{N}(A)$  for which the numerator in (2.4) is nonzero, the supremum in (2.4) is infinite.

the edges on the boundary  $\Gamma_D$  are eliminated from the system matrix since the corresponding unknowns are then known; the corresponding indices are therefore not mapped by  $T_e$ , where  $e = 1, \dots, n^{(e)}$ .

For the considered boundary value problem (1.2) with  $\beta = 0$ , the element matrix for a non-degenerate element  $e$  can be written [8] as

$$(2.9) \quad A_e = C_e M_e C_e^T,$$

where  $C_e = (\mathbf{c}_{f_e}^{(e)})$  is the transpose of the local discrete curl matrix for the considered element and  $M_e$  is an SPD matrix. In particular, each column  $\mathbf{c}_{f_e}^{(e)}$  of the transpose of the local discrete curl matrix is given by

$$(2.10) \quad (\mathbf{c}_{f_e}^{(e)})_i = \begin{cases} 1 & \text{if the local edge } i \text{ belongs to the boundary of the face } f_e \\ & \text{and follows its orientation,} \\ -1 & \text{if the local edge } i \text{ belongs to the boundary of the face } f_e \\ & \text{but has opposite orientation,} \\ 0 & \text{otherwise,} \end{cases}$$

and is associated with a particular oriented face  $f_e$  of the element, being nonzero only for the edges belonging to this face. Note that the orientation of a face and of its boundary are supposed compatible in the sense of Stokes' theorem.

From (2.9) we deduce the contribution of the element  $e$  to the decomposition for the denominator of (2.7). Indeed, since  $M_e$  is SPD, there exist real numbers  $\gamma_{f_e}^{(e)} > 0$ ,  $f_e = 1 \dots, n_e^{(f)}$ , where  $n_e^{(f)}$  is the number of faces in the element  $e$  such that the matrix

$$M_e - \text{diag}(\gamma_{f_e}^{(e)})$$

is positive semi-definite. Then, the aforementioned contribution corresponds to

$$(2.11) \quad \mathbf{v}_e^T A_e \mathbf{v}_e \geq \sum_{f_e=1}^{n_e^{(f)}} \gamma_{f_e}^{(e)} (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2.$$

In particular, for the two-dimensional problems each element “coincides” with its only oriented face, with hence  $n_e^{(f)} = 1$ , whereas the value of  $\gamma_{f_e}^{(e)}$  is then the inverse of the area of the face multiplied by  $\mu$ .

The combination of (2.8) and (2.11) yields the decomposition for the denominator of (2.7):

$$(2.12) \quad \mathbf{v}^T A \mathbf{v} \geq \sum_f \gamma_f (\mathbf{v}^T \mathbf{c}_f)^2,$$

where  $f$  is the global face index, which may correspond to the element indices  $f_e$  of several elements sharing this face,  $\gamma_f = \sum_e \gamma_{f_e}^{(e)} > 0$  is the sum of the contributions of the elements sharing the face  $f$  not included into  $\Gamma_D$ , and  $\mathbf{c}_f = \pm T_e \mathbf{c}_{f_e}^{(e)}$  is the global curl vector associated with the oriented face  $f$ ,  $f = 1, \dots, n^{(f)}$ , with the  $\pm 1$  factor accounting for the match between the orientations of the local face  $f_e$  and the corresponding global face  $f$ . In particular, this latter can be rewritten as

$$(\mathbf{c}_f)_i = \begin{cases} 1 & \text{if edge } i \text{ belongs to the boundary of face } f \text{ and follows its orientation,} \\ -1 & \text{if edge } i \text{ belongs to the boundary of face } f \text{ but has opposite orientation,} \\ 0 & \text{otherwise,} \end{cases}$$

where, due to the chosen  $T_e$ , only the edges and the faces not included into  $\Gamma_D$  are considered. Note in particular that the vectors  $\mathbf{c}_f$ ,  $f = 1 \dots, n^{(f)}$ , span the range  $\mathcal{R}(A)$  of the system matrix.

The projector  $\pi$  is chosen so that the numerator of (2.7) is bounded above by a decomposition similar to the one in (2.12). This is achieved if for any vector  $\mathbf{v}$  the entry  $(\pi\mathbf{v})_i$  is zero for every index  $i$  in some subset  $\mathcal{E}_\pi^0$ , whereas for every index  $i$  outside this subset it is given by

$$(2.13) \quad (\pi\mathbf{v})_i = \mathbf{v}^T \left( \sum_f \alpha_{i,f} \mathbf{c}_f \right),$$

with typically only a handful, say  $m_i$ , of the  $\alpha_{i,f}$  being nonzero for a given index  $i$ . Note that for  $\pi$  to be a projector it is enough to require that the  $i$ th entry of the vector  $\sum_f \alpha_{i,f} \mathbf{c}_f$  is set to 1 whereas the other entries of this vector whose indices are outside  $\mathcal{E}_\pi^0$  are set to 0. These requirements also ensure that the range of the projector is  $\{\mathbf{v} \mid (\mathbf{v})_i = 0 \text{ for } i \in \mathcal{E}_\pi^0\}$  and therefore that

$$(2.14) \quad \dim(\mathcal{N}(\pi)) = |\mathcal{E}_\pi^0|,$$

where  $|\mathcal{E}_\pi^0|$  is the number of elements in  $\mathcal{E}_\pi^0$ , that is, the number of entries set to zero by the projector. If such a projector  $\pi$  exists<sup>3</sup> and satisfies the requirements of Theorem 2.1, then for  $D = \text{diag}(d_i)$  the numerator of (2.7) satisfies

$$(2.15) \quad \begin{aligned} \mathbf{v}^T \pi^T D \pi \mathbf{v} &= \sum_{i \notin \mathcal{E}_\pi^0} d_i \left( \mathbf{v}^T \left( \sum_f \alpha_{i,f} \mathbf{c}_f \right) \right)^2 \\ &\leq \sum_{i \notin \mathcal{E}_\pi^0} d_i m_i \sum_f \alpha_{i,f}^2 (\mathbf{v}^T \mathbf{c}_f)^2 = \sum_f \alpha_f^2 (\mathbf{v}^T \mathbf{c}_f)^2, \end{aligned}$$

where  $\alpha_f^2 = \sum_{i \notin \mathcal{E}_\pi^0} m_i d_i \alpha_{i,f}^2$  and where we have used the inequality

$$\left( \sum_{k=1}^{m_i} x_k \right)^2 \leq m_i \sum_{k=1}^{m_i} (x_k)^2$$

for some real  $x_k$ ,  $k = 1, \dots, m_i$ .

The last step amounts to combining the inequalities (2.15) and (2.12) together with (2.7) and to replacing the quotient of two decompositions by the maximum of the quotients of the terms corresponding to the same oriented face; that is, in all generality,

$$(2.16) \quad \tilde{\kappa}_\pi \leq \max_f \frac{\alpha_f^2}{\gamma_f}.$$

**2.4. Toy problem.** To make the above (and the following) discussion less abstract we illustrate it with the following *toy problem*. It corresponds to the boundary value problem (1.2) stated on a square domain  $\Omega = [0, 1]^2$  with  $\mu = 1$ ,  $\beta = 0$ , with  $\Gamma_D$  corresponding to the top ( $y = 1$ ) and right ( $x = 1$ ) edges of the boundary  $\partial\Omega$ , and with  $\Gamma_N$  corresponding to the bottom ( $y = 0$ ) and left ( $x = 0$ ) edges. The problem is further discretized on a Cartesian  $4 \times 4$  grid (with hence  $h = 1/3$ ) using the lowest order edge elements on a square mesh; the boundary condition on  $\Gamma_D$  is imposed by elimination. The corresponding domain, the associated grid, the orientation, and the ordering of the edge unknowns are represented in Figure 2.1 (left).

<sup>3</sup>The projector  $\pi$  may not exist if for a given  $i \notin \mathcal{E}_\pi^0$  there is no vector of the form  $\sum_f \alpha_{i,f} \mathbf{c}_f$  whose  $i$ th entry is 1 and the other entries with indices outside  $\mathcal{E}_\pi^0$  are 0.



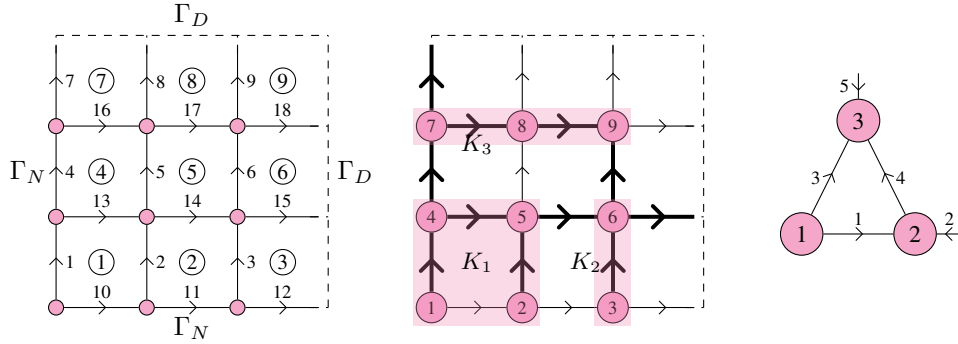


FIG. 2.1. The domain  $\Omega$ , the associated grid, the orientation and the ordering of the edges (and the corresponding edge unknowns), and the ordering of the faces (circled numbers) for the toy problem (left), the ordering of the nodes and the nodal aggregates (center), and the corresponding coarse grid (right). The bold edges on the central figure are edges whose indices form the set  $\mathcal{E}_\pi^0$ .

Regarding the resulting linear system, the system matrix  $A$  is assembled as in (2.8) with the element matrix for the 9 rectangular elements given by

$$A_e = \gamma \mathbf{c} \mathbf{c}^T, \quad \mathbf{c} = [-1 \quad 1 \quad 1 \quad -1]^T,$$

where  $\gamma = 1/h^2$  is the inverse of the area of the element and where the local unknowns are oriented as in Figure 2.1 (left) and ordered with vertical edges first: left one, then right one, and horizontal edges last: bottom one, then top one. Since the unknowns on the  $\Gamma_D$  boundary are eliminated, the resulting system matrix is a  $18 \times 18$  matrix. This matrix satisfies for  $\mathbf{v} = (v_i)$  the following variant of the decomposition<sup>4</sup> (2.12):

$$\begin{aligned}
 (2.17) \quad \mathbf{v}^T A \mathbf{v} &= \gamma \sum_f (\mathbf{v}^T \mathbf{c}_f)^2 \\
 &= \gamma \left( (-v_1 + v_2 + v_{10} - v_{13})^2 + (-v_2 + v_3 + v_{11} - v_{14})^2 + (-v_3 + v_{12} - v_{15})^2 \right. \\
 &\quad \left. + (-v_4 + v_5 + v_{13} - v_{16})^2 + (-v_5 + v_6 + v_{14} - v_{17})^2 + (-v_6 + v_{15} - v_{18})^2 \right. \\
 &\quad \left. + (-v_7 + v_8 + v_{16})^2 + (-v_8 + v_9 + v_{17})^2 + (-v_9 + v_{18})^2 \right),
 \end{aligned}$$

in which every term corresponds to a face, and the terms are ordered in the lexicographical order of the faces; this order is further given in Figure 2.1 (left) with circled numbers. The second equality above also specifies (up to a sign) the vectors  $\mathbf{c}_f$ ,  $f = 1, \dots, 9$ , via (the squares of) their scalar product with  $\mathbf{v}$ ; we fix the sign by assuming that the faces are oriented towards the reader.

Regarding the decomposition (2.15), we postpone until the next section the discussion on how to choose the projector  $\pi$  for a given Reitzinger-Schöberl prolongation matrix. We note, however, that for our toy problem a valid choice that fits the description of the previous section is given, via the product with a vector  $\mathbf{v}$ , by

$$\begin{aligned}
 \pi \mathbf{v} &= [0, 0, 0, 0, \mathbf{v}^T \mathbf{c}_4, 0, 0, \mathbf{v}^T \mathbf{c}_7, \mathbf{v}^T (\mathbf{c}_7 + \mathbf{c}_8), \\
 &\quad \mathbf{v}^T \mathbf{c}_1, \mathbf{v}^T \mathbf{c}_2, \mathbf{v}^T \mathbf{c}_3, 0, 0, 0, 0, 0, \mathbf{v}^T (\mathbf{c}_7 + \mathbf{c}_8 + \mathbf{c}_9)]^T.
 \end{aligned}$$

<sup>4</sup>The use of the equality (and not an inequality) at the first line of (2.17) is due to the fact that each element has only one face.

With  $D = \text{diag}(d_i)$  being the diagonal of the matrix  $A$ , the corresponding decomposition (2.15) is given by (with the terms after the inequality sign ordered as in (2.17))

$$\begin{aligned} \mathbf{v}^T \pi^T D \pi \mathbf{v} &= d_5(\mathbf{v}^T \mathbf{c}_4)^2 + d_8(\mathbf{v}^T \mathbf{c}_7)^2 + d_9(\mathbf{v}^T (\mathbf{c}_7 + \mathbf{c}_8))^2 \\ &\quad + d_{10}(\mathbf{v}^T \mathbf{c}_1)^2 + d_{11}(\mathbf{v}^T \mathbf{c}_2)^2 + d_{12}(\mathbf{v}^T \mathbf{c}_3)^2 + d_{18}(\mathbf{v}^T (\mathbf{c}_7 + \mathbf{c}_8 + \mathbf{c}_9))^2 \\ &\leq d_{10}(\mathbf{v}^T \mathbf{c}_1)^2 + d_{11}(\mathbf{v}^T \mathbf{c}_2)^2 + d_{12}(\mathbf{v}^T \mathbf{c}_3)^2 + d_5(\mathbf{v}^T \mathbf{c}_4)^2 \\ &\quad + (d_8 + 2d_9 + 3d_{18})(\mathbf{v}^T \mathbf{c}_7)^2 + (2d_9 + 3d_{18})(\mathbf{v}^T \mathbf{c}_8)^2 + 3d_{18}(\mathbf{v}^T \mathbf{c}_9)^2, \end{aligned}$$

where  $d_i = \gamma$  for the edges on  $\Gamma_N$ , that is, for  $i = 1, 4, 7, 10, 11, 12$ , and  $d_i = 2\gamma$  for the other edges. Hence, combining this latter inequality with (2.17) and bounding the result by a maximum over the quotients corresponding to individual faces one gets

$$\tilde{\kappa}_\pi = \sup_{\mathbf{v} \notin \mathcal{N}(A)} \frac{\mathbf{v}^T \pi^T D \pi \mathbf{v}}{\mathbf{v}^T A \mathbf{v}} \leq \frac{d_8 + 2d_9 + 3d_{18}}{\gamma} = 12.$$

**2.5. Reitzinger-Schöberl projector.** The actual construction of the projector  $\pi$  depends on the null space  $\mathcal{N}(A)$  of the system matrix and on the prolongation matrix  $P$  of the considered two-level multigrid method. The above analysis is applied here with a prolongation matrix of Reitzinger-Schöberl type. Therefore, we first introduce a suitable basis for the null space  $\mathcal{N}(A)$ , then present the Reitzinger-Schöberl prolongation scheme, then detail the projector construction, and eventually explain why such a projector satisfies the requirements of Theorem 2.1.

Regarding the null space  $\mathcal{N}(A)$  of the system matrix  $A$ , it is spanned by the discrete gradient vectors  $\mathbf{g}_i, i = 1, \dots, n^{(n)}$ , each of which is associated to a node of the discretization grid and given by

$$(\mathbf{g}_i)_j = \begin{cases} 1 & \text{if the edge } j \text{ starts at node } i, \\ -1 & \text{if the edge } j \text{ ends at node } i, \\ 0 & \text{otherwise.} \end{cases}$$

Note that, since the matrix  $A$  is symmetric, its range  $\mathcal{R}(A)$  is orthogonal to its null space  $\mathcal{N}(A)$ , and therefore  $\mathbf{c}_f^T \mathbf{g}_i = 0$  for all  $f = 1, \dots, n^{(f)}$  and  $i = 1, \dots, n^{(n)}$ .

A Reitzinger-Schöberl prolongation matrix is determined in two steps. First, the auxiliary nodal aggregates are chosen by grouping all the grid nodes outside  $\Gamma_D$  into disjoint sets  $K_i, i = 1, \dots, n_c^{(n)}$ , called here auxiliary nodal aggregates; each of these aggregates represents a node on a coarse grid. A possible set of auxiliary nodal aggregates for the toy problem is represented with shaded rectangles in Figure 2.1 (center). Here we further assume that the part of the grid restricted to each auxiliary nodal aggregate is connected<sup>5</sup>; this assumption is satisfied by all known nodal aggregation schemes. The set  $K_0$  gathers the grid nodes belonging to  $\Gamma_D$ .

Second, the actual (edge) aggregates are formed from the auxiliary nodal aggregates. More precisely, each edge aggregate  $M_i, i = 1, \dots, n_c$ , is associated to a couple of connected auxiliary nodal aggregates, say  $K_{i_1}$  and  $K_{i_2}$ . On a coarse grid this means that the coarse edge  $i$  has  $i_1$  as a starting coarse node and  $i_2$  as an ending coarse node. Further, the aggregate  $M_i$  gathers the edges that connect  $K_{i_1}$  and  $K_{i_2}$ . That is, denoting by  $j = (j_1, j_2)$  the fine edge  $j$

<sup>5</sup>Said otherwise, every two nodes belonging to an auxiliary nodal aggregate are connected with a path of edges which also belongs to this aggregate.

whose starting point is  $j_1$  and whose ending point is  $j_2$ , we have

$$\begin{aligned}
 M_i^+ &= \{j = (j_1, j_2) \mid j_1 \in K_{i_1}, j_2 \in K_{i_2}\}, \\
 M_i^- &= \{j = (j_1, j_2) \mid j_1 \in K_{i_2}, j_2 \in K_{i_1}\}, \\
 M_i &= M_i^+ \cup M_i^-.
 \end{aligned}$$

The edge aggregates for the edges not included in  $\Gamma_D$  but such that one of the end nodes is in  $\Gamma_D$  (pending edges) are also defined as above but with either  $i_1 = 0$  or  $i_2 = 0$ ; such aggregates correspond to pending coarse edges. As an example, for our toy problem the edge aggregate that connects  $K_1$  to  $K_2$  is given by  $M_1 = M_1^+ = \{11, 14\}$ , whereas the one that connects  $K_0$  to  $K_2$  corresponds to  $M_2 = M_2^- = \{12, 15\}$ ; see Figure 2.1. For this problem the other edge aggregates are  $M_3 = \{4, 5\}$ ,  $M_4 = \{6\}$ , and  $M_5 = \{7, 8, 9, 18\}$ .

Once the edge aggregates are formed, the Reitzinger-Schöberl prolongation matrix is given by

$$(P)_{ij} = \begin{cases} 1 & \text{if } i \in M_j^+, \\ -1 & \text{if } i \in M_j^-, \\ 0 & \text{otherwise.} \end{cases}$$

Note that although the prolongation depends on the relative orientation of every coarse edge and the corresponding fine edges and although the coarse edge orientation can be chosen arbitrarily, the actual two-level multigrid method does not depend on this latter choice<sup>6</sup>.

It is important to note that every edge belongs either to an edge aggregate or to an auxiliary nodal aggregate. This is because in the considered setting every edge connects either two different nodal aggregates or a nodal aggregate with the nodes in  $\Gamma_D$ , and in both of these cases it belongs to an edge aggregate or it connects two nodes of the same nodal aggregate and therefore belongs to it.

As pointed out in Section 2.3, the projector  $\pi$  is determined by the set  $\mathcal{E}_\pi^0$  of edge indices  $i$  for which  $(\pi\mathbf{v})_i = 0$  for every  $\mathbf{v}$ , and the coefficients  $\alpha_{i,f}$  in (2.13) that define  $(\pi\mathbf{v})_i$  for the remaining indices. The index set  $\mathcal{E}_\pi^0$  is constructed here by picking one edge from every edge aggregate (e.g., edge 14 for aggregate  $M_1 = \{11, 14\}$  from our toy problem), and further, by picking the edges associated to a spanning tree of every auxiliary nodal aggregate (e.g., edges 1, 2 and 13 for the auxiliary aggregate  $K_1$  from our toy problem). A possible set of edges corresponding to  $\mathcal{E}_\pi^0$  for our toy problem is given by

$$\mathcal{E}_\pi^0 = \{1, 2, 3, 4, 6, 7, 13, 14, 15, 16, 17\}.$$

These edges are further represented in bold in Figure 2.1 (center).

This construction of  $\mathcal{E}_\pi^0$  ensures that for any edge outside  $\mathcal{E}_\pi^0$  there is at least one *closed local* path that consists of this edge and the edges in  $\mathcal{E}_\pi^0$ . More precisely, for every edge outside  $\mathcal{E}_\pi^0$  and that belongs to an auxiliary nodal aggregate, a closed path can be constructed that consists of this edge and some edges in  $\mathcal{E}_\pi^0$  from the spanning tree of this auxiliary nodal aggregate. For our toy problem edge 10 belongs to  $K_1$  and forms a closed path together with edges 1, 2, and 13. Likewise, for every edge outside  $\mathcal{E}_\pi^0$  and that belongs to an edge aggregate a closed path can be constructed that consists of this edge, the edge in  $\mathcal{E}_\pi^0$  that belongs to this edge aggregate, and some edges in  $\mathcal{E}_\pi^0$  that belong to spanning trees of the auxiliary nodal

---

<sup>6</sup>More precisely, two prolongations  $P_1$  and  $P_2$  that differ only by the choice of the coarse edge orientation satisfy  $P_1 = P_2 O$ , where  $O = \text{diag}(o_i)$ ,  $o_i = \pm 1$ . Such prolongations lead to the same preconditioner  $B_{\text{TG}}$  in (2.1) provided that the Galerkin formula is used for the coarse grid matrix.

aggregates connected by this edge. For our toy problem edge 11 belongs to  $M_1$  and forms a closed path together with edges 2, 3, and 14.

The coefficients  $\alpha_{i,f}$  in (2.13) that define  $(\pi\mathbf{v})_i$  for the edge  $i$  outside  $\mathcal{E}_\pi^0$  are chosen based on the closed local path associated to  $i$ . More precisely, we consider the set of faces  $F_i$  that form a surface of which this closed path is the boundary. For instance, for our toy problem edge 9 is associated with the closed path<sup>7</sup>  $9 \rightarrow (-7) \rightarrow 16 \rightarrow 17$  which represents the boundary of the union of faces 7 and 8, and hence  $F_9 = \{7, 8\}$ . The coefficients  $\alpha_{i,f}$  for the faces  $f \in F_i$  are set to  $\pm 1$ , the sign depending on the relative orientation of the face and the orientation of the surface as induced (in the sense of Stokes' theorem) by the orientation of the corresponding closed path; this latter orientation is chosen compatible with the orientation of the edge  $i$ . The coefficients  $\alpha_{i,f}$  for the faces  $f$  outside  $F_i$  are set to 0. Hence, for the toy problem the fact that  $F_9 = \{7, 8\}$  implies, taking into account the orientation, that

$$(\pi\mathbf{v})_9 = \mathbf{v}^T \left( \sum_f \alpha_{9,f} \mathbf{c}_f \right) = \mathbf{v}^T (\mathbf{c}_8 + \mathbf{c}_9) = -v_7 + v_9 + v_{15} + v_{16}.$$

Note that by construction, the indices of the nonzero entries of the vector  $\sum_f \alpha_{i,f} \mathbf{c}_f$  correspond to the edges of the closed path associated to  $i$ , and the values of the nonzero entries are  $\pm 1$  depending on the relative orientation of the corresponding edges and the closed path. Since the edges of the closed path belong to  $\{i\} \cup \mathcal{E}_\pi^0$  by construction, the corresponding  $\pi$  is a projector. On the other hand, the fact that the path is local, typically entails that  $\alpha_{i,f}$  are nonzero only for a handful of faces  $f$ . We note that the closed path can always be chosen so that it crosses each grid node at most once since otherwise a smaller path can be extracted from it that also contains a given edge and the edges from  $\mathcal{E}_\pi^0$ .

The following theorem confirms that the projector  $\pi$ , obtained as explained in the preceding paragraphs, fits the requirements of Theorem 2.1.

**THEOREM 2.2.** *Let  $A$  be a symmetric matrix arising from the discretization of the boundary value problem (1.2) with  $\beta = 0$ , and let  $P$  be the prolongation of Reitzinger-Schöberl type as introduced in this section. Then the projector  $\pi$ , constructed as described earlier in this section, satisfies the requirements of Theorem 2.1.*

*Sketch of the proof.* We need to show that the condition (2.2) of Theorem 2.1 is satisfied or, alternatively, that  $\mathcal{N}(A) \subset \mathcal{N}(\pi)$ ,  $\mathcal{R}(P) \subset \mathcal{N}(\pi)$  and the dimension of  $\mathcal{N}(\pi)$  is given by (2.5).

The condition  $\mathcal{N}(A) \subset \mathcal{N}(\pi)$  follows from (2.13), the fact that  $\mathcal{N}(A)$  is spanned by the discrete gradients  $\mathbf{g}_i$ ,  $i = 1, \dots, n^{(n)}$ , and the fact that  $\mathbf{g}_i^T \mathbf{c}_f = 0$ , for all  $i = 1, \dots, n^{(n)}$ , and all  $f = 1, \dots, n^{(f)}$ . Regarding the condition  $\mathcal{R}(P) \subset \mathcal{N}(\pi)$ , we note that the column  $\mathbf{p}_j$  of  $P$  is nonzero only for edges belonging to a corresponding edge aggregate  $M_j$ , whereas  $\sum_f \alpha_{i,f} \mathbf{c}_f$  is nonzero only on the closed local path that includes, by construction, zero or two edges of  $M_j$ . The latter case arises if  $i \in M_j \setminus \mathcal{E}_\pi^0$ , and in this case

$$(\pi\mathbf{p}_j)_i = \mathbf{p}_j^T \left( \sum_f \alpha_{i,f} \mathbf{c}_f \right) = 0,$$

since the only two nonzero contributions to the vector product of the middle term, corresponding to the edge  $i$  and the only edge in  $M_j \cap \mathcal{E}_\pi^0$ , cancel out; see Figure 3.1 (left) for an illustration.

The dimension of  $\mathcal{N}(\pi)$  according to (2.14) is given by the number  $|\mathcal{E}_\pi^0|$  of edge indices set to zero by the projector, which corresponds to one edge per edge aggregate ( $n_c$  edges in

<sup>7</sup> $(-7)$  means here that the orientation of the edge 7 is opposite to the one of the path.

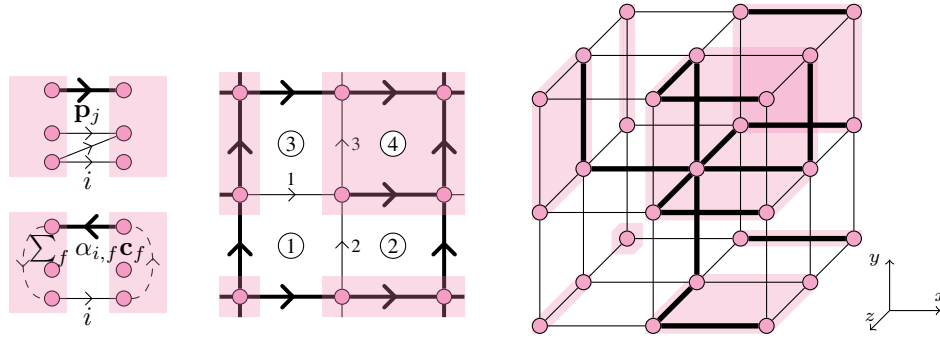


FIG. 3.1. Representation of the edges corresponding to typical nonzero entries of  $\mathbf{p}_j$  (left, upper) and of  $\sum_f \alpha_{i,f} \mathbf{c}_f$  for an edge  $i$  belonging to an edge aggregate  $M_j$  (left, lower) as well as the aggregation pattern for square  $2 \times 2$  (center) and cube  $2 \times 2 \times 2$  (left) auxiliary nodal aggregates on the Cartesian grid. The local face indices are circled, whereas the edges whose indices belong to the considered subset  $\mathcal{E}_\pi^0$  are bold. For the left figure, the edge orientation highlights the sign of the nonzero entries. The central and right figures depict a repetitive pattern, and some redundant elements are present on both figures.

total), augmented by the number of edges in the spanning trees of all auxiliary node aggregates ( $n^{(n)} - n_c^{(n)}$  edges in total since every node aggregate has a connected graph), that is,

$$\dim(\mathcal{N}(\pi)) = n^{(n)} + n_c - n_c^{(n)}.$$

This corresponds to the condition (2.5) since  $\dim(\mathcal{N}(A)) = n^{(n)}$  (or  $n^{(n)} - 1$  if  $\Gamma_N = \partial\Omega$ ),  $\text{rank}(A_c) = n_c - \dim(\mathcal{N}(A_c))$  and, for the Reitzinger-Schöberl prolongation,  $\dim(\mathcal{N}(A_c)) = n_c^{(n)}$  (or  $n_c^{(n)} - 1$  if  $\Gamma_N = \partial\Omega$ ).  $\square$

**3. Outcomes for structured grids.** In this and the following section we present a few outcomes of the just stated analysis. The considered problems include the discretizations of the boundary value problem (1.2) with  $\beta = 0$  in two and three dimensions on structured Cartesian and (possibly) unstructured simplex grids. In particular, in this section we consider the discretizations on structured Cartesian grids, including grids with square/cubic meshes, as well as stretched meshes, problems with constant isotropic and anisotropic coefficients, as well as coefficients with jumps. The outcomes are stated with respect to specific variants of the Reitzinger-Schöberl two-level multigrid method.

To simplify the analysis, the discretizations on structured Cartesian grids are considered here with the prolongation based on auxiliary nodal aggregates that have identical shape, are tiled in a periodic fashion, and fill the whole discretization grid outside  $\Gamma_D$ ; see Figure 3.1 (center) and (right) for the possible aggregates patterns. This implies, in particular, that the considered grids are such that the number of their nodes in one coordinate direction is divisible (up to a boundary effect) by the number of nodes of the nodal aggregate in this direction.

The above simplification allows us to restrict the analysis to the neighborhood of one nodal aggregate. Indeed, since both the grid structure and the nodal aggregates (coarse nodes) are arranged periodically, the different contributions in the decompositions (2.12) and (2.15) can be regrouped into local contributions which are identical (up to a translation) from one aggregate (coarse node) to another. Of course, this approach does not explicitly account for the contributions at the boundary; such contributions are ignored here since they do not bring any additional insight, whereas at the same time they are not expected to significantly change

the resulting estimate<sup>8</sup>. If required, the contributions at the boundary can be obtained similarly to what is done for the toy problem.

The results of the analysis below are regularly complemented with numerical experiments. In these experiments we report the effective condition number

$$(3.1) \quad \kappa_{\text{eff}} = \max_{\lambda \in \sigma(B_{\text{TG}}A)} \lambda / \min_{\substack{\lambda \in \sigma(B_{\text{TG}}A) \\ \lambda \neq 0}} \lambda,$$

where  $\sigma(B_{\text{TG}}A)$  is the set of all eigenvalues of  $B_{\text{TG}}A$ . This condition number satisfies  $\kappa_{\text{eff}} \leq \kappa_{\pi}$  and can be used in the same way as  $\kappa_{\pi}$  in the estimate (2.3) which bounds the convergence rate of a two-level multigrid method; see [25] for details. The parameter  $\kappa_{\text{eff}}$  is computed here using the Matlab `eigs` routine.

The numerical results highlight, amongst others, the fact that  $\kappa_{\text{eff}}$  changes little when going from the small values of  $\beta > 0$  to  $\beta = 0$  despite the fact that the eigenvalues associated with the discrete gradients are excluded from the evaluation of  $\kappa_{\text{eff}}$  if  $\beta = 0$ . Since a Hiptmair smoother [11] is considered for the numerical experiments, this actually means that the extra smoothing iteration it performs in the case  $\beta > 0$  on the system projected on the subspace of the discrete gradients is efficient enough. To be a bit more specific, we note that the Hiptmair smoother is used in the numerical experiments with Gauss-Seidel as a one-level iteration for both the original and the projected system; since this smoother is not symmetric, to keep the symmetry of the preconditioner, we used a forward sweep for the pre-smoothing and a backward sweep for the post-smoothing; see [20] for details.

**3.1. Isotropic case.** Here we estimate the convergence parameter  $\tilde{\kappa}_{\pi}$  for a two-level multigrid method applied to the discretized isotropic boundary value problem (1.2) with constant coefficients  $\beta = 0$  and  $\mu = 1$ . The discretization is performed on a Cartesian grid with a square (in 2D) or a cube (in 3D) mesh of size  $h$ . The multigrid prolongation is of Reitzinger-Schöberl type with square (in 2D) or cube (in 3D) auxiliary nodal aggregates tiled periodically as shown in Figure 3.1 (center) and (right). The analysis is performed locally for a typical aggregate not at the boundary.

Considering first the two-dimensional case, we note that the decomposition (2.12) for the denominator of (2.7) holds as equality with  $\gamma_f = \gamma := 1/h^2$ . In particular, here we consider the local contribution corresponding to the faces locally numbered 1–4 in Figure 3.1 (center). The local contribution can thus be written as

$$(3.2) \quad (\mathbf{v}^T A \mathbf{v})|_{\text{loc}} = \gamma \sum_{f=1}^4 (\mathbf{v}^T \mathbf{c}_f)^2.$$

Regarding the decomposition (2.15) for the numerator of (2.7), the edges corresponding to a possible choice of  $\mathcal{E}_{\pi}^0$  for the considered tiling of square auxiliary nodal aggregates are represented in bold in Figure 3.1 (center). On the other hand, the edges outside  $\mathcal{E}_{\pi}^0$  are locally numbered from 1 to 3 in the figure. The local contribution to the decomposition (2.15) is given by

$$(3.3) \quad \begin{aligned} (\mathbf{v}^T \pi^T D \pi \mathbf{v})|_{\text{loc}} &= d_1 (\mathbf{v}^T (\mathbf{c}_3 + \mathbf{c}_4))^2 + d_2 (\mathbf{v}^T \mathbf{c}_2)^2 + d_3 (\mathbf{v}^T \mathbf{c}_4)^2 \\ &\leq 2d_1 (\mathbf{v}^T \mathbf{c}_3)^2 + d_2 (\mathbf{v}^T \mathbf{c}_2)^2 + (2d_1 + d_3) (\mathbf{v}^T \mathbf{c}_4)^2, \end{aligned}$$

<sup>8</sup>Note that the local Fourier analysis technique used in [7] for the discrete boundary value problems (1.2) also amounts to ignore the boundary conditions, and the resulting estimates still accurately reproduce the actual convergence rate.

TABLE 3.1

*The parameter  $\kappa_{\text{eff}}$  and the associated upper bound (2.6) for the problem and the two-level multigrid methods considered in Section 3.1 for different values of  $\beta$  and  $h$  in two and three dimensions.*

| 2D       |             |                   |             |      |                  |
|----------|-------------|-------------------|-------------|------|------------------|
| $h^{-1}$ | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ | Jac  | $\omega_J^{-1}6$ |
| 101      | 4.26        | 4.26              | 4.14        | 4.90 | 24               |
| 201      | 4.33        | 4.32              | 4.26        | 4.91 | 24               |
| 401      | 4.35        | 4.34              | 4.31        | 4.91 | 24               |

| 3D       |             |                   |             |      |                   |
|----------|-------------|-------------------|-------------|------|-------------------|
| $h^{-1}$ | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ | Jac  | $\omega_J^{-1}72$ |
| 11       | 2.25        | 2.25              | 2.20        | 3.91 | 216               |
| 21       | 3.15        | 3.15              | 3.08        | 4.52 | 216               |
| 41       | 3.73        | 3.73              | 3.65        | 4.79 | 216               |

where  $d_i = 2\gamma$  is the diagonal entry of  $A$  associated to the local edge  $i$ ,  $i = 1, 2, 3$ . Therefore, since the other groups of faces have the same type of contributions (possibly except at the boundary), one ends up with

$$(3.4) \quad \tilde{\kappa}_\pi \approx \tilde{\kappa}_\pi|_{\text{loc}} \leq \frac{2d_1 + d_3}{\gamma} = 6.$$

Regarding the three-dimensional case, we consider the local contribution of the 24 faces that play the same role as the local faces 1–4 in two-dimensions. These faces are depicted in Figure 3.1 (right) along with some other faces (those belonging to the left, bottom, and back surfaces of the represented region). The pattern of auxiliary nodal  $2 \times 2 \times 2$  aggregates is depicted with shaded cubes on the figure, whereas the corresponding set  $\mathcal{E}_\pi^0$  of edges is represented in bold. For this configuration one may similarly show (see [19] for more details) that

$$\tilde{\kappa}_\pi \approx \tilde{\kappa}_\pi|_{\text{loc}} \leq 72.$$

We report in Table 3.1 the results for the two-level multigrid method in the considered setting. More specifically, the problem is defined here on the unit square  $[0, 1]^2$  (in 2D) or the unit cube  $[0, 1]^3$  (in 3D) with  $\Gamma_D = \partial\Omega$ . The reported quantities are the effective condition number  $\kappa_{\text{eff}}$  for the two-level multigrid method with a symmetrized Gauss-Seidel iteration ( $\beta = 0$ ), symmetrized Hiptmair smoother based on a Gauss-Seidel iteration ( $\beta > 0$ ), symmetrized weighted Jacobi iteration for  $\beta = 0$  (labeled Jac), as well as the upper bound (2.6) combined with  $\tilde{\kappa}_\pi \leq 6$  (2D) or  $\tilde{\kappa}_\pi \leq 72$  (3D). For the Jacobi smoother, the considered weighting is given by  $\omega_J^{-1} = 4$  in 2D and  $\omega_J^{-1} = 3$  in 3D; in both cases  $\omega_J^{-1} \approx \lambda_{\max}(D^{-1}A)$  is satisfied.

Regarding the reported results, we note that the values of  $\kappa_{\text{eff}}$  for  $\beta = 0$  and  $\beta > 0$  are almost identical and that these values decrease (slightly) with increasing  $\beta$ . We also note that  $\kappa_{\text{eff}}$  seems to be bounded independently of the mesh size  $h$  for both the Gauss-Seidel and the weighted Jacobi smoothers. Note that a weighted Jacobi iteration is considered here both to show that  $\kappa_{\text{eff}}$  is similar for this smoother and for the Gauss-Seidel one, and because the reported upper bound is for a two-level multigrid method based on a weighted Jacobi smoother. Comparing this latter bound with the actual value of  $\kappa_{\text{eff}}$  for the Jacobi smoother, we also note that the latter is significantly better.

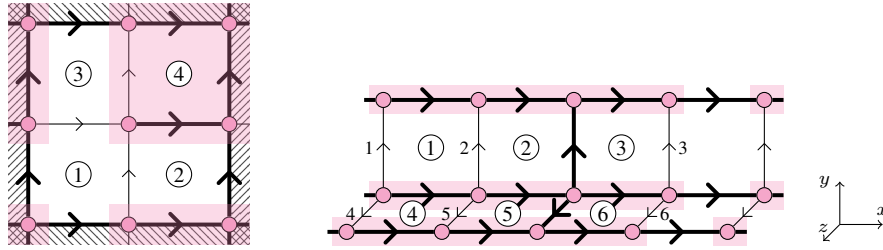


FIG. 3.2. The considered configuration of the square auxiliary nodal aggregates with respect to possible jumps (left) and the considered pattern of line auxiliary nodal aggregates (right) on a Cartesian grid. The local face indices are circled, whereas the edges whose indices correspond to the considered subset  $\mathcal{E}_\pi^0$  are bold.

**3.2. Jumps.** Here we consider the effect of jumps in the coefficient  $\mu$  on the convergence of the two-level multigrid of Reitzinger-Schöberl type. We restrict our analysis to two-dimensional problems with piecewise constant coefficients discretized on a Cartesian grid of mesh size  $h$  and assume that the jumps are located on the grid lines. We also assume that the jumps are resolved well enough in that no two jumps lay on two parallel adjacent grid lines. Regarding the prolongation, we consider the one based on the same tiling of square aggregates as in the previous section. Here also, the analysis is made locally.

We further restrict our attention to the group of four faces depicted in Figure 3.2 (left), with the hatched regions on the figure being the ones where the coefficient  $\mu$  has a value which is possibly different from the one for the four faces, which may further differ from the region with one hatched pattern to the region with another. This situation is quite general since, according to the assumptions made in the previous paragraph, there are at most two discontinuity lines crossing at a node, such a node necessarily belongs to an auxiliary nodal aggregate, and some other discontinuities may only be located at least two grid lines further.

For the configuration depicted in Figure 3.2 (left), the local contributions to the decompositions (2.12) and (2.15) are also given (up to a  $\mu^{-1}$  factor, with  $\mu$  associated to the region of faces 1–4) by (3.2) and (3.3), respectively. This is because the edges along the jumps are included into the set  $\mathcal{E}_\pi^0$ , and therefore all the contributions are proportional to the same  $\mu^{-1}$  coefficient associated to the region of faces 1–4. As a result, the convergence estimate (3.4) for the two-dimensional isotropic case also holds here. More generally, in the considered setting for two-dimensional problems, the set  $\mathcal{E}_\pi^0$  can always be chosen to lay along the jumps, hereby “hiding” their effect.

The impact of the jumps in the coefficient  $\mu$  is further illustrated with the numerical experiments reported in Table 3.2. The corresponding problems are defined on a unit square (2D) or a unit cube (3D) with jumps located in 2D along the grid lines  $x = (1 + h)/2$  and  $y = (1 + h)/2$  and in 3D along the grid planes  $x = (1 + h)/2$ ,  $y = (1 + h)/2$ , and  $z = (1 + h)/2$ . The value of the coefficient  $\mu^{-1}$  for the 2D problem changes by a factor  $10^{\pm 1}$  when crossing the line  $x = (1 + h)/2$  and by a factor  $10^{\pm 2}$  when crossing the line  $y = (1 + h)/2$ , the overall magnitude of  $\mu$  varying from 1 to  $10^{\pm 3}$  through the domain; the + sign in the exponent corresponds to the results in the top part of the table, whereas the – sign corresponds to the results in the bottom part (labeled reversed). Likewise, for the 3D problem the value of the coefficient  $\mu^{-1}$  changes by a factor  $10^{\pm 1}$  when crossing the plane  $x = (1 + h)/2$ , by a factor  $10^{\pm 2}$  when crossing the plane  $y = (1 + h)/2$ , and by a factor  $10^{\pm 4}$  when crossing the plane  $z = (1 + h)/2$ ; in absolute terms  $\mu^{-1}$  varies from 1 to  $10^{\pm 7}$ .

The main conclusion from the results in Table 3.2 is that the presence of jumps have little impact on the convergence in both two and three dimensions despite the analysis was carried out only in the former case. Moreover, the relative location of the aggregates with respect to



TABLE 3.2

*The parameter  $\kappa_{\text{eff}}$  for the problems and the two-level multigrid method considered in Section 3.2 for different values of  $\beta$  and  $h$  in two and three dimensions.*

| 2D            |             |                   |             | 3D            |             |                   |             |
|---------------|-------------|-------------------|-------------|---------------|-------------|-------------------|-------------|
| $h^{-1}$      | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ | $h^{-1}$      | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ |
| 101           | 4.51        | 4.51              | 4.47        | 11            | 2.28        | 2.28              | 2.24        |
| 201           | 4.59        | 4.59              | 4.56        | 21            | 3.48        | 3.47              | 3.40        |
| 401           | 4.60        | 4.60              | 4.59        | 41            | 4.01        | 4.01              | 3.94        |
| 2D (reversed) |             |                   |             | 3D (reversed) |             |                   |             |
| $h^{-1}$      | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ | $h^{-1}$      | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ |
| 101           | 4.23        | 4.23              | 4.12        | 11            | 2.35        | 2.35              | 2.31        |
| 201           | 4.31        | 4.31              | 4.25        | 21            | 3.03        | 3.03              | 2.98        |
| 401           | 4.33        | 4.33              | 4.31        | 41            | 3.66        | 3.66              | 3.59        |

the discontinuity is also of little importance as follows from the comparison of the results from the top and bottom (reversed) parts of the table.

**3.3. Anisotropic case.** Here we consider the convergence of the Reitzinger-Schöberl two-level multigrid method in the presence of anisotropy and/or stretched Cartesian grids. By anisotropy we mean that the boundary value problem (1.2) with  $\tilde{\mu} = \text{diag}(\mu_x, \mu_y, \mu_z)$  is such that the coefficients  $\mu_x, \mu_y$ , and  $\mu_z$  are possibly different, whereas a Cartesian grid of mesh size  $h_x \times h_y \times h_z$  is considered stretched if the mesh sizes  $h_x, h_y$ , and  $h_z$  are possibly different. In what follows we further assume that both the problem coefficients and the mesh sizes are uniform over  $\Omega$  and therefore do not vary from one node (or mesh) of the grid to another.

In this setting, the system matrix  $A$  is assembled from the element matrix  $A_e$  corresponding to a  $h_x \times h_y \times h_z$  brick element. The element matrix is given by

$$A_e = C_e M_e C_e^T, \quad M_e = \frac{\gamma_v}{3} \text{blockdiag}(h_x^2 \mu_x^{-1} \tilde{M}, h_y^2 \mu_y^{-1} \tilde{M}, h_z^2 \mu_z^{-1} \tilde{M}),$$

$$\tilde{M} = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix},$$

where  $\gamma_v = 1/(h_x h_y h_z)$  is the inverse of the volume of the brick element and  $C_e = (\mathbf{c}_{f_e}^{(e)})$  is a  $12 \times 6$  matrix whose columns  $\mathbf{c}_{f_e}^{(e)}$ , defined as in (2.10), are associated with the 6 faces of the brick element with the faces orthogonal to the  $x$  axis being ordered first, those orthogonal to the  $y$  axis being ordered next, and those orthogonal to the  $z$  axis being ordered last. Note that the effects of the coefficient anisotropy and the grid stretching on the element matrix  $A_e$  are equivalent, the overall impact depending on the difference in magnitude of  $h_x^2 \mu_x^{-1}, h_y^2 \mu_y^{-1}$ , and  $h_z^2 \mu_z^{-1}$ .

The expression of the element contribution (2.11) follows from the above expression of the element matrix  $A_e$  and the fact that the smallest eigenvalue of  $\tilde{M}$  is  $1/2$ . The contribution is given by

$$\mathbf{v}_e^T A_e \mathbf{v}_e \geq \frac{\gamma_v}{6} \left( h_x^2 \mu_x^{-1} \sum_{f_e=1}^2 (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2 + h_y^2 \mu_y^{-1} \sum_{f_e=3}^4 (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2 + h_z^2 \mu_z^{-1} \sum_{f_e=5}^6 (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2 \right).$$

The decomposition (2.12) mainly differs from the above element's variant in that each face  $f$  not at the boundary  $\partial\Omega$  corresponds to the local faces  $f_e$  of two different adjacent elements.

Therefore, for the aggregates not at the boundary, the local contribution to the decomposition (2.12) is given by

$$(3.5) \quad (\mathbf{v}^T A \mathbf{v})|_{\text{loc}} \geq \frac{\gamma_v}{3} \left( h_x^2 \mu_x^{-1} \sum_{f \in F_x} (\mathbf{v}^T \mathbf{c}_f)^2 + h_y^2 \mu_y^{-1} \sum_{f \in F_y} (\mathbf{v}^T \mathbf{c}_f)^2 + h_z^2 \mu_z^{-1} \sum_{f \in F_z} (\mathbf{v}^T \mathbf{c}_f)^2 \right),$$

where  $F_x$ ,  $F_y$ , and  $F_z$  are the index sets of the local faces orthogonal to the  $x$ ,  $y$ , and  $z$  axis, respectively.

Let us now specify a suitable two-level multigrid method. Assuming in what follows that  $h_x^2 \mu_x^{-1} \leq h_y^2 \mu_y^{-1} \leq h_z^2 \mu_z^{-1}$ , the considered two-level method is based on the Reitzinger-Schöberl prolongation with auxiliary nodal aggregates aligned in the direction of the  $x$  axis; see Figure 3.2 (right). For non-stretched Cartesian grids (with  $h_x = h_y = h_z$ ) this corresponds to the direction of the weakest coefficient  $\mu_x^{-1}$ . For simplicity, we pick a nodal aggregate of length 4; the derivation of the results for aggregates of different length follows the same lines.

For the considered nodal aggregates the possible set of edges whose indices are in  $\mathcal{E}_\pi^0$  is marked in bold in Figure 3.2 (right) in the neighborhood of one line aggregate. As of the indices outside  $\mathcal{E}_\pi^0$ , they correspond to the edges numbered 1–6 in the figure. The corresponding local contribution to the decomposition (2.15) is given by

$$\begin{aligned} (\mathbf{v}^T \pi^T D \pi \mathbf{v})|_{\text{loc}} &= d_1 (\mathbf{v}^T (\mathbf{c}_1 + \mathbf{c}_2))^2 + d_2 (\mathbf{v}^T \mathbf{c}_2)^2 + d_3 (\mathbf{v}^T \mathbf{c}_3)^2 \\ &\quad + d_4 (\mathbf{v}^T (\mathbf{c}_4 + \mathbf{c}_5))^2 + d_5 (\mathbf{v}^T \mathbf{c}_5)^2 + d_6 (\mathbf{v}^T \mathbf{c}_6)^2, \end{aligned}$$

where the vectors  $\mathbf{c}_f$ ,  $f = 1, \dots, 6$ , correspond to the six local faces in the figure (circled numbers), whereas  $d_1 = d_2 = d_3 = d_y$  and  $d_4 = d_5 = d_6 = d_z$  are the diagonal entries of the system matrix for edges along the  $y$  and  $z$  axis, respectively, with  $d_y$  and  $d_z$  given by

$$d_y = 4/3 \gamma_v (h_x^2 \mu_x^{-1} + h_z^2 \mu_z^{-1}), \quad d_z = 4/3 \gamma_v (h_x^2 \mu_x^{-1} + h_y^2 \mu_y^{-1}).$$

Noting that the local faces 1–3 belong to  $F_z$ , whereas the faces 4–6 belong to  $F_y$ , one has

$$(3.6) \quad (\mathbf{v}^T \pi^T D \pi \mathbf{v})|_{\text{loc}} \leq 3d_z \sum_{f \in F_y} (\mathbf{v}^T \mathbf{c}_f)^2 + 3d_y \sum_{f \in F_z} (\mathbf{v}^T \mathbf{c}_f)^2,$$

where the factor 3 is determined by the aggregate length; for an aggregate of length  $\ell$  it is given by  $c(\ell)(c(\ell) + 1)/2$  with  $c(\ell) = \lceil (\ell - 1)/2 \rceil$  being the maximal number of times a local face  $f$  is enclosed by a path associated to an edge outside  $\mathcal{E}_\pi^0$ ; here  $c(\ell) = c(4) = 2$ . Combining the local contributions (3.5) and (3.6) one has

$$\tilde{\kappa}_\pi \approx \tilde{\kappa}_\pi|_{\text{loc}} \leq \frac{9}{\gamma_v} \max \left( \frac{d_z}{h_y^2 \mu_y^{-1}}, \frac{d_y}{h_z^2 \mu_z^{-1}} \right) \leq 24,$$

where the latter inequality stems from

$$d_y \leq 8/3 \gamma_v h_z^2 \mu_z^{-1}, \quad d_z \leq 8/3 \gamma_v h_y^2 \mu_y^{-1},$$

which itself comes from the assumption  $h_x^2 \mu_x^{-1} \leq h_y^2 \mu_y^{-1} \leq h_z^2 \mu_z^{-1}$  made earlier.

We report in Table 3.3 the numerical experiments with the considered Reitzinger-Schöberl two-level multigrid based on line nodal aggregates of length 4 aligned along the  $x$  axis; this method is here applied to the boundary value problem (1.2) defined on a unit cube  $[0, 1]^3$

TABLE 3.3

*The parameter  $\kappa_{\text{eff}}$  for the problems and the two-level multigrid method considered in Section 3.3 for different values of  $\beta$  and  $h$  in three dimensions.*

| $h^{-1}$ | 3D (weak)   |                   |             | 3D (strong) |                   |             |
|----------|-------------|-------------------|-------------|-------------|-------------------|-------------|
|          | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ |
| 9        | 3.32        | 3.32              | 3.32        | 8.55        | 8.55              | 8.07        |
| 17       | 4.26        | 4.26              | 4.26        | 28.7        | 28.7              | 26.7        |
| 33       | 4.82        | 4.82              | 4.82        | 98.0        | 97.9              | 90.8        |

with  $\Gamma_D = \partial\Omega$  and discretized on a Cartesian grid of mesh size  $h_x = h_y = h_z = h$ , with  $h = 1/(4p + 1)$  for some integer  $p$ . The problem coefficients are  $\mu_x^{-1} = 1$ ,  $\mu_y^{-1} = 10^{\pm 2}$ , and  $\mu_z^{-1} = 10^{\pm 4}$ ; the + sign in the exponent corresponds to the results in the left part of the table (labeled weak since then  $\mu_x^{-1} < \mu_y^{-1} < \mu_z^{-1}$ ), whereas the – sign corresponds to the results in the right part (labeled strong). Note that in the first case the aggregates are aligned in the direction of the weakest coefficient, whereas they are aligned with the strongest coefficient in the second case, hence the labeling. Now, the results reported in the table do not only corroborate the above analysis in that the two-level convergence seems indeed bounded if the nodal aggregates are aligned in the direction of the weakest coefficient but further indicate that if the aggregates are not aligned in the direction of the weakest coefficient, then the convergence may deteriorate.

**4. Outcomes for unstructured grids.** Here we analyze the convergence of a Reitzinger-Schöberl two-level multigrid method for the boundary value problem (1.2) discretized on an unstructured simplex grid. The boundary value problem is considered here in two and three dimensions with  $\beta = 0$ ,  $\mu = 1$ , and with  $\Gamma_N = \partial\Omega$ ; the assumption on the boundary condition is only made to keep the discussion simple.

The main convergence result is stated in Theorem 4.7 below. The theorem shows that in the considered setting the convergence parameter  $\tilde{\kappa}_\tau$  associated with any two-level method of Reitzinger-Schöberl type can be bounded independently of the mesh size. The main assumptions are the admissibility and the uniformity of the discretization grid and the bounded size of the auxiliary nodal aggregates. An additional technical assumption is on the topology of the discretization grid.

Let us now clarify some of the just used terms. A simplex grid is *admissible* if any two simplexes are either disjoint or have in common a mesh node, a mesh edge, or (in 3D) a mesh face. The grid of mesh size  $h$  is  $\tau$ -uniform if every simplex of the grid is contained in a disk (in 2D) or a ball (in 3D) of radius  $h$  and contains a disk (in 2D) or a ball (in 3D) of radius  $\tau h$ .

To prove Theorem 4.7 we use the three following lemmas. The proof of the first lemma can be found in Appendix A of [19]. It follows from (2.9) and the fact that in two dimensions  $M_e$  is a scalar given by the inverse of the area of the element, whereas in three dimensions the  $4 \times 4$  matrix  $M_e$  is a Raviart-Thomas mass matrix of the element  $e$  whose conditioning can be related to the regularity parameter  $\tau$ .

LEMMA 4.1 ([19]). *The element matrix  $A_e$  corresponding to the boundary value problem (1.2) in two ( $d = 2$ ) or three ( $d = 3$ ) dimensions with  $\beta = 0$ ,  $\mu = 1$ , and for a simplex element  $e$  of a  $\tau$ -uniform grid of mesh size  $h$  satisfies*

$$(4.1) \quad c_m(\tau)h^{d-4} \sum_{f_e=1}^{n_e^{(f)}} (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2 \leq \mathbf{v}_e^T A_e \mathbf{v}_e \leq c_M(\tau)h^{d-4} \sum_{f_e=1}^{n_e^{(f)}} (\mathbf{v}_e^T \mathbf{c}_{f_e}^{(e)})^2$$

for some positive  $c_m(\tau)$ ,  $c_M(\tau)$ , and with  $n_e^{(f)} = 1$  if  $d = 2$ , and  $n_e^{(f)} = 4$  if  $d = 3$ .

The second lemma gives an upper bound on the number of faces sharing an edge in a three-dimensional  $\tau$ -uniform grid.

LEMMA 4.2. *The number of faces sharing an edge in a three-dimensional admissible  $\tau$ -uniform grid is at most  $2\pi\tau^{-1}$ .*

*Proof.* It is enough to show that for a three-dimensional  $\tau$ -uniform grid the angle between two faces of a simplex that share an edge is at least  $\tau$ . To see this, one can orthogonally project the simplex on a plane perpendicular to this edge. The projected simplex is a triangle, and the angle between the faces is also one of the angles of the triangle. Moreover, due to the projection, the resulting triangle has a diameter not greater than that of the simplex. On the other hand, the projection of any ball contained in the simplex is a disk contained in the triangle. The result then follows from the observation that for a triangle with a diameter at most  $2h$  and with a radius of the inscribed disk at least  $\tau h$ , any angle of the triangle is at least  $\tau$ .  $\square$

The third lemma gives an upper bound for the magnitude of the diagonal entries of the system matrix.

LEMMA 4.3. *Let  $d_i$  be a diagonal entry of a matrix  $A$  arising from the discretization of the boundary value problem (1.2) in two ( $d = 2$ ) and three ( $d = 3$ ) dimensions with  $\beta = 0$ ,  $\mu = 1$ , and  $\Gamma_N = \partial\Omega$  on a  $\tau$ -uniform admissible simplex grid. Then*

$$(4.2) \quad d_i \leq c_M(\tau)h^{d-4} (d-1)2\pi\tau^{-1}.$$

*Proof.* Let  $\mathbf{e}_i$  be the  $i$ th canonical basis vector. Using equations (2.8),  $\mathbf{c}_f = \pm T_e \mathbf{c}_{f_e}^{(e)}$  with  $f_e$  being the element face index corresponding to the global face index  $f$ , and (4.1), one has

$$\begin{aligned} d_i &= \sum_{e=0}^{n^{(e)}} \mathbf{e}_i^T T_e A_e T_e^T \mathbf{e}_i \leq c_M(\tau)h^{d-4} \sum_{e=0}^{n^{(e)}} \sum_{f_e=1}^{n_e^{(f)}} (\mathbf{e}_i^T T_e \mathbf{c}_{f_e}^{(e)})^2 \\ &= c_M(\tau)h^{d-4} \sum_{e=0}^{n^{(e)}} \sum_{f_e=1}^{n_e^{(f)}} (\mathbf{e}_i^T \mathbf{c}_f)^2 \leq c_M(\tau)h^{d-4} (d-1) \sum_{f=1}^{n^{(f)}} (\mathbf{e}_i^T \mathbf{c}_f)^2, \end{aligned}$$

where the last inequality comes from the fact that each global face  $f$  corresponds to at most  $d-1$  element faces  $f_e$ . The sum in the last term of the above inequality represents the number of faces sharing a given edge  $i$ ; it is bounded by 2 in two dimensions and, according to Lemma 4.2, by  $2\pi\tau^{-1}$  in three dimensions. Since  $\pi\tau^{-1} > \tau^{-1} > 1$ , combining the above yields the desired result.  $\square$

Let us now state the aforementioned technical assumption.

ASSUMPTION 4.4. *The considered discretization grid is such that any closed path of at most  $\ell$  edges represents the boundary of a surface formed by a union of at most  $c_a(\ell)$  faces.*

Note that the above assumption is not unrealistic as is highlighted with the following lemma. In this lemma, a two-dimensional grid is said to be *without holes* if any closed path of edges enclose a surface of faces of which the path is the boundary; that is, there is no extra boundary due to “holes” in the grid.

LEMMA 4.5. *Assumption 4.4 holds*

- (a) *with  $c_a(\ell) = \ell^2/(\pi\tau)^2$  for any two-dimensional  $\tau$ -uniform grid without holes of mesh size  $h$ ;*
- (b) *with  $c_a(\ell) = 4 \cdot k(\ell+1)^3$  for any three-dimensional grid topologically equivalent to a structured simplex grid obtained by subdividing every cube into  $k$  simplexes in a standard way, with hence  $k = 5, 6$ .*

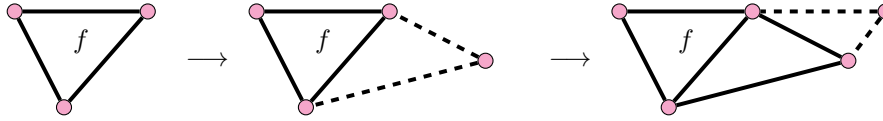


FIG. 4.1. Appending faces of the grid to a given face  $f$ ; the possible first two steps.

*Sketch of the proof.* We prove (a) by noting that a closed path of  $\ell$  edges has a length at most  $2h\ell$ , hence it “encloses” an area of at most  $(h\ell)^2\pi^{-1}$  and therefore has at most  $c_a(\ell) = (h\ell)^2\pi^{-1}/(\tau^2h^2\pi) = \ell^2/(\pi\tau)^2$  grid triangles. Since the grid is without holes, such triangles form a surface of which the path is a boundary.

Regarding (b), we note that the proof can be limited to a three-dimensional structured simplex grid itself excluding the topologically equivalent grids since Assumption 4.4 only involves the grid topology. For such a grid any closed path of  $\ell$  edges can always be included in a cube whose edges are parallel to the grid lines and are of length at most  $h\ell$ , where  $h$  is the mesh size. Moreover, for a standard subdivision, there is always a surface formed of faces for which the closed path is the boundary and which is also included in the same cube. Assuming that the structured grid is obtained by subdividing every mesh cube into  $k$  simplexes, a cube of edge length  $h\ell$  contains at most  $k(\ell + 1)^3$  elements and therefore at most  $c_a(\ell) = 4 \cdot k(\ell + 1)^3$  faces.  $\square$

The following lemma is also useful for the proof of our main result.

LEMMA 4.6. *In an admissible  $\tau$ -uniform grid in two or three dimensions, the number of different surfaces that can be formed starting from a given face and using at most  $s$  faces is bounded above by*

$$c_b(\tau, s) = \sum_{k=1}^s \prod_{j=1}^k (2j + 1) 2\pi\tau^{-1}.$$

*Proof.* We determine the latter number by counting all possible surfaces; this is further illustrated in Figure 4.1.

Since by Lemma 4.2 every edge can be shared by at most  $2\pi\tau^{-1} > 2$  faces of the grid (including the case of two dimensions), and every face has 3 edges, we can append to the face  $f$  at most  $3 \cdot 2\pi\tau^{-1}$  different faces of the grid, and to every such construct having 5 edges (one edge is common), we can append at most  $5 \cdot 2\pi\tau^{-1}$  other faces of the grid (see Figure 4.1), and so on until  $s$  faces are reached; the overall number of possible surfaces is at most

$$3 \cdot 2\pi\tau^{-1} + 3 \cdot 5 \cdot (2\pi\tau^{-1})^2 + \dots + \prod_{j=1}^s (2j + 1) 2\pi\tau^{-1} = c_b(\tau, s). \quad \square$$

We now prove the main convergence result.

THEOREM 4.7. *Let  $A$  be a matrix arising from the discretization of the boundary value problem (1.2) in two ( $d = 2$ ) and three ( $d = 3$ ) dimensions with  $\beta = 0$ ,  $\mu = 1$ , and  $\Gamma_N = \partial\Omega$  on a  $\tau$ -uniform admissible simplex grid satisfying Assumption 4.4, and set  $D = \text{diag}(A)$ . Let  $\pi$  be the projector constructed as in Section 2.5 based on auxiliary nodal aggregates of size at most  $k$ . Then the convergence parameter  $\tilde{\kappa}_\pi$  defined in (2.7) satisfies*

$$\tilde{\kappa}_\pi \leq C(\tau, k, c_a(2k)),$$

where  $C(\tau, k, c_a(2k))$  is defined in (4.7).

*Proof.* The starting point of the proof is the inequality (2.16), which computes the maximum over all the faces  $f$  of the grid of the quotients of the quantities  $\alpha_f^2$  and  $\gamma_f$  defined in

Section 2.3. The proof proceeds by bounding above  $\alpha_f^2$  and bounding below  $\gamma_f$  as a function of  $\tau$ ,  $k$ ,  $c_a(2k)$ , and  $h$ . Besides, the considered assumptions enable the use of Lemmas 4.1, 4.2, 4.3, and 4.6.

Considering first the value of  $\gamma_f$ , we note by comparing (2.11) and (4.1) that here one can choose  $\gamma_{f_e}^{(e)} = c_m(\tau)h^{d-4}$ . Then

$$(4.3) \quad \gamma_f = \sum_e \gamma_{f_e}^{(e)} \geq c_m(\tau)h^{d-4},$$

as there is at least one element  $e$  whose local face  $f_e$  corresponds to the global face  $f$ . Regarding the value of  $\alpha_f^2$ , it satisfies (see Section 2.3)

$$(4.4) \quad \alpha_f^2 = \sum_{i \notin \mathcal{E}_\pi^0} m_i d_i \alpha_{i,f}^2 \leq \max_{i \notin \mathcal{E}_\pi^0} (m_i d_i) \sum_{i \notin \mathcal{E}_\pi^0} \alpha_{i,f}^2,$$

where the meaning of (and a bound for)  $m_i$ ,  $d_i$ , and  $\sum_{i \notin \mathcal{E}_\pi^0} \alpha_{i,f}^2$  is now discussed.

First,  $d_i$  is the  $i$ th diagonal entry of  $A$  and is bounded above with the inequality (4.2) from Lemma 4.3. Next, since  $\pi$  is constructed as in Section 2.5,  $m_i$  is the number of faces in  $F_i$ , that is, a number of faces that form a surface for which the boundary is given by the closed path corresponding to the  $i$ th edge,  $i \notin \mathcal{E}_\pi^0$ . Since the considered auxiliary nodal aggregates have at most  $k$  nodes, any closed path involves the nodes of at most two connected nodal aggregates and any node appears at most once in this closed path, the number of edges composing this closed path is at most  $2k$ . Then, according to Assumption 4.4, it holds that

$$(4.5) \quad m_i \leq c_a(2k).$$

As of the sum  $\sum_{i \notin \mathcal{E}_\pi^0} \alpha_{i,f}^2$ , assuming again that  $\pi$  is constructed as in Section 2.5, it is given by  $|\{i \notin \mathcal{E}_\pi^0 | f \in F_i\}|$ , that is, the number of edges  $i$  outside  $\mathcal{E}_\pi^0$  whose associated closed path is the boundary of the surface of mesh faces  $F_i$  one of which is  $f$ . Note that, as already mentioned, each set  $F_i$  has at most  $c_a(2k)$  faces. Moreover, two surfaces corresponding to  $F_i$  and  $F_j$  are necessarily different since the only edges outside  $\mathcal{E}_\pi^0$  that form their boundary are  $i$  and  $j$ , respectively. Hence, it follows from Lemma 4.6 that

$$(4.6) \quad \sum_{i \notin \mathcal{E}_\pi^0} \alpha_{i,f}^2 \leq c_b(\tau, c_a(2k)).$$

Using the bounds (4.2), (4.5) and (4.6) in (4.4) gives the upper bound for the value of  $\alpha_f^2$  which, together with (4.3) and (2.16), gives the convergence estimate

$$(4.7) \quad \tilde{\kappa}_\pi \leq c_M(\tau)c_m(\tau)^{-1} (d-1) 2\pi\tau^{-1} c_a(2k)c_b(\tau, c_a(2k)) =: C(\tau, k, c_a(2k)). \quad \square$$

The  $h$ -independent convergence of a two-level method is further illustrated with the numerical experiments reported in Table 4.1. The problem considered is defined on a unit square or a unit cube and discretized on an unstructured grid. The two-level method of Reitzinger-Schöberl is based on auxiliary nodal aggregates containing at most 8 nodes; see [20] for further details.

The analysis above highlights the importance of a uniform grid or at least of a grid formed with shape-regular simplexes. This is further illustrated with the numerical results reported in Table 4.2 for  $\beta = 0.01$ , where we use the two-dimensional setting of the previous numerical experiment except for the discretization grid, which is now structured but partially

TABLE 4.1

*The parameter  $\kappa_{\text{eff}}$  for the problems and the two-level multigrid method considered in Section 4 for different values of  $\beta$  and  $h$  in two and three dimensions.*

| $h_{\text{max}}^{-1}$ | 2D          |                   |             | $h^{-1}$ | 3D          |                   |             |
|-----------------------|-------------|-------------------|-------------|----------|-------------|-------------------|-------------|
|                       | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ |          | $\beta = 0$ | $\beta = 10^{-2}$ | $\beta = 1$ |
| 100                   | 8.73        | 8.73              | 8.69        | 10       | 5.91        | 5.91              | 5.79        |
| 200                   | 8.84        | 8.84              | 8.82        | 20       | 6.02        | 6.02              | 5.92        |
| 400                   | 9.00        | 9.00              | 8.99        | 40       | 6.18        | 6.18              | 6.11        |

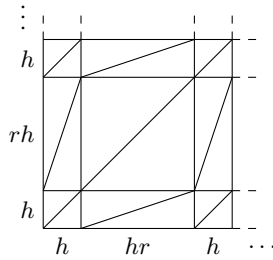


FIG. 4.2. Structured and partially stretched discretization grid as considered in Section 4.

TABLE 4.2

*The parameter  $\kappa_{\text{eff}}$  for the problem discretized on a partially stretched grid and the two-level multigrid method as considered in Section 4; the experiments correspond to  $\beta = 0.01$ .*

| $n_x$ | $n (\times 10^4)$ | $r = 1$ | $r = 10$ | $r = 100$ |
|-------|-------------------|---------|----------|-----------|
| 100   | 3.0               | 9.53    | 18.0     | 153       |
| 200   | 11.9              | 9.57    | 18.0     | 153       |
| 400   | 47.8              | 9.57    | 18.0     | 153       |

stretched, as shown in Figure 4.2. More precisely, the mesh size in the  $x$  direction is  $h$  for the first layer of triangles along  $y$ , then  $rh$  for the second, then  $h$  for the third,  $rh$  for the fourth, and so on, the same applying also for the mesh sizes in the  $y$  direction. The size of the problem is now reported via the number  $n_x$  of grid nodes in every coordinate direction, with  $n = (n_x - 1)^2 + 2n_x(n_x - 1)$ . Note that, although the application of the approach from [20] yields slightly different nodal aggregates for different values of  $r$ , using the same aggregates as for  $r = 1$  in all cases yields the same results.

In such a setting, roughly half of the triangles are badly shaped if  $r$  differs significantly from 1. The results in Table 4.2 show that the condition number increases significantly with  $r$ , and therefore the convergence indeed deteriorates in the presence of badly shaped simplexes.

**5. Conclusions.** We have presented an algebraic analysis of two-level multigrid methods for the solution of linear systems arising from the discretization of the boundary value problem (1.2) with the lowest order edge element method. The analysis is restricted to the case where  $\beta = 0$ , and the resulting linear systems are then singular. The system singularity allows us to simplify the Hiptmair smoother, and therefore makes the analysis possible. We further exploit the knowledge of the range and the null space of the resulting singular system in the design of the projector  $\pi$  which enters the two-level convergence estimate. For the two-level multigrid methods of Reitzinger-Schöberl type, we further explain how to build a possible projector  $\pi$  systematically, based on the auxiliary nodal aggregates. As of the solution of the

boundary value problem (1.2) with  $\beta > 0$ , the numerical experiments confirm that for small values of  $\beta > 0$  the convergence properties of the Reitzinger-Schöberl two-level multigrid are essentially the same as for  $\beta = 0$ .

Regarding the outcomes of the analysis, we show that in a number of situations the variants of the Reitzinger-Schöberl two-level method have an  $h$ -independent convergence. Moreover, for some variants the convergence is independent of the jumps or anisotropy in the problem coefficients, and of the uniform stretch in the grid. On the other hand, the convergence may deteriorate in the presence of simplex elements with low shape regularity.

**Acknowledgments.** The first author acknowledges the *Professeur Visiteur* (Visiting Professor) program of INP Toulouse (France) for supporting his stay in Toulouse during the work on the manuscript. His work was also partially supported by the Fonds de Recherche Scientifique-FNRS (Belgium) under grant n° J.0084.16.

## REFERENCES

- [1] R. ALBANESE AND G. RUBINACCI, *Integral formulation for 3D eddy-current computation using edge elements*, IEE Proc. A, 137 (1998), pp. 457–462.
- [2] D. N. ARNOLD, R. S. FALK, AND R. WINTHER, *Multigrid in  $H(\text{div})$  and  $H(\text{curl})$* , Numer. Math., 85 (2000), pp. 197–217.
- [3] R. BECK, P. DEUFLHARD, R. HIPTMAIR, R. H. W. HOPPE, AND B. WOHLMUTH, *Adaptive multilevel methods for edge element discretizations of Maxwell's equations*, Surveys Math. Indust., 8 (1999), pp. 271–312.
- [4] R. BECK, *Algebraic multigrid by components splitting for edge elements on simplicial triangulations*, Tech. Report SC 99–40, ZIB, Berlin, December 1999.
- [5] P. B. BOCHEV, C. J. GARASI, J. J. HU, A. C. ROBINSON, AND R. S. TUMINARO, *An improved algebraic multigrid method for solving Maxwell's equations*, SIAM J. Sci. Comput., 25 (2003), pp. 623–642.
- [6] T. BOONEN, G. DELIÉGE, AND S. VANDEWALLE, *On algebraic multigrid methods derived from partition of unity nodal prolongators*, Numer. Linear Algebra Appl., 13 (2006), pp. 105–131.
- [7] T. BOONEN, J. VAN LENT, AND S. VANDEWALLE, *Local Fourier analysis of multigrid for the curl-curl equation*, SIAM J. Sci. Comput., 30 (2008), pp. 1730–1755.
- [8] A. BOSSAVIT, *Computational Electromagnetism*, Academic Press, San Diego, 1998.
- [9] A. BRANDT, *Algebraic multigrid theory: the symmetric case*, Appl. Math. Comput., 19 (1986), pp. 23–56.
- [10] R. D. FALGOUT, P. S. VASSILEVSKI, AND L. T. ZIKATANOV, *On two-grid convergence estimates*, Numer. Linear Algebra Appl., 12 (2005), pp. 471–494.
- [11] R. HIPTMAIR, *Multigrid method for Maxwell's equations*, SIAM J. Numer. Anal., 36 (1999), pp. 204–225.
- [12] R. HIPTMAIR AND J. XU, *Nodal auxiliary space preconditioning in  $\mathbf{H}(\text{curl})$  and  $\mathbf{H}(\text{div})$  spaces*, SIAM J. Numer. Anal., 45 (2007), pp. 2483–2509.
- [13] J. J. HU, R. S. TUMINARO, P. B. BOCHEV, C. J. GARASI, AND A. C. ROBINSON, *Toward an  $h$ -independent algebraic multigrid method for Maxwell's equations*, SIAM J. Sci. Comput., 27 (2006), pp. 1669–1688.
- [14] H. IGARASHI, *On the property of the curl-curl matrix in finite element analysis with edge elements*, IEEE Trans. Magnetics, 37 (2001), pp. 3129–3132.
- [15] T. V. KOLEV AND P. S. VASSILEVSKI, *Parallel auxiliary space AMG for  $H(\text{curl})$  problems*, J. Comput. Math., 27 (2009), pp. 604–623.
- [16] J. B. MANGES AND Z. J. CENDES, *A generalized tree-cotree gauge for magnetic field computation*, IEEE Trans. Magnetics, 31 (1995), pp. 1342–1347.
- [17] A. NAPOV AND Y. NOTAY, *Algebraic analysis of aggregation-based multigrid*, Numer. Linear Algebra Appl., 18 (2011), pp. 539–564.
- [18] ———, *An algebraic multigrid method with guaranteed convergence rate*, SIAM J. Sci. Comput., 34 (2012), pp. A1079–A1109.
- [19] A. NAPOV AND R. PERRUSSEL, *Algebraic analysis of two-level multigrid methods for edge elements*, Technical Report, ULB, 2018. Available at <http://homepages.ulb.ac.be/anapov>.
- [20] ———, *Revisiting aggregation-based multigrid for edge elements*, Electron. Trans. Numer. Anal., 51 (2019), pp. 118–134.  
<http://etna.ricam.oeaw.ac.at/vol.51.2019/pp118-134.dir/pp118-134.pdf>
- [21] J.-C. NÉDÉLEC, *Mixed finite elements in  $\mathbb{R}^3$* , Numer. Math., 35 (1980), pp. 315–341.
- [22] Y. NOTAY, *Algebraic analysis of two-grid methods: The nonsymmetric case*, Numer. Linear Algebra Appl., 17 (2010), pp. 73–96.



- [23] ———, *Aggregation-based algebraic multigrid for convection-diffusion equations*, SIAM J. Sci. Comput., 34 (2012), pp. A2288–A2316.
- [24] ———, *Algebraic theory of two-grid methods*, Numer. Math. Theory Methods Appl., 8 (2015), pp. 168–198.
- [25] ———, *Algebraic two-level convergence theory for singular systems*, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 1419–1439.
- [26] R. PERRUSSEL, L. NICOLAS, F. MUSY, L. KRÄHENBÜHL, M. SCHATZMAN, AND C. POIGNARD, *Algebraic multilevel methods for edge elements*, IEEE Trans. Magnetics, 42 (2006), pp. 619–622
- [27] S. REITZINGER AND J. SCHÖBERL, *An algebraic multigrid method for finite element discretizations with edge elements*, Numer. Linear Algebra Appl., 9 (2002), pp. 223–238.
- [28] Z. REN, *Influence of the RHS on the convergence behaviour of the curl-curl equation*, IEEE Trans. Magnetics, 32 (1996), pp. 655–658.
- [29] K. STÜBEN, *An introduction to algebraic multigrid*, in Multigrid, U. Trottenberg, C. W. Oosterlee, and A. Schüller, eds., Academic Press, San Diego, 2001, pp. 413–532.
- [30] H. A. VAN DER VORST, *Iterative Krylov Methods for Large Linear Systems*, Cambridge University Press, Cambridge, 2003.