

## STABILITY OF EXPLICIT ONE-STEP METHODS FOR P1-FINITE ELEMENT APPROXIMATION OF LINEAR DIFFUSION EQUATIONS ON ANISOTROPIC MESHES\*

WEIZHANG HUANG<sup>†</sup>, LENNARD KAMENSKI<sup>‡</sup>, AND JENS LANG<sup>§</sup>

**Abstract.** We study the stability of explicit one-step integration schemes for the linear finite element approximation of linear parabolic equations. The derived bound on the largest permissible time step is tight for any mesh and any diffusion matrix within a factor of  $2(d + 1)$ , where  $d$  is the spatial dimension. Both full mass matrix and mass lumping are considered. The bound reveals that the stability condition is affected by two factors. The first depends on the number of mesh elements and corresponds to the classic bound for the Laplace operator on a uniform mesh. The second factor reflects the effects of the interplay of the mesh geometry and the diffusion matrix. It is shown that it is not the mesh geometry itself but the mesh geometry in relation to the diffusion matrix that is crucial to the stability of explicit methods. When the mesh is uniform in the metric specified by the inverse of the diffusion matrix, the stability condition is comparable to the situation with the Laplace operator on a uniform mesh. Numerical results are presented to verify the theoretical findings.

**Key words.** finite element method, anisotropic mesh, stability condition, parabolic equation, explicit one-step method

**AMS subject classifications.** 65M60, 65M50, 65F15

**DOI.** 10.1137/130949531

**1. Introduction.** Adaptive meshes are commonly used for the numerical solution of partial differential equations (PDEs) to enhance computational efficiency, but there are still gaps in the mathematical understanding of the effects of the variation of element size and shape on the properties of numerical schemes for solving PDEs.

In this paper, we are concerned with the stability of explicit one-step time integration of linear finite element approximation with general nonuniform simplicial meshes for the initial-boundary value problem (IBVP)

$$(1) \quad \begin{cases} \partial_t u = \nabla \cdot (\mathbb{D} \nabla u), & \mathbf{x} \in \Omega, \quad t \in (0, T], \\ u(\mathbf{x}, t) = 0, & \mathbf{x} \in \Gamma_D, \quad t \in (0, T], \\ \mathbb{D} \nabla u(\mathbf{x}, t) \cdot \mathbf{n} = 0, & \mathbf{x} \in \Gamma_N, \quad t \in (0, T], \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), & \mathbf{x} \in \Omega, \end{cases}$$

where  $\Omega \subset \mathbb{R}^d$  ( $d \geq 1$ ) is an interval, a bounded polygonal or polyhedral domain,  $\Gamma_D \cup \Gamma_N = \partial\Omega$ ,  $\Gamma_D$  has a positive  $(d - 1)$ -volume,  $u_0$  is a given initial function, and  $\mathbb{D}$  is the diffusion matrix which is assumed to be symmetric and uniformly positive definite on  $\Omega$ . In this study, we also assume that  $\mathbb{D}$  is time independent, i.e.,  $\mathbb{D} = \mathbb{D}(\mathbf{x})$ . Problem (1) is isotropic when  $\mathbb{D}(\mathbf{x}) = \alpha(\mathbf{x})I$  for all  $\mathbf{x}$  in  $\Omega$ , where  $\alpha$  is a scalar function

---

\*Received by the editors December 17, 2013; accepted for publication (in revised form) March 29, 2016; published electronically May 26, 2016. The research of the authors was supported in part by the NSF (USA) under grant DMS-1115118, the DFG (Germany) under grant KA 3215/2-1, and the Darmstadt Graduate Schools of Excellence *Computational Engineering and Energy Science and Engineering*.

<http://www.siam.org/journals/sinum/54-3/94953.html>

<sup>†</sup>Department of Mathematics, The University of Kansas, Lawrence, KS 66045 (whuang@ku.edu).

<sup>‡</sup>Weierstrass Institute, Berlin 10117, Germany (kamenski@wias-berlin.de).

<sup>§</sup>Department of Mathematics, TU Darmstadt, Darmstadt D-64289, Germany (lang@mathematik.tu-darmstadt.de).

and  $I$  is the  $d \times d$  identity matrix. Otherwise, the problem is an anisotropic diffusion problem, which we shall consider in this work. Anisotropic diffusion arises in various areas of science and engineering, including plasma physics [7], petroleum reservoir simulation [3, 20], and image processing [17, 25].

Assume that  $u_0 \in H_D^1(\Omega) = \{v \in H^1(\Omega) : v = 0 \text{ on } \Gamma_D\}$ . Then, if  $u$  is sufficiently smooth, we have the stability estimates

$$(2) \quad \begin{cases} \|u(\cdot, t)\|_{L^2(\Omega)} \leq \|u_0\|_{L^2(\Omega)}, & t \in (0, T], \\ |||u(\cdot, t)|||_{H^1(\Omega)} \leq |||u_0|||_{H^1(\Omega)}, & t \in (0, T], \end{cases}$$

where  $|||u(\cdot, t)|||_{H^1(\Omega)} \equiv \|\mathbb{D}^{1/2}\nabla u\|_{L^2(\Omega)}$  is the energy norm of  $u(\cdot, t)$ . It is essential that a numerical scheme for (1) preserves the stability estimates. The stability of the time integration depends on the largest eigenvalue of the system related to the numerical scheme, which, in turn, depends on the underlying meshes and the coefficients of the IBVP.

For a uniform mesh and the Laplace operator, it is well known that the largest permissible time step is proportional to the square of the element diameter.

In the case of a nonuniform mesh or a variable diffusion matrix, the situation becomes more complicated. Essentially, one needs to estimate the largest eigenvalues of  $M^{-1}A$ , where  $M$  and  $A$  are the mass and stiffness matrices corresponding to the discretization of the IBVP. This can be done by estimating the extreme eigenvalues of  $M$  and  $A$ . Tight bounds on those of the mass matrix  $M$  for linear finite elements with locally quasi-uniform meshes are available in the literature and are typically proportional to the extremal mesh element volumes [4, 5, 24], whereas those for the stiffness matrix  $A$  are more difficult to obtain, and only a few results are available in the literature for the case of nonuniform meshes. For example, Fried [4] shows how to obtain these bounds for the finite element approximation of the Laplace operator for general nonuniform meshes using local element mass and stiffness matrices. A similar argument was used by Shewchuk [23] to develop a bound on the largest eigenvalue of  $M^{-1}A$  in terms of the maximum eigenvalues of local element matrices for the case of a lumped mass matrix. Graham and McLean [5] study the finite/boundary element approximation of a general differential/integral operator on locally quasi-uniform meshes in terms of patch volumes and aspect ratios. Du, Wang, and Zhu [1] obtain bounds on the extreme eigenvalues of the stiffness matrix for the Galerkin approximation of a general diffusion operator in terms of element geometry. Zhu and Du [26, 27] further develop bounds on the largest permissible time step for time dependent problems. It is worth mentioning that these existing works allow anisotropic meshes. However, the interplay between the mesh geometry and the diffusion matrix is not really taken into account, which, as we will see, is crucially important for the stability of explicit integration schemes. A notable exception is the bound obtained by Shewchuk [23], which takes the effects of the diffusion coefficients fully into account; see Remark 8 for details and Example 12 for a numerical example. Moreover, the existing analysis either employs some mesh regularity assumptions such as the local uniformity or involves parameters in final estimates that are related to mesh regularity, such as the maximum ratio of volumes of neighboring elements and/or the maximum number of elements in a patch.

The objective of this work is to provide estimates for the largest permissible time step which are accurate and tight for *any mesh* and *any diffusion matrix*. We utilize bounds recently obtained by Kamenski, Huang, and Xu [16] on the extreme eigenvalues of  $M$  and the largest eigenvalue of  $A$  for a general diffusion operator with arbitrary

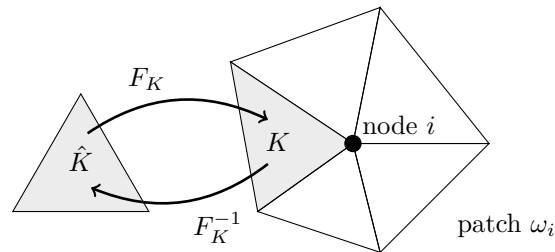


FIG. 1. Reference and mesh elements, mapping  $F_K$ ,  $i$ th node, and its patch  $\omega_i$ .

meshes. The obtained stability condition expressed in terms of matrix entries is tight within a constant factor which is independent of the mesh and the diffusion matrix. No assumption on the mesh regularity is made in the development. We show that the alignment of the mesh with the diffusion matrix plays a crucial role in the stability condition: the largest permissible time step depends only on the number of mesh elements and the mesh geometry in relation to the diffusion matrix. In particular, if the mesh is uniform in the metric specified by  $\mathbb{D}^{-1}$ , the stability condition is essentially the same as that for the Laplace operator with a uniform mesh. Although we consider only linear finite elements, the presented analysis is applicable to high order finite elements without major modifications [13].

The paper is organized as follows. We start in section 2 with the problem setting and a detailed description of mesh quality measures which are needed for the geometric interpretations of stability estimates. The main results on stability are given in section 3; both the full mass matrix and mass lumping are considered. Numerical examples to demonstrate the theoretical findings are presented in section 4, including a two-dimensional (2D) groundwater flow problem. Conclusions are drawn in section 5.

**2. Linear finite element approximation.** We consider the standard linear finite element method for the spatial discretization of IBVP (1).

We assume that a family  $\{\mathcal{T}_h\}$  of simplicial meshes is given for  $\Omega$ . While having adaptive meshes in mind, we consider the meshes to be general nonuniform ones, which may contain elements of small size and large aspect ratio. Let  $K$  be an arbitrary element of  $\mathcal{T}_h$ ,  $\hat{K}$  the *reference element*, and  $\omega_i$  the *element patch* of the  $i$ th vertex (Figure 1). Element and patch volumes are denoted by

$$|K| \quad \text{and} \quad |\omega_i| = \sum_{K \in \omega_i} |K|.$$

For each mesh element  $K \in \mathcal{T}_h$  let  $F_K$  be the invertible affine mapping from  $\hat{K}$  to  $K$  (Figure 1) and  $F'_K$  its Jacobian matrix. Note that  $F'_K$  is a constant matrix with  $\det(F'_K) = |K|$  (for simplicity, we assume that  $\hat{K}$  is equilateral with  $|\hat{K}| = 1$ ).

Let  $V^h$  be the linear finite element space associated with mesh  $\mathcal{T}_h$ . Defining  $V_D^h = V^h \cap H_D^1(\Omega) = \{v^h \in V^h : v^h = 0 \text{ on } \Gamma_D\}$ , the piecewise linear finite element solution  $u^h(t) \in V_D^h$ ,  $t \in (0, T]$ , is defined by

$$(3) \quad \int_{\Omega} v^h \partial_t u^h \, d\mathbf{x} = - \int_{\Omega} \nabla v^h \cdot \mathbb{D} \nabla u^h \, d\mathbf{x} \quad \forall v^h \in V_D^h, \quad t \in (0, T],$$

subject to the initial condition

$$(4) \quad \int_{\Omega} u^h(\mathbf{x}, 0)v^h \, d\mathbf{x} = \int_{\Omega} u_0(\mathbf{x})v^h \, d\mathbf{x} \quad \forall v^h \in V_D^h.$$

We denote the number of the elements of  $\mathcal{T}_h$  by  $N$  and the number of the interior vertices plus the vertices associated with the Neumann boundary condition by  $N_{vi}$ . If we express  $u^h$  as

$$u^h(\mathbf{x}, t) = \sum_{j=1}^{N_{vi}} u_j^h(t)\phi_j(\mathbf{x}),$$

where  $\phi_j$  is the linear basis function associated to the  $j$ th vertex ( $j = 1, \dots, N_{vi}$ ), from (3) and (4) we obtain

$$(5) \quad M\mathbf{U}_t = -A\mathbf{U}, \quad \mathbf{U}(0) = \mathbf{U}_0,$$

where  $\mathbf{U} = (u_1^h, \dots, u_{N_{vi}}^h)^T$  and  $M$  and  $A$  are the mass and the stiffness matrices,

$$(6) \quad M_{ij} = \int_{\Omega} \phi_i\phi_j \, d\mathbf{x}, \quad A_{ij} = \int_{\Omega} \nabla\phi_i \cdot \mathbb{D}\nabla\phi_j \, d\mathbf{x}, \quad i, j = 1, \dots, N_{vi}.$$

We shall investigate how the geometry of the mesh and the anisotropy of the diffusion matrix affect the stability of explicit one-step methods for integrating (5). In the following we assume that the mesh is fixed for all time steps.

**2.1. Mathematical description of nonuniform meshes; mesh quality measures.** An adaptive mesh, which is typically nonuniform, can be generated as a uniform one in the metric specified by a given metric tensor, which is always assumed to be symmetric and uniformly positive definite in  $\Omega$  [11]. On the other hand, a metric tensor can be defined for any given mesh such that all elements are uniform in the metric specified by this tensor [14]. Thus, it is natural to consider nonuniform meshes in relation to a given metric tensor. In the following, we describe several quality measures and mathematical characterizations for (nonuniform) meshes in terms of a given metric tensor  $\mathbb{M} = \mathbb{M}(\mathbf{x})$ . As we will see in section 3, the matching between the mesh metric tensor and the diffusion matrix plays a crucial role for the stability condition. In our analysis, we slightly adjust the original definitions of the mesh quality measures in [10] (see also [12, 14]).

Let

$$(7) \quad \mathbb{M}_K = \frac{1}{|K|} \int_K \mathbb{M} \, d\mathbf{x}, \quad |K|_{\mathbb{M}} = |K|\det(\mathbb{M}_K)^{\frac{1}{2}}, \quad |\Omega|_{\mathbb{M},h} = \sum_{K \in \mathcal{T}_h} |K|_{\mathbb{M}}.$$

Note that  $\mathbb{M}_K$  is the average of  $\mathbb{M}$  over the element  $K$  and  $|K|_{\mathbb{M}}$  and  $|\Omega|_{\mathbb{M},h}$  are approximate volumes of  $K$  and  $\Omega$  in the metric  $\mathbb{M}$ , viz.,

$$|K|_{\mathbb{M}} \approx \int_K \det(\mathbb{M}(\mathbf{x}))^{\frac{1}{2}} \, d\mathbf{x} \quad \text{and} \quad |\Omega|_{\mathbb{M},h} \approx \sum_{K \in \mathcal{T}_h} \int_K \det(\mathbb{M}(\mathbf{x}))^{\frac{1}{2}} \, d\mathbf{x} = |\Omega|_{\mathbb{M}}.$$

Hereafter, without confusion we will call  $|K|_{\mathbb{M}}$  and  $|\Omega|_{\mathbb{M},h}$  the volumes of  $K$  and  $\Omega$  in the metric  $\mathbb{M}$ , respectively. We also define the *average diameter* of element  $K$  and the *global average element diameter* with respect to  $\mathbb{M}$  as

$$h_{K,\mathbb{M}} = |K|_{\mathbb{M}}^{\frac{1}{d}} \quad \text{and} \quad h_{\mathbb{M}} = \left( \frac{1}{N} |\Omega|_{\mathbb{M},h} \right)^{\frac{1}{d}}.$$

The diameter  $h_K$  of  $K$  is defined as the length of the longest edge of  $K$ .

With this notation established, we now are ready to describe the mesh quality measures. The first one, the *equidistribution* quality measure, is defined as the ratio of the average element volume to the volume of  $K$ , both measured in the metric specified by  $\mathbb{M}_K$ ,

$$(8) \quad Q_{\text{eq},\mathbb{M}}(K) = \frac{\frac{1}{N}|\Omega|_{\mathbb{M},h}}{|K|_{\mathbb{M}}} = \left( \frac{h_{\mathbb{M}}}{h_{K,\mathbb{M}}} \right)^d.$$

It satisfies

$$(9) \quad 0 < Q_{\text{eq},\mathbb{M}}(K) < \infty, \quad \frac{1}{N} \sum_{K \in \mathcal{T}_h} \frac{1}{Q_{\text{eq},\mathbb{M}}(K)} = 1, \quad \max_{K \in \mathcal{T}_h} Q_{\text{eq},\mathbb{M}}(K) \geq 1.$$

The second one, the *alignment* quality measure, is local (elementwise) and measures how closely the principal directions of the circumscribed ellipsoid of  $K$  are aligned with the eigenvectors of  $\mathbb{M}_K$ , and the semilengths of the principal axes are inversely proportional to the square root of the eigenvalues of  $\mathbb{M}_K$ . It is defined as

$$(10) \quad Q_{\text{ali},\mathbb{M}}(K) = \frac{\left\| (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right\|_2}{\det \left( (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right)^{\frac{1}{d}}} = h_{K,\mathbb{M}}^2 \left\| (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right\|_2.$$

Since  $\left\| (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right\|_2 \geq \det \left( (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right)^{\frac{1}{d}}$ ,  $Q_{\text{ali},\mathbb{M}}$  always satisfies

$$1 \leq Q_{\text{ali},\mathbb{M}}(K) < \infty$$

with  $Q_{\text{ali},\mathbb{M}}(K) = 1$  if and only if  $K$  is equilateral with respect to  $\mathbb{M}_K$ . The alignment quality measure can be seen as an alternative to the aspect ratio of  $K$  in the metric specified by  $\mathbb{M}_K$ , and it satisfies

$$(11) \quad Q_{\text{ali},\mathbb{M}}(K) \leq \hat{h}^2 \cdot \left( \frac{h_{K,\mathbb{M}}}{\rho_{K,\mathbb{M}}} \right)^2,$$

where  $\hat{h}$  is the length of the longest edge of  $\hat{K}$  and  $\rho_{K,\mathbb{M}}$  is the diameter of the largest sphere inscribed in the element  $K$  viewed in the metric  $M_K$ . To show this, we consider two points  $\mathbf{x}_1, \mathbf{x}_2 \in K$  and the corresponding points  $\boldsymbol{\xi}_1 = F_K^{-1}(\mathbf{x}_1)$  and  $\boldsymbol{\xi}_2 = F_K^{-1}(\mathbf{x}_2)$  in  $\hat{K}$ . The distance between  $\mathbf{x}_1$  and  $\mathbf{x}_2$  in the metric  $\mathbb{M}_K$  is

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}_2\|_{\mathbb{M}_K}^2 &= (\mathbf{x}_1 - \mathbf{x}_2)^T \mathbb{M}_K (\mathbf{x}_1 - \mathbf{x}_2) = (\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)^T (F'_K)^T \mathbb{M}_K F'_K (\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2) \\ &= \|\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2\|_2^2 \cdot \frac{(\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)^T (F'_K)^T \mathbb{M}_K F'_K (\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)}{\|\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2\|_2^2} \\ &\leq \hat{h}^2 \cdot \frac{(\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)^T (F'_K)^T \mathbb{M}_K F'_K (\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)}{\|\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2\|_2^2}. \end{aligned}$$

If we take the minimum over all pairs of opposing points on the largest sphere inscribed in the element  $K$  viewed in the metric  $M_K$ , then

$$\rho_{K,\mathbb{M}}^2 \leq \hat{h}^2 \lambda_{\min} \left( (F'_K)^T \mathbb{M}_K F'_K \right).$$

Hence,

$$(12) \quad \left\| (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right\|_2 = \frac{1}{\lambda_{\min}((F'_K)^T \mathbb{M}_K F'_K)} \leq \left( \frac{\hat{h}}{\rho_{K, \mathbb{M}}} \right)^2,$$

which, together with (10), gives (11).

The *element quality* measure is defined as a combination of  $Q_{\text{ali}, \mathbb{M}}$  and  $Q_{\text{eq}, \mathbb{M}}$ ,

$$(13) \quad Q_{\mathbb{M}}(K) = Q_{\text{ali}, \mathbb{M}}(K) \cdot (Q_{\text{eq}, \mathbb{M}}(K))^{\frac{2}{d}} = h_{\mathbb{M}}^2 \left\| (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} \right\|_2.$$

It measures how far  $K$  is from being equilateral with a constant volume when viewed in the metric specified by  $\mathbb{M}$ . By definition and from (12) it follows that

$$(14) \quad 0 < Q_{\mathbb{M}}(K) \leq \hat{h}^2 \left( \frac{h_{\mathbb{M}}}{\rho_{K, \mathbb{M}}} \right)^2 < \infty.$$

When a mesh is uniform with respect to  $\mathbb{M}$  (we will refer to it as an  $\mathbb{M}$ -uniform mesh), it satisfies

$$(15) \quad Q_{\text{ali}, \mathbb{M}}(K) = 1 \quad \text{and} \quad Q_{\text{eq}, \mathbb{M}}(K) = 1 \quad \forall K \in \mathcal{T}_h,$$

which is equivalent to

$$(16) \quad Q_{\mathbb{M}}(K) = 1 \quad \forall K \in \mathcal{T}_h.$$

Indeed, (16) follows directly from (15). On the other hand, since  $Q_{\text{ali}, \mathbb{M}} \geq 1$ , (16) implies  $Q_{\text{eq}, \mathbb{M}}(K) \leq 1$  for all  $K$ . Due to the property (9), the latter is only possible if  $Q_{\text{eq}, \mathbb{M}}(K) = 1$  for all  $K$ , which, in turn, implies  $Q_{\text{ali}, \mathbb{M}}(K) = 1$  for all  $K$ .

It is worth mentioning that an  $\mathbb{M}$ -uniform mesh satisfies

$$(17) \quad (F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T} = h_{\mathbb{M}}^{-2} I \quad \forall K \in \mathcal{T}_h,$$

since (15) implies that all eigenvalues of  $(F'_K)^{-1} \mathbb{M}_K^{-1} (F'_K)^{-T}$  are equal to  $h_{\mathbb{M}}$ . On the other hand, when a mesh is far from being  $\mathbb{M}$ -uniform, then

$$Q_{\text{ali}, \mathbb{M}}(K) \gg 1 \quad \text{or} \quad \max_K Q_{\text{eq}, \mathbb{M}}(K) \gg 1,$$

and therefore

$$\max_K Q_{\mathbb{M}}(K) \gg 1.$$

**2.2. Preliminary results.** In this subsection we present a few properties of the mass matrix  $M$  and the stiffness matrix  $A$  of linear finite elements, which will be used repeatedly in our analysis. Throughout the paper the less-than-or-equal-to sign between matrix terms means that the difference between the right-hand side and left-hand side terms is positive semidefinite.

LEMMA 2.1 (see [16, sect. 3]). *The linear finite element mass matrix  $M$  and its diagonal part  $M_D$  satisfy*

$$(18) \quad \frac{1}{2} M_D \leq M \leq \frac{d+2}{2} M_D \quad \text{and} \quad M_{ii} = \frac{2|\omega_i|}{(d+1)(d+2)}, \quad i = 1, \dots, N_{vi}.$$

LEMMA 2.2. Let  $M_{lump}$  be the lumped linear finite element mass matrix defined through

$$M_{ii,lump} = \int_{\Omega} \phi_i(\mathbf{x}) \cdot \sum_{j=1}^{N_{vi}} \phi_j(\mathbf{x}) \, d\mathbf{x}, \quad i = 1, \dots, N_{vi}.$$

Then

$$(19) \quad \frac{2|\omega_i|}{(d+1)(d+2)} \leq M_{ii,lump} \leq \frac{|\omega_i|}{d+1}.$$

*Proof.* Since

$$\phi_i(\mathbf{x}) \leq \sum_{j=1}^{N_{vi}} \phi_j(\mathbf{x}) \leq 1,$$

we have

$$M_{ii,lump} \geq \int_{\Omega} \phi_i^2(\mathbf{x}) \, d\mathbf{x} = \sum_{K \in \omega_i} \int_K \phi_i^2(\mathbf{x}) \, d\mathbf{x} = \sum_{K \in \omega_i} \frac{2|K|}{(d+1)(d+2)} = \frac{2|\omega_i|}{(d+1)(d+2)}$$

and

$$M_{ii,lump} \leq \int_{\Omega} \phi_i(\mathbf{x}) \, d\mathbf{x} = \sum_{K \in \omega_i} \int_K \phi_i(\mathbf{x}) \, d\mathbf{x} = \sum_{K \in \omega_i} \frac{|K|}{d+1} = \frac{|\omega_i|}{d+1}. \quad \square$$

LEMMA 2.3. The linear finite element mass matrix  $M$  and the lumped mass matrix  $M_{lump}$  satisfy

$$\frac{1}{d+2} M_{lump} \leq M \leq \frac{d+2}{2} M_{lump}.$$

*Proof.* Since  $M_D \leq M_{lump}$ , we get the upper bound directly from (18). Combining the lower bound in (18) with the upper bound in (19) gives

$$\frac{1}{d+2} M_{lump} \leq \frac{1}{(d+2)(d+1)} \text{diag}(|\omega_1|, \dots, |\omega_{N_{vi}}|) = \frac{1}{2} M_D \leq M. \quad \square$$

LEMMA 2.4 (see [16, sect. 4]). The linear finite element stiffness matrix  $A$  and its diagonal part  $A_D$  satisfy

$$(20) \quad A \leq (d+1)A_D.$$

LEMMA 2.5. Let  $\mathbb{D}_K$  be the average of the diffusion matrix  $\mathbb{D}$  over  $K$ ,

$$\mathbb{D}_K = \frac{1}{|K|} \int_K \mathbb{D}(\mathbf{x}) \, d\mathbf{x}.$$

Then the diagonal entries of the linear finite element stiffness matrix  $A$  are bounded by

$$(21) \quad C_{\hat{\nabla}} \sum_{K \in \omega_i} |K| \cdot \lambda_{\min}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}) \leq A_{ii} \leq C_{\hat{\nabla}} \sum_{K \in \omega_i} |K| \cdot \lambda_{\max}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}),$$

where  $C_{\hat{\nabla}} = \frac{d}{d+1} \left( \frac{\sqrt{d+1}}{d} \right)^{\frac{2}{d}}$ .

*Proof.* From (6) we have

$$A_{ii} = \int_{\Omega} \nabla \phi_i^T \mathbb{D} \nabla \phi_i \, d\mathbf{x} = \sum_{K \in \omega_i} \int_K \nabla \phi_i^T \mathbb{D} \nabla \phi_i \, d\mathbf{x} = \sum_{K \in \omega_i} |K| \nabla \phi_i^T \mathbb{D}_K \nabla \phi_i.$$

Denote the gradient operator in  $\hat{K}$  by  $\hat{\nabla} = \frac{\partial}{\partial \hat{\xi}}$ . By the chain rule,  $\nabla = (F'_K)^{-T} \hat{\nabla}$  and

$$(22) \quad \begin{aligned} A_{ii} &= \sum_{K \in \omega_i} |K| \hat{\nabla} \hat{\phi}_i^T (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \hat{\nabla} \hat{\phi}_i \\ &\leq \sum_{K \in \omega_i} |K| \hat{\nabla} \hat{\phi}_i^T \hat{\nabla} \hat{\phi}_i \lambda_{\max}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}). \end{aligned}$$

Recall that  $\hat{K}$  is taken to be equilateral. Thus,  $\hat{\nabla} \hat{\phi}_i^T \hat{\nabla} \hat{\phi}_i = C_{\hat{\nabla}}$  for all  $i = 1, \dots, d + 1$ . Consequently,

$$A_{ii} \leq C_{\hat{\nabla}} \sum_{K \in \omega_i} |K| \lambda_{\max}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}).$$

Similarly, we can obtain the left inequality of (21). □

*Remark 1.* From (13), with  $\mathbb{M}$  being replaced by  $\mathbb{D}^{-1}$ , the bound (21) on  $A_{ii}$  can be expressed in terms of the element quality measure  $Q_{\mathbb{D}^{-1}}(K)$  as

$$(23) \quad A_{ii} \leq C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_i} |K| Q_{\mathbb{D}^{-1}}(K).$$

*Remark 2* ( $\mathbb{D}^{-1}$ -nonobtuse meshes). Note that Lemma 2.4 is very general and valid for any given mesh. It implies that

$$(24) \quad \lambda_{\max}(A) \leq (d + 1) \max_i A_{ii}.$$

This bound can be sharpened for some special types of meshes. For example, if a mesh has no obtuse angles with respect to  $\mathbb{D}^{-1}$ , then  $A$  is an  $M$ -matrix (its off-diagonal entries are nonpositive) and  $\sum_j A_{ij} \geq 0$  for all  $i$  (e.g., see the proof of Theorem 2.1 of [18]). From the Gershgorin circle theorem we have

$$\lambda_{\max}(A) \leq \max_i \left( A_{ii} + \sum_{j \neq i} |A_{ij}| \right) = \max_i \left( A_{ii} - \sum_{j \neq i} A_{ij} \right) = \max_i \left( 2A_{ii} - \sum_j A_{ij} \right),$$

and thus

$$(25) \quad \lambda_{\max}(A) \leq 2 \max_i A_{ii}.$$

If, further, the mesh is  $\mathbb{D}^{-1}$ -uniform, then from (16) and (23) we have

$$(26) \quad \lambda_{\max}(A) \leq 2 \max_i A_{ii} \leq 2C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \max_i \sum_{K \in \omega_i} |K| Q_{\mathbb{D}^{-1}}(K) = 2C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \max_i |\omega_i|.$$

**3. Explicit time stepping and the stability condition.** In this section we study stability conditions for explicit one-step methods applied to the finite element system (5) and obtain estimates for the maximum time step.

Suppose we are given a constant time step  $\tau$ . Then an explicit one-step integration scheme with  $s$  stages of order  $p$  computes approximations  $\mathbf{U}_n \approx \mathbf{U}(n\tau)$  from

$$(27) \quad \mathbf{U}_n = R(-\tau M^{-1}A)\mathbf{U}_{n-1},$$

where the stability function  $R(z)$  is a polynomial in  $z$  and satisfies

$$(28) \quad R(z) = 1 + z + \cdots + \frac{z^p}{p!} + \sum_{i=p+1}^s \alpha_i z^i = e^z + \mathcal{O}(z^{p+1}).$$

Classical explicit one-step methods have severe step size restrictions when solving stiff problems such as (5) for  $N_{vi} \gg 1$ . An interesting alternative is stabilized explicit Runge–Kutta (RK) methods, which have an extended stability domain along the negative real axis and therefore allow for larger time steps than classical explicit one-step methods. Parameters  $\alpha_{p+1}, \dots, \alpha_s \in \mathbb{R}$  in (28) are chosen such that  $|R(z)| \leq 1$  for  $z \in [-r_s, 0]$  and  $r_s > 0$  is as large as possible. Explicit methods have low memory demand and can be considered as a good alternative to implicit methods when the solution of algebraic equations arising from the latter is difficult and/or costly. Impressive examples and comparison results with VODEPK (a stiff ODE solver with Krylov iterations) are documented in [15]. Commonly used explicit methods include DUMKA, Runge–Kutta–Chebyshev (RKC), and the orthogonal RKC (ROCK) methods. A common practical choice is  $p = 2$ , but there exist also DUMKA and ROCK methods of higher order [8].

In the following we first study stability estimates for the approximate solutions  $\mathbf{U}_n$  obtained from (27), assuming that  $M$  is a full mass matrix. However, the decomposition of a consistent mass matrix as a part of an explicit time integration method is in general not affordable, since an implicit scheme with a much larger step may be performed at the same cost. Hence, we mainly discuss consequences of lumping the mass matrix as a routine procedure for (linear) finite elements. Although appropriate mass lumping does not affect the overall accuracy, it is well known that lumping the consistent mass induces dispersion errors that can affect the quality of the numerical solution. More generally, we consider symmetric positive definite, surrogate matrices  $\tilde{M}$  that satisfy

$$(29) \quad c_1 \tilde{M} \leq M \leq c_2 \tilde{M}$$

and have nearly the same complexity as the diagonal lumped mass matrix  $M_{lump}$ . Correction techniques for the dispersive effects of mass lumping and several efficient choices for  $\tilde{M}$  can be found in [6]. Note that, due to Lemma 2.3, we have  $c_1 = 1/(d+2)$  and  $c_2 = (d+2)/2$  for the special case  $\tilde{M} = M_{lump}$ .

**3.1. Stability of explicit one-step integration schemes.** The investigation of the stability is based on the following main observation: if  $B$  is a normal matrix and  $R$  is a rational function, then

$$(30) \quad \|R(B)\|_2 = \max_i |R(\lambda_i(B))|.$$

This fundamental relation is a direct consequence of the existence of a factorization  $B = Q \operatorname{diag}(\lambda_1(B), \dots, \lambda_{N_{vi}}(B)) Q^T$  with a unitary matrix  $Q$ .

Using the fact that  $M^{-\frac{1}{2}} A M^{-\frac{1}{2}}$  and  $A^{\frac{1}{2}} M^{-1} A^{\frac{1}{2}}$  are normal matrices, we can prove the stability of the linear finite element approximation computed with an explicit one-step method.

**THEOREM 3.** *For a given explicit one-step method with the polynomial stability function  $R$ , the linear finite element approximation  $u_n^h$  satisfies*

$$\|u_n^h\|_{L^2(\Omega)} \leq \|u_0^h\|_{L^2(\Omega)} \quad \text{and} \quad \|u_n^h\|_{H^1(\Omega)} \leq \|u_0^h\|_{H^1(\Omega)}$$

if the time step  $\tau$  is chosen such that

$$\max_i |R(-\tau\lambda_i(M^{-1}A))| \leq 1.$$

*Proof.* Since  $R$  is a polynomial function, we have

$$R(-\tau M^{-1}A) = M^{-\frac{1}{2}}R(-\tau M^{-\frac{1}{2}}AM^{-\frac{1}{2}})M^{\frac{1}{2}} = A^{-\frac{1}{2}}R(-\tau A^{\frac{1}{2}}M^{-1}A^{\frac{1}{2}})A^{\frac{1}{2}}.$$

From this, it is easy to see that (27) can be written as

$$(31) \quad M^{\frac{1}{2}}\mathbf{U}_n = R(-\tau M^{-\frac{1}{2}}AM^{-\frac{1}{2}})M^{\frac{1}{2}}\mathbf{U}_{n-1},$$

$$(32) \quad A^{\frac{1}{2}}\mathbf{U}_n = R(-\tau A^{\frac{1}{2}}M^{-1}A^{\frac{1}{2}})A^{\frac{1}{2}}\mathbf{U}_{n-1}.$$

Since  $M$  and  $A$  are symmetric and positive definite,  $M^{-\frac{1}{2}}AM^{-\frac{1}{2}}$  and  $A^{\frac{1}{2}}M^{-1}A^{\frac{1}{2}}$  are symmetric and therefore normal. From (30), our assumption on the time step, and the fact that  $M^{-1}A$ ,  $M^{-\frac{1}{2}}AM^{-\frac{1}{2}}$ , and  $A^{\frac{1}{2}}M^{-1}A^{\frac{1}{2}}$  are similar to each other, we get

$$\|R(-\tau M^{-\frac{1}{2}}AM^{-\frac{1}{2}})\|_2 = \|R(-\tau A^{\frac{1}{2}}M^{-1}A^{\frac{1}{2}})\|_2 = \max_i |R(-\tau\lambda_i(M^{-1}A))| \leq 1.$$

Thus, (31) and (32) imply

$$\begin{aligned} \|u_n^h\|_{L^2(\Omega)} &= \|M^{\frac{1}{2}}\mathbf{U}_n\|_2 \leq \|M^{\frac{1}{2}}\mathbf{U}_{n-1}\|_2 = \|u_{n-1}^h\|_{L^2(\Omega)}, \\ \|u_n^h\|_{H^1(\Omega)} &= \|A^{\frac{1}{2}}\mathbf{U}_n\|_2 \leq \|A^{\frac{1}{2}}\mathbf{U}_{n-1}\|_2 = \|u_{n-1}^h\|_{H^1(\Omega)}. \end{aligned}$$

Successive application of these inequalities yields the assertion. □

We next consider the case where the linear finite element mass matrix  $M$  is replaced by a symmetric positive definite, surrogate matrix  $\tilde{M}$  of lower complexity. That means that from now on we compute approximations  $\mathbf{U}_n \approx \mathbf{U}(n\tau)$  from

$$(33) \quad \mathbf{U}_n = R(-\tau \tilde{M}^{-1}A)\mathbf{U}_{n-1}.$$

**THEOREM 4.** *For a given explicit one-step method with the polynomial stability function  $R$  and a symmetric positive definite, surrogate matrix  $\tilde{M}$  that satisfies  $c_1\tilde{M} \leq M \leq c_2\tilde{M}$  for some positive constants  $c_1$  and  $c_2$ , the linear finite element approximation  $u_n^h$  satisfies*

$$\|u_n^h\|_{L^2(\Omega)} \leq \sqrt{\frac{c_2}{c_1}} \|u_0^h\|_{L^2(\Omega)} \quad \text{and} \quad \|u_n^h\|_{H^1(\Omega)} \leq \|u_0^h\|_{H^1(\Omega)}$$

if the time step  $\tau$  is chosen such that

$$\max_i |R(-\tau\lambda_i(\tilde{M}^{-1}A))| \leq 1.$$

TABLE 1  
 $C_*$  in Theorem 5.

	General meshes	Nonobtuse w.r.t. $\mathbb{D}^{-1}$
$\tilde{M} = M$	$2(d+1)$	4
$\tilde{M} = M_{lump}$	$d+1$	2

*Proof.* Replacing  $M$  by  $\tilde{M}$  in the proof of Theorem 3 does not change the arguments and gives

$$\begin{aligned} \|u_n^h\|_{H^1(\Omega)} &= \|A^{\frac{1}{2}} \mathbf{U}_n\|_2 \leq \|A^{\frac{1}{2}} \mathbf{U}_{n-1}\|_2 = \|u_{n-1}^h\|_{H^1(\Omega)}, \\ \|\tilde{M}^{\frac{1}{2}} \mathbf{U}_n\|_2 &\leq \|\tilde{M}^{\frac{1}{2}} \mathbf{U}_{n-1}\|_2. \end{aligned}$$

From the first inequality, stability in the energy norm follows. To derive stability in the  $L^2$ -norm, we make use of the assumption on  $\tilde{M}$ :

$$\begin{aligned} \|u_n^h\|_{L^2(\Omega)}^2 &= (\mathbf{U}_n)^T M \mathbf{U}_n \leq c_2 (\mathbf{U}_n)^T \tilde{M} \mathbf{U}_n = c_2 \|\tilde{M}^{\frac{1}{2}} \mathbf{U}_n\|_2^2 \leq c_2 \|\tilde{M}^{\frac{1}{2}} \mathbf{U}_{n-1}\|_2^2 \\ &\leq \dots \leq c_2 \|\tilde{M}^{\frac{1}{2}} \mathbf{U}_0\|_2^2 = c_2 (\mathbf{U}_0)^T \tilde{M} \mathbf{U}_0 \leq \frac{c_2}{c_1} (\mathbf{U}_0)^T M \mathbf{U}_0 = \frac{c_2}{c_1} \|u_0^h\|_{L^2(\Omega)}^2, \end{aligned}$$

which gives the desired result.  $\square$

In the special case  $\tilde{M} = M_{lump}$  we have the following result.

COROLLARY 4.1. *Under the assumptions of Theorem 4 and  $\tilde{M} = M_{lump}$ , we have*

$$\|u_n^h\|_{L^2(\Omega)} \leq \frac{d+2}{\sqrt{2}} \|u_0^h\|_{L^2(\Omega)} \quad \text{and} \quad \|u_n^h\|_{H^1(\Omega)} \leq \|u_0^h\|_{H^1(\Omega)}.$$

**3.2. Estimates on the largest eigenvalue of  $\tilde{M}^{-1}A$ .** The above results show that the contractivity of any given explicit one-step method is guaranteed if all eigenvalues of  $-\tau \tilde{M}^{-1}A$  are in the corresponding stability domain  $\mathcal{S} = \{z \in \mathbb{C} : |R(z)| \leq 1\}$ . As a consequence, the key to the stability analysis of a given scheme is the estimation of the eigenvalues of  $\tilde{M}^{-1}A$ . The following theorem provides such an estimate for two choices of  $\tilde{M}$ :  $\tilde{M} = M$  and  $\tilde{M} = M_{lump}$ . It turns out that in these cases the largest eigenvalue of  $\tilde{M}^{-1}A$  is equivalent to the largest ratio between the corresponding diagonal entries of  $A$  and  $\tilde{M}$ .

THEOREM 5. *The eigenvalues of  $\tilde{M}^{-1}A$  with  $\tilde{M}$  being either  $M$  or  $M_{lump}$  are real and positive. Moreover, the largest eigenvalue is bounded by*

$$(34) \quad \max_i \frac{A_{ii}}{\tilde{M}_{ii}} \leq \lambda_{\max}(\tilde{M}^{-1}A) \leq C_* \max_i \frac{A_{ii}}{\tilde{M}_{ii}},$$

where  $C_*$  is given in Table 1.

*Proof.* Since  $\tilde{M}$  and  $A$  are symmetric and positive definite and  $\tilde{M}^{-1}A$  is similar to the symmetric matrix  $\tilde{M}^{-\frac{1}{2}} A \tilde{M}^{-\frac{1}{2}}$ , the eigenvalues of  $\tilde{M}^{-1}A$  are real and positive.

Using the canonical basis vectors  $\mathbf{e}_i$  gives

$$\lambda_{\max}(\tilde{M}^{-1}A) = \max_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{v}^T \tilde{M} \mathbf{v}} \geq \max_i \frac{\mathbf{e}_i^T A \mathbf{e}_i}{\mathbf{e}_i^T \tilde{M} \mathbf{e}_i} = \max_i \frac{A_{ii}}{\tilde{M}_{ii}}.$$

Let us first have a look at the case  $\tilde{M} = M$ . Lemmas 2.1 and 2.4 yield

$$(35) \quad \lambda_{\max}(M^{-1}A) = \max_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{v}^T M \mathbf{v}} \leq \max_{\mathbf{v} \neq 0} \frac{(d+1)\mathbf{v}^T A_D \mathbf{v}}{\frac{1}{2}\mathbf{v}^T M_D \mathbf{v}} = 2(d+1) \max_i \frac{A_{ii}}{M_{ii}}.$$

For the special case of meshes with nonobtuse angles with respect to  $\mathbb{D}^{-1}$ , the above bound can be sharpened by replacing the factor  $d+1$  in (35) with 2 (see Remark 2). If  $\tilde{M} = M_{lump}$ , then the factor 1/2 in the denominator of (35) can be replaced by 1 since  $M_{lump}$  is already diagonal.  $\square$

*Example 6* (stabilized RK methods). The stability region of a stabilized RK method of order  $p = 1$  with  $s$  stages extends along the negative real axis of the complex plane, including the interval  $[-2s^2, 0]$  [8, p. 31f]. Thus, the method is stable if all eigenvalues of  $-\tau \tilde{M}^{-1}A$  are between  $-2s^2$  and 0. This leads to the stability condition

$$(36) \quad \tau \leq \frac{2s^2}{\lambda_{\max}(\tilde{M}^{-1}A)}.$$

Using Theorem 5 and noticing that  $(\max_i \frac{A_{ii}}{M_{ii}})^{-1} = \min_i \frac{\tilde{M}_{ii}}{A_{ii}}$ , we obtain a bound for the largest permissible time step  $\tau_{\max}$  as

$$(37) \quad \frac{2s^2}{C_*} \min_i \frac{\tilde{M}_{ii}}{A_{ii}} \leq \tau_{\max} \leq 2s^2 \min_i \frac{\tilde{M}_{ii}}{A_{ii}}.$$

Clearly, if

$$\tau > 2s^2 \min_i \frac{\tilde{M}_{ii}}{A_{ii}},$$

we have  $|R(-\tau \lambda_{\max}(\tilde{M}^{-1}A))| > 1$  and the scheme becomes unstable. In order to guarantee stability, the step size has to be chosen such that

$$\tau \leq \frac{2s^2}{C_*} \min_i \frac{\tilde{M}_{ii}}{A_{ii}}.$$

Note that here  $\tilde{M} = M$  or  $\tilde{M} = M_{lump}$ .

In practice, a few steps of a nonlinear power method are often sufficient to estimate the spectral radius automatically, especially if the eigenvalues are close to the negative real axis. However, the power method can degenerate in many ways, so precaution has to be taken and theoretical bounds can be helpful. Such bounds are also necessary for gaining insight into the effects of mesh geometry on the stability of explicit integration schemes and the maximum allowed step size. The estimate in Theorem 5 is easily computable, but it does not explain how the mesh geometry affects the time step. To reveal these effects, we provide several geometric formulations of the estimate in the following. First, substituting (18) and (23) for  $\tilde{M}_{ii}$  and  $A_{ii}$  in Theorem 5 gives the following corollary.

COROLLARY 6.1. *The largest eigenvalue of  $\tilde{M}^{-1}A$  is bounded by*

$$(38) \quad \lambda_{\max}(\tilde{M}^{-1}A) \leq C_* C_{\#} \max_i \sum_{K \in \omega_i} \frac{|K|}{|\omega_i|} \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2$$

$$(39) \quad = C_* C_{\#} h_{\mathbb{D}^{-1}}^{-2} \max_i \sum_{K \in \omega_i} \frac{|K|}{|\omega_i|} Q_{\mathbb{D}^{-1}}(K),$$

where  $C_{\#} = \frac{1}{2}C_{\nabla}(d+1)(d+2)$ ,  $C_{\nabla}$  and  $C_*$  are as given in Lemma 2.5 and Table 1, and the element quality  $Q_{\mathbb{D}^{-1}}(K)$  is as defined in (13) (with  $\mathbb{M}$  being replaced by  $\mathbb{D}^{-1}$ ).

The factor  $h_{\mathbb{D}^{-1}}^{-2}$  in (38) corresponds to  $h^2$  in the classic stability condition  $\tau \sim h^2$  for uniform meshes with the Laplace operator. Since

$$h_{\mathbb{D}^{-1}} = (|\Omega|_{\mathbb{D}^{-1},h}/N)^{\frac{1}{d}} \rightarrow (|\Omega|_{\mathbb{D}^{-1}}/N)^{\frac{1}{d}}$$

as the mesh is being refined,  $h_{\mathbb{D}^{-1}}$  can be considered independent of the mesh geometry, and therefore it essentially depends only on  $N$ ,  $\mathbb{D}$ , and  $\Omega$ .

The effect of the mesh geometry is reflected mainly through the patch-average of  $\|(F'_K)^{-1}\mathbb{D}_K(F'_K)^{-T}\|_2$  or, alternatively, the element quality measure  $Q_{\mathbb{D}^{-1}}(K)$ . Recall from (14) that  $Q_{\mathbb{D}^{-1}}(K)$  can be seen as a ratio of the average element size to the diameter of the largest sphere inscribed in  $K$ , both measured in the metric  $\mathbb{D}_K^{-1}$ .

Hence, we can conclude that the largest possible time step depends on *the number of mesh elements and the correspondence of the geometry of the mesh elements to  $\mathbb{D}^{-1}$* . In other words, it is not the mesh geometry itself but *the mesh geometry in relation to the diffusion matrix that matters for the stability of explicit schemes*.

We now study the situation with an  $\mathbb{M}$ -uniform mesh for a general metric tensor  $\mathbb{M}$ . Recall that such a mesh satisfies (17), which can be rewritten as

$$(F'_K)^{-T}(F'_K)^{-1} = h_{\mathbb{M}}^{-2}\mathbb{M}_K \quad \forall K \in \mathcal{T}_h.$$

Then,

$$Q_{\mathbb{D}^{-1}}(K) = h_{\mathbb{D}^{-1}}^2 \left\| (F'_K)^{-1}\mathbb{D}_K(F'_K)^{-T} \right\|_2 = \left( \frac{h_{\mathbb{D}^{-1}}}{h_{\mathbb{M}}} \right)^2 \|\mathbb{M}_K\mathbb{D}_K\|_2.$$

Inserting this into (38), we get

$$(40) \quad \lambda_{\max}(\tilde{M}^{-1}A) \leq C_*C_{\#}h_{\mathbb{M}}^{-2} \max_i \sum_{K \in \omega_i} \frac{|K|}{|\omega_i|} \cdot \|\mathbb{M}_K\mathbb{D}_K\|_2.$$

Once again, this shows that the largest eigenvalue of  $\tilde{M}^{-1}A$  and, consequently, the largest permissible time step depend on the number of elements and the matching between the mesh (essentially determined by  $\mathbb{M}$ ) and the diffusion matrix. If mesh adaptation and the major diffusion directions match, the largest permissible time step depends mainly on the number of mesh elements. A mismatch between (anisotropic) mesh adaptation and the diffusion directions can lead to a drastic reduction of the time step (see Example 10 in section 4). In particular, it implies that one gets both accuracy and stability with the same grid if the solution anisotropy is in correspondence with the diffusion and one would have to trade off accuracy for stability if the demands of accuracy and stability contradict each other (see also remarks by Shewchuk [23, sect. 4.3]). To some extent, the demands of accuracy and stability can be combined using a metric tensor in the form

$$\mathbb{M}_K = \theta_K\mathbb{D}_K^{-1} \quad \forall K \in \mathcal{T}_h,$$

where  $\theta_K$  is a scalar function based on some (isotropic) error estimate; a similar idea was used in [18] to combine mesh adaptation with satisfaction of the maximum principle. This will not provide full mesh adaptation but will provide at least some degree of it.

*Remark 7* (special cases). For a uniform mesh ( $\mathbb{M} = I$ ), we have

$$\lambda_{\max}(\tilde{M}^{-1}A) \leq C_* C_{\#} h^{-2} \max_i \sum_{K \in \omega_i} \frac{|K|}{|\omega_i|} \cdot \|\mathbb{D}_K\|_2 \approx C_* C_{\#} h^{-2} \max_i \|\mathbb{D}_{\omega_i}\|_2,$$

where  $\mathbb{D}_{\omega_i}$  denotes the average of  $\mathbb{D}$  over a patch  $\omega_i$ .

In case of coefficient-adaptive ( $\mathbb{D}^{-1}$ -uniform) meshes ( $\mathbb{M} = \mathbb{D}^{-1}$ ), mesh adaptation and diffusion directions match exactly, and (17) and (22) yield  $A_{ii} = C_{\nabla} |\omega_i| / h_{\mathbb{D}^{-1}}^2$ . Thus, using (18) and Theorem 5 gives

$$\lambda_{\max}(\tilde{M}^{-1}A) \sim h_{\mathbb{D}^{-1}}^{-2} \sim N^{\frac{2}{d}}.$$

*Remark 8* (comparison to results available in the literature). For the full mass matrix, Zhu and Du [27, Theorem 3.1] developed an estimate in terms of the element geometry and the eigenvalues of the diffusion matrix which is valid for  $d \geq 2$  and  $P_k$  finite elements. For the linear finite elements it becomes

$$(41) \quad \frac{\max_K \lambda_{\min}(\mathbb{D}_K) \mathbf{Z}_K}{d(1 + c_1 p_{\max}(d + 2))} \leq \lambda_{\max}(M^{-1}A) \leq (d + 2) \max_K \lambda_{\max}(\mathbb{D}_K) \mathbf{Z}_K,$$

$$\mathbf{Z}_K = \frac{d + 1}{d^2} \sum_{i_K} \frac{|V_{i_K}|^2}{|K|^2},$$

where  $|V_{i_K}|$  is the volume of a  $(d - 1)$ -dimensional face opposing the  $i_K$ th vertex of  $K$ ,  $p_{\max}$  is the maximum number of elements in a patch, and  $c_1$  is the maximum ratio between the volumes of neighboring elements. The ratio of the upper bound to the lower one is approximately  $d(d + 2)^2 c_1 p_{\max} \kappa(\mathbb{D})$ , where  $\kappa(\mathbb{D}) = \lambda_{\max}(\mathbb{D}_K) / \lambda_{\min}(\mathbb{D}_K)$ .

Geometric bound (38) is similar to (41), but there is a significant difference. Since  $\mathbf{Z}_K \sim \|(F'_K)^{-1}(F'_K)^{-T}\|_2$ , the interplay between the mesh geometry and the diffusion matrix in (38) and (41) is mainly reflected by

$$\left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \quad \text{and} \quad \lambda_{\max}(\mathbb{D}_K) \left\| (F'_K)^{-1} (F'_K)^{-T} \right\|_2,$$

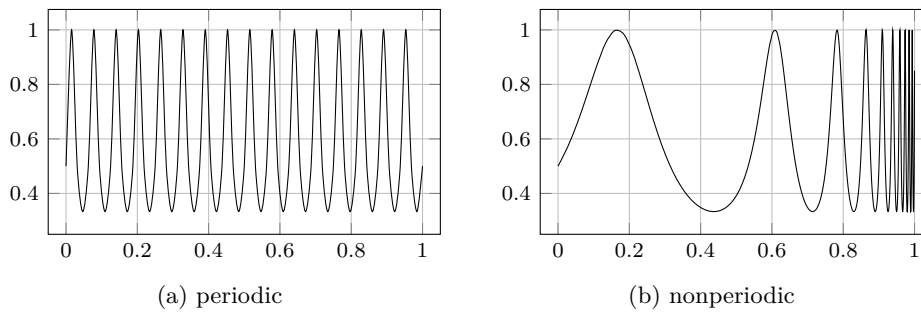
respectively. If either the mesh or  $\mathbb{D}$  is isotropic, then the factors are comparable. However, if both the mesh and  $\mathbb{D}$  are anisotropic, then the former factor can be much smaller than the latter. In the worst situation, the accuracy of (41) can deteriorate proportionally to  $\kappa(\mathbb{D})$  (see Example 12 in section 4), whereas the bound (34) in Theorem 5 in terms of matrix entries is sharp within a factor of at most  $2(d + 1)$ , independently of the mesh and  $\mathbb{D}$ .

In the case of mass lumping, Shewchuk [23, sect. 3] obtained geometric bounds in 2D and three dimensions (3D). The bounds can be generalized to any dimension as

$$(42) \quad \frac{1}{d} \max_K \mathbf{S}_K \leq \lambda_{\max}(\tilde{M}^{-1}A) \leq p_{\max} \max_K \mathbf{S}_K, \quad \mathbf{S}_K = \frac{1}{d^2} \sum_{i_K} \frac{|K|}{\tilde{M}_{i_K i_K}} \frac{|V_{i_K}|_{\mathbb{D}^{-1}}^2}{|K|_{\mathbb{D}^{-1}}^2},$$

where  $|V_{i_K}|_{\mathbb{D}^{-1}}$  is the volume of a  $(d - 1)$ -dimensional face opposing the  $i_K$ th vertex of  $K$  with respect to  $\mathbb{D}^{-1}$  and  $\tilde{M}_{i_K i_K}$  is the entry of the (global) lumped mass matrix corresponding to the node  $i_K$ . The bound takes the interplay between the mesh shape and  $\mathbb{D}$  fully into account and is tight within a factor of  $d p_{\max}$ , independently of  $\mathbb{D}$ , but it still has a weak mesh dependence through  $p_{\max}$  (typically,  $p_{\max} \geq 6$  in 2D and can be much larger in higher dimensions). Numerical examples in section 4 show that it is comparable but less accurate than bound (34) obtained in this paper.

For the lumped case there is also an earlier result by Zhu and Du [26, Theorem 3.1], but we omit it in this study since it is less accurate than Shewchuk's bound (42).

FIG. 2. Diffusion coefficients  $\mathbb{D}$  in 1D (Example 9).

**4. Numerical examples.** To test the developed estimates, we continue Example 6 (stabilized RK methods) and compare the exact value of the largest permissible time step (36) with the lower bound (37),

$$\tau_{\max} = \frac{2s^2}{\lambda_{\max}(M^{-1}A)} \quad \text{and} \quad \tau_h = \frac{2s^2}{C_*} \min_i \frac{M_{ii}}{A_{ii}},$$

and compute the ratio  $\tau_{\max}/\tau_h$  to evaluate the accuracy of the estimate. Since  $\tau_{\max}/\tau_h$  is independent of the number of stages  $s$ , we rescale the values of  $\tau_{\max}$  and  $\tau_h$  by  $s^{-2}$  to stay general; i.e., in the following we compare

$$(43) \quad \frac{\tau_{\max}}{s^2} = \frac{2}{\lambda_{\max}(M^{-1}A)} \quad \text{with} \quad \frac{\tau_h}{s^2} = \frac{2}{C_*} \min_i \frac{M_{ii}}{A_{ii}}.$$

Note that (37) implies that  $1 \leq \tau_{\max}/\tau_h \leq C_*$  for any mesh and any diffusion matrix  $\mathbb{D}$ . Moreover, from (40),

$$(44) \quad \frac{\tau_{\max}}{s^2} \geq \frac{2h_{\mathbb{M}}^2}{C_* C_{\#} \max_i \frac{1}{|\omega_i|} \sum_{K \in \omega_i} |K| \cdot \|\mathbb{M}_K \mathbb{D}_K\|_2}.$$

*Example 9* (one-dimensional (1D) example [21, sects. 6.1 and 6.2]). As a first example we consider the heat diffusion  $u_t = (\mathbb{D}u_x)_x$  in  $\Omega = (0, 1)$  with the diffusion coefficients

$$\mathbb{D}(x) = \left(2 - \sin\left(2\pi \frac{x}{\varepsilon}\right)\right)^{-1} \quad \text{and} \quad \mathbb{D}(x) = \left(2 - \sin\left(2\pi \tan\left(\frac{(1-\varepsilon)\pi x}{2}\right)\right)\right)^{-1},$$

where  $\varepsilon$  is a positive parameter (Figure 2). We choose  $\varepsilon = 2^{-4}$  for our tests.

Numerical results in Table 2 show that  $1.00 \leq \tau_{\max}/\tau_h \leq 1.45$  for all considered meshes and cases, which is consistent with the theoretical prediction  $1 \leq \tau_{\max}/\tau_h \leq 2$  (with mass lumping) and  $1 \leq \tau_{\max}/\tau_h \leq 4$  (without mass lumping). Interestingly, for this example, the estimate appears to be even asymptotically exact ( $\tau_{\max}/\tau_h \rightarrow 1$  as  $N \rightarrow \infty$ ) except for the case of  $\mathbb{D}^{-1}$ -uniform meshes with mass lumping.

Table 2 further shows that in case of mass lumping  $\tau_{\max}$  is roughly three times as large as  $\tau_{\max}$  without mass lumping. The largest permissible time step  $\tau_{\max}$  for  $\mathbb{D}^{-1}$ -uniform meshes is approximately 1.4 to 1.8 times as large as for uniform meshes.

TABLE 2  
Numerical results in 1D (Example 9).

(a) periodic  $\mathbb{D}$  (Figure 2(a))

N	With mass lumping			Without mass lumping		
	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
uniform meshes						
64	$1.84 \times 10^{-4}$	$1.27 \times 10^{-4}$	1.45	$6.57 \times 10^{-5}$	$5.10 \times 10^{-5}$	1.29
128	$3.79 \times 10^{-5}$	$3.26 \times 10^{-5}$	1.16	$1.45 \times 10^{-5}$	$1.09 \times 10^{-5}$	1.32
256	$8.66 \times 10^{-6}$	$7.76 \times 10^{-6}$	1.12	$3.11 \times 10^{-6}$	$2.59 \times 10^{-6}$	1.20
512	$2.04 \times 10^{-6}$	$1.91 \times 10^{-6}$	1.06	$7.08 \times 10^{-7}$	$6.37 \times 10^{-7}$	1.11
1024	$4.93 \times 10^{-7}$	$4.77 \times 10^{-7}$	1.03	$1.68 \times 10^{-7}$	$1.59 \times 10^{-7}$	1.06
2048	$1.21 \times 10^{-7}$	$1.19 \times 10^{-7}$	1.02	$4.09 \times 10^{-8}$	$3.97 \times 10^{-8}$	1.03
$\mathbb{D}^{-1}$ -uniform meshes						
64	$2.30 \times 10^{-4}$	$1.86 \times 10^{-4}$	1.23	$7.67 \times 10^{-5}$	$7.54 \times 10^{-5}$	1.02
128	$5.86 \times 10^{-5}$	$4.86 \times 10^{-5}$	1.21	$1.96 \times 10^{-5}$	$1.94 \times 10^{-5}$	1.01
256	$1.47 \times 10^{-5}$	$1.22 \times 10^{-5}$	1.21	$4.91 \times 10^{-6}$	$4.90 \times 10^{-6}$	1.00
512	$3.69 \times 10^{-6}$	$3.06 \times 10^{-6}$	1.21	$1.23 \times 10^{-6}$	$1.23 \times 10^{-6}$	1.00
1024	$9.22 \times 10^{-7}$	$7.67 \times 10^{-7}$	1.20	$3.07 \times 10^{-7}$	$3.07 \times 10^{-7}$	1.00
2048	$2.31 \times 10^{-7}$	$1.92 \times 10^{-7}$	1.20	$7.68 \times 10^{-8}$	$7.68 \times 10^{-8}$	1.00

(b) nonperiodic  $\mathbb{D}$  (Figure 2(b))

N	With mass lumping			Without mass lumping		
	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
uniform meshes						
64	$1.25 \times 10^{-4}$	$1.19 \times 10^{-4}$	1.05	$4.31 \times 10^{-5}$	$3.96 \times 10^{-5}$	1.09
128	$3.09 \times 10^{-5}$	$3.01 \times 10^{-5}$	1.03	$1.05 \times 10^{-5}$	$1.00 \times 10^{-5}$	1.04
256	$7.67 \times 10^{-6}$	$7.57 \times 10^{-6}$	1.01	$2.58 \times 10^{-6}$	$2.52 \times 10^{-6}$	1.02
512	$1.91 \times 10^{-6}$	$1.90 \times 10^{-6}$	1.01	$6.41 \times 10^{-7}$	$6.33 \times 10^{-7}$	1.01
1024	$4.78 \times 10^{-7}$	$4.76 \times 10^{-7}$	1.00	$1.60 \times 10^{-7}$	$1.59 \times 10^{-7}$	1.01
2048	$1.19 \times 10^{-7}$	$1.19 \times 10^{-7}$	1.00	$3.98 \times 10^{-8}$	$3.97 \times 10^{-8}$	1.00
$\mathbb{D}^{-1}$ -uniform meshes						
64	$2.04 \times 10^{-4}$	$1.68 \times 10^{-4}$	1.22	$7.09 \times 10^{-5}$	$6.59 \times 10^{-5}$	1.08
128	$5.28 \times 10^{-5}$	$4.18 \times 10^{-5}$	1.26	$1.82 \times 10^{-5}$	$1.67 \times 10^{-5}$	1.09
256	$1.32 \times 10^{-5}$	$1.10 \times 10^{-5}$	1.21	$4.53 \times 10^{-6}$	$4.24 \times 10^{-6}$	1.07
512	$3.43 \times 10^{-6}$	$2.76 \times 10^{-6}$	1.25	$1.15 \times 10^{-6}$	$1.12 \times 10^{-6}$	1.02
1024	$8.65 \times 10^{-7}$	$6.98 \times 10^{-7}$	1.24	$2.89 \times 10^{-7}$	$2.86 \times 10^{-7}$	1.01
2048	$2.17 \times 10^{-7}$	$1.77 \times 10^{-7}$	1.22	$7.23 \times 10^{-8}$	$7.20 \times 10^{-8}$	1.00

Example 10 (2D example,  $\mathbb{D} = I$ ). In this example we consider the most simple case of  $\mathbb{D} = I$ . Mesh examples are taken from [26, 27]; they are uniform isotropic, uniform anisotropic, and strongly refined toward the boundary (Figure 3). Since these meshes have no obtuse angles, we can use sharper bounds with  $C_* = 2$  (mass lumping) or  $C_* = 4$  (no mass lumping), and therefore we expect that  $1 \leq \tau_{\max}/\tau_h \leq 2$  or  $1 \leq \tau_{\max}/\tau_h \leq 4$ , respectively.

Table 3 shows that  $1.14 \leq \tau_{\max}/\tau_h \leq 1.69$  (mass lumping) and  $1.18 \leq \tau_{\max}/\tau_h \leq 2.33$  (no mass lumping). In comparison, the same ratio if using (41) and (42) ranges

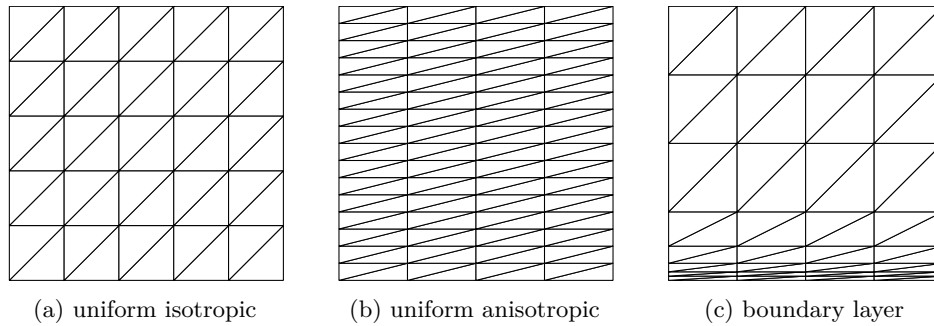


FIG. 3. Mesh examples in 2D (Example 10).

from 1.78 to 3.50<sup>1</sup> and 4.00 to 6.77, respectively. In this example  $\mathbb{D} = I$ , so that the difference is mainly due to the fact that estimates in terms of mesh geometry are generally less tight than those in terms of matrix entries since additional estimation steps decrease the accuracy.

Notice the significant reduction of  $\tau_{\max}$  when the mesh gets adapted in the “wrong” way, i.e., away from  $\mathbb{D}^{-1}$ . For example, a  $32 \times 32$  uniform mesh requires  $\tau_{\max} = 2.38 \times 10^{-4}$ , whereas the  $4 \times 256$  mesh with the same number of elements requires  $\tau_{\max} = 6.36 \times 10^{-6}$ , a reduction by a factor of 37. A strongly anisotropic  $4 \times 16$  mesh adapted toward the boundary with a much smaller number of elements leads to the further reduction of the step size by a factor of 3000. Thus, the matching between the element geometry and the diffusion matrix has significant effects on the time step size, and, depending on the anisotropy of the mesh and diffusion matrix, changes in the mesh alignment can result in changes in the time step size by orders of magnitude.

Again, mass lumping allows approximately 1.9 to 3.2 times larger time steps.

*Example 11* (2D groundwater flow with jumping coefficients [19]). As the next example we consider groundwater flow through an aquifer. The problem is given by the IBVP (1) with  $\Omega = (0, 100) \times (0, 100)$  and two impermeable subdomains  $\Omega_1 = (0, 80) \times (64, 68)$  and  $\Omega_2 = (20, 100) \times (40, 44)$ . Figure 4(a) shows the diffusion matrix  $\mathbb{D}$  and the boundary conditions. Although  $\mathbb{D}$  is isotropic, it has a jump between the subdomains, leading to the anisotropic behavior of the solution.

We compute the solution by  $h$ -refinement in the standard way and use Hessian recovery based mesh adaptation to obtain adaptive meshes at particular time points and compare the exact  $\tau_{\max}$  with the lower bound  $\tau_h$ . For our computation we used KARDOS [2] to solve the PDE and BAMG [9] for mesh generation. Examples of adaptive meshes are shown in Figure 4 for the time points  $t = 1.0 \times 10^4$  and  $t = 1.0 \times 10^5$ . Note that these meshes have oblique elements and angles close to 180°: the maximum angles in Figures 4(b) and 4(c) are 175° and 177°, respectively.

Table 4a shows that the ratio  $\tau_{\max}/\tau_h$  is about 2.13 to 2.48 with mass lumping and 3.25 to 3.87 without mass lumping, which is consistent with the theoretical upper bounds  $d + 1 = 3$  and  $2(d + 1) = 6$ . In this example, mass lumping would allow 2.6 to 2.8 times larger time steps, which is similar to Example 10 (a factor of 1.9 to 3.2 there).

<sup>1</sup>In our tests, the estimate by Zhu and Du [27] seems to provide better results than that in the numerical examples of the original paper.

TABLE 3  
 Numerical results in 2D (Example 10).

Without mass lumping			New estimate (43)		Zhu & Du [27]	
Mesh	$N$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
uniform isotropic (Figure 10(a))						
$8 \times 8$	128	$1.31 \times 10^{-3}$	$9.77 \times 10^{-4}$	1.34	$6.51 \times 10^{-4}$	2.01
$16 \times 16$	512	$3.09 \times 10^{-4}$	$2.44 \times 10^{-4}$	1.27	$1.63 \times 10^{-4}$	1.90
$32 \times 32$	2048	$7.60 \times 10^{-5}$	$6.10 \times 10^{-5}$	1.24	$4.07 \times 10^{-5}$	1.87
$64 \times 64$	8192	$1.89 \times 10^{-5}$	$1.53 \times 10^{-5}$	1.24	$1.02 \times 10^{-5}$	1.86
$128 \times 128$	32768	$4.72 \times 10^{-6}$	$3.81 \times 10^{-6}$	1.24	$2.54 \times 10^{-6}$	1.86
uniform anisotropic (Figure 10(b))						
$32 \times 32$	2048	$7.60 \times 10^{-5}$	$6.10 \times 10^{-5}$	1.24	$4.07 \times 10^{-5}$	1.87
$16 \times 64$	2048	$3.40 \times 10^{-5}$	$2.87 \times 10^{-5}$	1.18	$1.91 \times 10^{-5}$	1.78
$8 \times 128$	2048	$9.00 \times 10^{-6}$	$7.60 \times 10^{-6}$	1.18	$5.07 \times 10^{-6}$	1.78
$4 \times 256$	2048	$2.38 \times 10^{-6}$	$1.91 \times 10^{-6}$	1.25	$1.27 \times 10^{-6}$	1.87
$2 \times 512$	2048	$6.36 \times 10^{-7}$	$4.77 \times 10^{-7}$	1.33	$3.18 \times 10^{-7}$	2.00
boundary layer (Figure 10(c))						
$4 \times 8$	64	$7.08 \times 10^{-5}$	$3.04 \times 10^{-5}$	2.33	$2.03 \times 10^{-5}$	3.49
$4 \times 10$	80	$4.45 \times 10^{-6}$	$1.91 \times 10^{-6}$	2.33	$1.27 \times 10^{-6}$	3.50
$4 \times 12$	96	$2.78 \times 10^{-7}$	$1.19 \times 10^{-7}$	2.33	$7.95 \times 10^{-8}$	3.50
$4 \times 14$	112	$1.74 \times 10^{-8}$	$7.45 \times 10^{-9}$	2.33	$4.97 \times 10^{-9}$	3.50
$4 \times 16$	128	$1.09 \times 10^{-9}$	$4.66 \times 10^{-10}$	2.33	$3.10 \times 10^{-10}$	3.50
With mass lumping			New estimate (43)		Shewchuk [23]	
Mesh	$N$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
uniform isotropic (Figure 10(a))						
$8 \times 8$	128	$3.79 \times 10^{-3}$	$2.60 \times 10^{-3}$	1.46	$6.51 \times 10^{-4}$	5.82
$16 \times 16$	512	$9.53 \times 10^{-4}$	$6.51 \times 10^{-4}$	1.46	$1.63 \times 10^{-4}$	5.86
$32 \times 32$	2048	$2.38 \times 10^{-4}$	$1.63 \times 10^{-4}$	1.46	$4.07 \times 10^{-5}$	5.86
$64 \times 64$	8192	$5.96 \times 10^{-5}$	$4.07 \times 10^{-5}$	1.46	$1.02 \times 10^{-5}$	5.86
$128 \times 128$	32768	$1.49 \times 10^{-5}$	$1.02 \times 10^{-5}$	1.46	$2.54 \times 10^{-6}$	5.86
uniform anisotropic (Figure 10(b))						
$32 \times 32$	2048	$2.38 \times 10^{-4}$	$1.63 \times 10^{-4}$	1.46	$4.07 \times 10^{-5}$	5.86
$16 \times 64$	2048	$9.86 \times 10^{-5}$	$7.66 \times 10^{-5}$	1.29	$1.91 \times 10^{-5}$	5.15
$8 \times 128$	2048	$2.54 \times 10^{-5}$	$2.03 \times 10^{-5}$	1.25	$5.07 \times 10^{-6}$	5.01
$4 \times 256$	2048	$6.36 \times 10^{-6}$	$5.09 \times 10^{-6}$	1.25	$1.27 \times 10^{-6}$	5.00
$2 \times 512$	2048	$1.27 \times 10^{-6}$	$1.11 \times 10^{-6}$	1.14	$3.18 \times 10^{-7}$	4.00
boundary layer (Figure 10(c))						
$4 \times 8$	64	$1.37 \times 10^{-4}$	$8.11 \times 10^{-5}$	1.69	$2.03 \times 10^{-5}$	6.77
$4 \times 10$	80	$8.61 \times 10^{-6}$	$5.09 \times 10^{-6}$	1.69	$1.27 \times 10^{-6}$	6.77
$4 \times 12$	96	$5.38 \times 10^{-7}$	$3.18 \times 10^{-7}$	1.69	$7.95 \times 10^{-8}$	6.77
$4 \times 14$	112	$3.36 \times 10^{-8}$	$1.99 \times 10^{-8}$	1.69	$4.97 \times 10^{-9}$	6.77
$4 \times 16$	128	$2.10 \times 10^{-9}$	$1.24 \times 10^{-9}$	1.69	$3.10 \times 10^{-10}$	6.77

In a practical computation, however, one would rather use a numerical approximation for  $\lambda_{\max}(M^{-1}A)$ . Typically, five steps of the Lanczos method with a random starting vector approximate the largest eigenvalue within 10%. Another practical alternative is the power method, for which it is reported [22, sect. 3.2] that, for the case of eigenvalues being close to the negative real axis, usually only a few iterations are required if the computed eigenvector from the previous step is used as a new

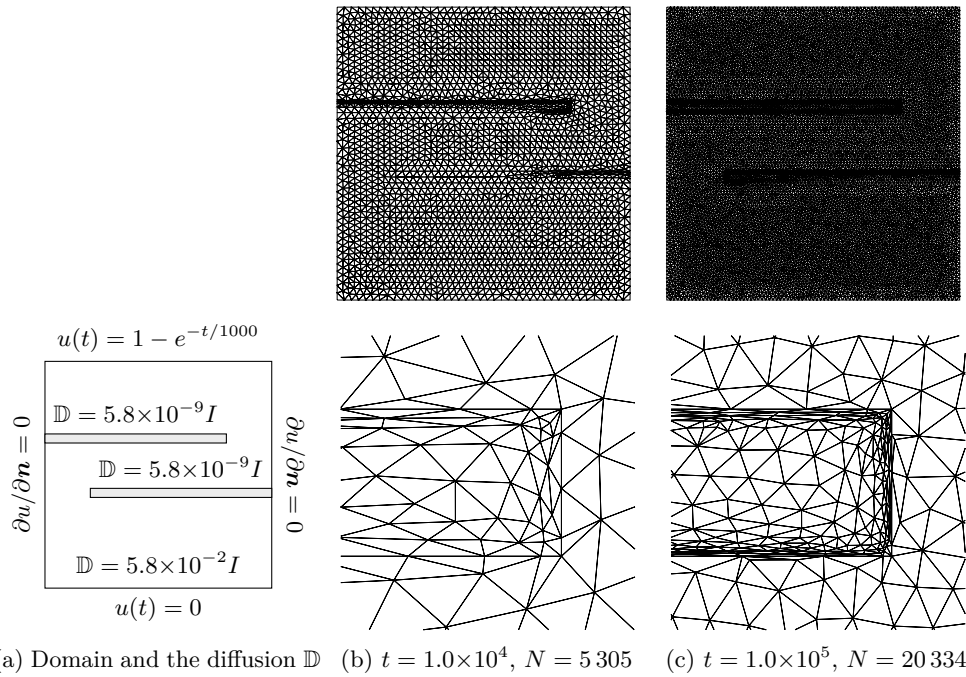


FIG. 4. Domain, mesh examples, and close-ups at  $[74, 82] \times [62, 70]$  (the upper right corner at the entrance of the tunnel) for the groundwater flow (Example 11).

starting vector. To compare it with our theoretical estimate, we additionally computed  $\tau_h$  using five steps of the Lanczos method with a random starting vector (divided by 1.1 as a security factor since Lanczos approximation is an approximation from below). Table 4b shows that the corresponding ratio  $\tau_{\max}/\tau_h$  is about 1.00 to 1.07; i.e., the computed time step approximation is within 7% from the optimal value. In our computations, the accuracy of our theoretical estimate (43) corresponds to about two to three steps of the Lanczos method.

We would like to also point out that the lower bound in (34) can be used as a practical security check for a numerical approximation: if the computed numerical approximation of  $\lambda_{\max}(M^{-1}A)$  is smaller than this bound, the time step is guaranteed to be out of the stability region of the time integration method.

*Example 12* (2D anisotropic diffusion). This example shows the importance of the interplay between the major diffusion directions and the mesh geometry.

Consider the IBVP (1) in  $\Omega = (0, 1)^2 \setminus [\frac{4}{9}, \frac{5}{9}]^2$  with the homogeneous Dirichlet boundary condition and

$$\mathbb{D} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} 1000 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \quad \theta = \pi \sin x \cos y.$$

First, we consider quasi-uniform meshes (Figure 5(a)), for which elements are close to being uniform in shape and size,  $F'_K \approx |K|^{\frac{1}{2}} I$ , and  $\|(F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}\|_2 \approx \lambda_{\max}(\mathbb{D}) \|(F'_K)^{-1} (F'_K)^{-T}\|_2$ . Hence, using (39) and (41) provides comparable results,

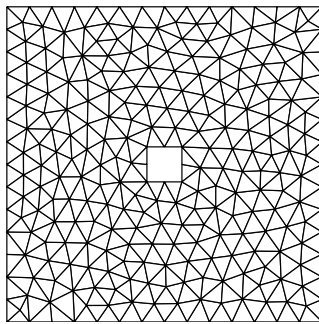
TABLE 4  
 Numerical results for the groundwater flow (Example 11).

(a) Computing  $\tau_h$  with (43)

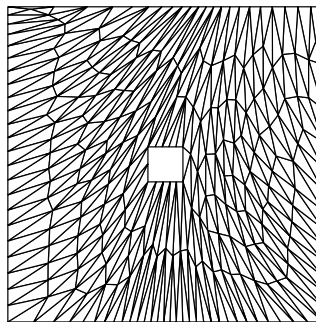
Time	N	With mass lumping		Without mass lumping			
		$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
$1.0 \times 10^2$	3 071	$1.48 \times 10^0$	$5.97 \times 10^{-1}$	2.48	$5.77 \times 10^{-1}$	$1.49 \times 10^{-1}$	3.87
$5.0 \times 10^3$	2 799	$4.74 \times 10^0$	$2.23 \times 10^0$	2.13	$1.81 \times 10^0$	$5.57 \times 10^{-1}$	3.25
$1.0 \times 10^4$	5 305	$1.80 \times 10^0$	$8.01 \times 10^{-1}$	2.25	$6.89 \times 10^{-1}$	$2.00 \times 10^{-1}$	3.44
$1.0 \times 10^5$	20 334	$2.05 \times 10^{-1}$	$9.11 \times 10^{-2}$	2.25	$7.45 \times 10^{-2}$	$2.28 \times 10^{-2}$	3.27

(b) Computing  $\tau_h$  with five steps of the Lanczos method using a random starting vector

Time	N	With mass lumping		Without mass lumping			
		$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_{\max}}{s^2}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
$1.0 \times 10^2$	3 071	$1.48 \times 10^0$	$1.48 \times 10^0$	1.00	$5.77 \times 10^{-1}$	$5.64 \times 10^{-1}$	1.02
$5.0 \times 10^3$	2 799	$4.74 \times 10^0$	$4.44 \times 10^0$	1.07	$1.81 \times 10^0$	$1.77 \times 10^0$	1.02
$1.0 \times 10^4$	5 305	$1.80 \times 10^0$	$1.69 \times 10^0$	1.07	$6.89 \times 10^{-1}$	$6.46 \times 10^{-1}$	1.07
$1.0 \times 10^5$	20 334	$2.05 \times 10^{-1}$	$1.98 \times 10^{-1}$	1.03	$7.45 \times 10^{-2}$	$7.05 \times 10^{-2}$	1.06



(a) quasi-uniform



(b)  $\mathbb{D}^{-1}$ -uniform

FIG. 5. Mesh examples for the anisotropic diffusion (Example 12).

which is confirmed by the numerical results in Table 5: for quasi-uniform grids, (39) and (41) or (42) are accurate within a factor of 4.04 to 6.35 and 4.52 to 6.02, respectively.

For  $\mathbb{D}^{-1}$ -uniform (coefficient-adaptive) meshes (Figure 5(b)) the situation is quite different, and, as mentioned in Remark 8, bound (39) should be more accurate than that obtained when using (41). This is indeed confirmed by the numerical results: bound (39) is accurate within a factor of 3.40 to 6.44, whereas (41) underestimates the real value by a factor of 347 to 1 020 (recalling that  $\kappa(\mathbb{D}) = 1000$ ). Note that Shewchuk’s bound (42) provides accurate results in any case, although not quite as accurate as (39). It is worth pointing out that the most accurate bound in all cases is (43) in terms of the matrix entries.

This example also shows that  $\mathbb{D}^{-1}$ -uniform meshes allow larger time steps even if their elements may have “bad quality” in the common sense. Hence, it is important to consider the quality of the mesh *in relation to the diffusion* and not by itself.

TABLE 5  
*Numerical results for the anisotropic diffusion (Example 12).*

(a) Without mass lumping

$N$	$\frac{\tau_{\max}}{s^2}$	New estimate (43)		Geometric (44)		Zhu & Du [27]	
		$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
quasi-uniform meshes (Figure 5(a))							
2 050	$1.06 \times 10^{-7}$	$3.23 \times 10^{-8}$	3.28	$1.67 \times 10^{-8}$	6.35	$2.34 \times 10^{-8}$	4.53
8 206	$2.67 \times 10^{-8}$	$8.59 \times 10^{-9}$	3.10	$4.31 \times 10^{-9}$	6.19	$5.90 \times 10^{-9}$	4.52
32 742	$6.18 \times 10^{-9}$	$2.03 \times 10^{-9}$	3.05	$1.15 \times 10^{-9}$	5.36	$1.18 \times 10^{-9}$	5.26
132 468	$1.34 \times 10^{-9}$	$4.49 \times 10^{-10}$	2.98	$2.25 \times 10^{-10}$	5.95	$2.71 \times 10^{-10}$	4.93
$\mathbb{D}^{-1}$ -uniform meshes (Figure 5(b))							
2 058	$6.17 \times 10^{-7}$	$2.11 \times 10^{-7}$	2.92	$9.58 \times 10^{-8}$	6.44	$6.05 \times 10^{-10}$	1 020
8 257	$1.77 \times 10^{-7}$	$8.64 \times 10^{-8}$	2.05	$5.22 \times 10^{-8}$	3.40	$2.42 \times 10^{-10}$	733
32 669	$5.97 \times 10^{-8}$	$3.03 \times 10^{-8}$	1.97	$1.63 \times 10^{-8}$	3.66	$6.37 \times 10^{-11}$	937
132 053	$7.43 \times 10^{-10}$	$2.31 \times 10^{-10}$	3.22	$1.64 \times 10^{-10}$	4.52	$2.14 \times 10^{-12}$	347

(b) With mass lumping

$N$	$\frac{\tau_{\max}}{s^2}$	New estimate (43)		Geometric (44)		Shewchuk [23]	
		$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$	$\frac{\tau_h}{s^2}$	$\frac{\tau_{\max}}{\tau_h}$
quasi-uniform meshes (Figure 5(a))							
2 050	$2.98 \times 10^{-7}$	$1.29 \times 10^{-7}$	2.31	$6.68 \times 10^{-8}$	4.46	$5.07 \times 10^{-8}$	5.87
8 206	$6.86 \times 10^{-8}$	$3.42 \times 10^{-8}$	2.01	$1.66 \times 10^{-8}$	4.13	$1.51 \times 10^{-8}$	4.54
32 742	$1.76 \times 10^{-8}$	$8.12 \times 10^{-9}$	2.16	$4.35 \times 10^{-9}$	4.04	$3.31 \times 10^{-9}$	5.31
132 468	$4.06 \times 10^{-9}$	$1.77 \times 10^{-9}$	2.30	$8.99 \times 10^{-10}$	4.51	$6.74 \times 10^{-10}$	6.02
$\mathbb{D}^{-1}$ -uniform meshes (Figure 5(b))							
2 058	$1.10 \times 10^{-6}$	$5.56 \times 10^{-7}$	1.98	$2.53 \times 10^{-7}$	4.35	$2.26 \times 10^{-7}$	4.87
8 257	$4.47 \times 10^{-7}$	$2.25 \times 10^{-7}$	1.99	$1.39 \times 10^{-7}$	3.22	$5.91 \times 10^{-8}$	7.56
32 669	$1.51 \times 10^{-7}$	$7.57 \times 10^{-8}$	2.00	$4.05 \times 10^{-8}$	3.74	$1.92 \times 10^{-8}$	7.87
132 053	$1.66 \times 10^{-9}$	$9.22 \times 10^{-10}$	1.79	$6.57 \times 10^{-10}$	2.52	$2.17 \times 10^{-10}$	7.64

**5. Conclusions.** Theorem 5 gives an easily computable bound on the largest eigenvalue of the system matrix  $\tilde{M}^{-1}A$  in terms of the diagonal entries of  $\tilde{M}$  and  $A$  with  $\tilde{M}$  being either  $M$  or  $M_{lump}$ . The bound is tight for *any mesh* and *any diffusion matrix*  $\mathbb{D}$  within a small constant which is given explicitly and depends only on the dimension of the domain. This allows efficient and accurate estimation of the largest permissible time step  $\tau_{\max}$ .

Moreover, estimates (38) and (40) in terms of the mesh geometry reveal how the mesh and the diffusion matrix affect the stability condition. Roughly speaking,  $\tau_{\max}$  depends only on the number of mesh elements and the matching between the element geometry with the diffusion matrix. Thus, it is not the element geometry itself but the *element geometry in relation to the diffusion matrix* that is important for the stability. The element quality measure  $Q_{\mathbb{D}^{-1}}$  provides a measure for the effect of a given element on the stability condition. As seen in Example 10, strong anisotropic adaptation in the “wrong” direction can cause a significant reduction of the time step size. Meanwhile, the result suggests that improvements in the element quality can significantly increase  $\tau_{\max}$ .

The achieved result can be extended for high order [13] or even  $p$ -adaptive finite elements without major modifications. Essentially, one only needs to recalculate the constants which depend on the choice of the basis functions.

Furthermore, numerical results suggest that, at least in 1D and 2D, mass lumping can increase the time step size by a factor of 2 to 3. This topic deserves more detailed investigations.

**Acknowledgments.** Lennard Kamenski is thankful to Klaus Gärtner for a helpful comment that led to Remark 2 and to Larissa Kaspar for providing parts of the code used in computations in Example 11.

The authors are grateful to an anonymous referee and particularly to Jed Brown for their valuable comments and suggestions which helped to improve the quality of this paper.

## REFERENCES

- [1] Q. DU, D. WANG, AND L. ZHU, *On mesh geometry and stiffness matrix conditioning for general finite element spaces*, SIAM J. Numer. Anal., 47 (2009), pp. 1421–1444, <http://dx.doi.org/10.1137/080718486>.
- [2] B. ERDMANN, J. LANG, AND R. ROITZSCH, *Kardos: User's Guide*, ZIB-Report 02-42, ZIB, 2002, <http://www.zib.de/software>.
- [3] T. ERTEKIN, J. H. ABOU-KASSEM, AND G. R. KING, *Basic Applied Reservoir Simulation*, SPE Textbook Series 7, SPE, Richardson, TX, 2001.
- [4] I. FRIED, *Bounds on the spectral and maximum norms of the finite element stiffness, flexibility and mass matrices*, Internat. J. Solids Structures, 9 (1973), pp. 1013–1034.
- [5] I. G. GRAHAM AND W. MCLEAN, *Anisotropic mesh refinement: The conditioning of Galerkin boundary element matrices and simple preconditioners*, SIAM J. Numer. Anal., 44 (2006), pp. 1487–1513, <http://dx.doi.org/10.1137/040621247>.
- [6] J.-L. GUERMOND AND R. PASQUETTI, *A correction technique for the dispersive effects of mass lumping for transport problems*, Comput. Methods Appl. Mech. Engrg., 253 (2013), pp. 186–198.
- [7] S. GÜNTER AND K. LACKNER, *A mixed implicit-explicit finite difference scheme for heat transport in magnetised plasmas*, J. Comput. Phys., 228 (2009), pp. 282–293.
- [8] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*, 2nd ed., Springer Ser. Comput. Math. 14, Springer-Verlag, Berlin, 1996.
- [9] F. HECHT, *BAMG: Bidimensional Anisotropic Mesh Generator*, <http://www.ann.jussieu.fr/hecht/ftp/bamg/>.
- [10] W. HUANG, *Measuring mesh qualities and application to variational mesh adaptation*, SIAM J. Sci. Comput., 26 (2005), pp. 1643–1666, <http://dx.doi.org/10.1137/S1064827503429405>.
- [11] W. HUANG, *Metric tensors for anisotropic mesh generation*, J. Comput. Phys., 204 (2005), pp. 633–665.
- [12] W. HUANG, *Anisotropic mesh adaptation and movement*, in Adaptive computations: Theory and Algorithms, T. Tang and J. Xu, eds., Mathematics Monograph Series 6, Science Press, Beijing, China, 2007, Ch. 3, pp. 68–158.
- [13] W. HUANG, L. KAMENSKI, AND J. LANG, *Stability of explicit Runge-Kutta methods for high order finite element approximation of linear parabolic equations*, in Numerical Mathematics and Advanced Applications - ENUMATH 2013, A. Abdulle, S. Deparis, D. Kressner, F. Nobile, and M. Picasso, eds., Lect. Notes Comput. Sci. Eng. 103, Springer, Berlin, 2015, pp. 165–173.
- [14] W. HUANG AND R. D. RUSSELL, *Adaptive Moving Mesh Methods*, Appl. Math. Sci. 174, Springer, New York, 2011.
- [15] W. HUNSDORFER AND J. G. VERWER, *Numerical Solution of Time-Dependent Advection-Diffusion-Reaction Equations*, Springer Ser. Comput. Math. 33, Springer, New York, 2003.
- [16] L. KAMENSKI, W. HUANG, AND H. XU, *Conditioning of finite element equations with arbitrary anisotropic meshes*, Math. Comp., 83 (2014), pp. 2187–2211.
- [17] D. A. KARRAS AND G. B. MERTZIOS, *New PDE-based methods for image enhancement using SOM and Bayesian inference in various discretization schemes*, Meas. Sci. Technol., 20 (2009), 104012.

- [18] X. LI AND W. HUANG, *An anisotropic mesh adaptation method for the finite element solution of heterogeneous anisotropic diffusion problems*, J. Comput. Phys., 229 (2010), pp. 8072–8094.
- [19] S. MICHELETTI AND S. PEROTTO, *Anisotropic mesh adaption for time-dependent problems*, Internat. J. Numer. Methods Fluids, 58 (2008), pp. 1009–1015.
- [20] M. J. MLACNIK AND L. J. DURLOFSKY, *Unstructured grid optimization for improved monotonicity of discrete solutions of elliptic equations with highly anisotropic coefficients*, J. Comput. Phys., 216 (2006), pp. 337–361.
- [21] D. PETERSEIM AND S. A. SAUTER, *Finite elements for elliptic problems with highly varying, nonperiodic diffusion matrix*, Multiscale Model. Simul., 10 (2012), pp. 665–695, <http://dx.doi.org/10.1137/10081839X>.
- [22] B. P. SOMMEIJER, L. F. SHAMPINE, AND J. G. VERWER, *RKC: An explicit solver for parabolic PDEs*, J. Comput. Appl. Math., 88 (1998), pp. 315–326.
- [23] J. R. SHEWCHUK, *What Is a Good Linear Finite Element? Interpolation, Conditioning, Anisotropy, and Quality Measures*, University of California - Berkeley, Berkeley, CA, 2002, available online from <http://www.cs.berkeley.edu/~jrs/jrspapers.html#quality>.
- [24] A. J. WATHEN, *Realistic eigenvalue bounds for the Galerkin mass matrix*, IMA J. Numer. Anal., 7 (1987), pp. 449–457.
- [25] J. WEICKERT, *Anisotropic Diffusion in Image Processing*, Teubner-Verlag, Stuttgart, Germany, 1998.
- [26] L. ZHU AND Q. DU, *Mesh-dependent stability for finite element approximations of parabolic equations with mass lumping*, J. Comput. Appl. Math., 236 (2011), pp. 801–811.
- [27] L. ZHU AND Q. DU, *Mesh dependent stability and condition number estimates for finite element approximations of parabolic problems*, Math. Comp., 83 (2014), pp. 37–64.