# OPTIMAL ORDER MULTILEVEL PRECONDITIONERS FOR REGULARIZED ILL-POSED PROBLEMS

ANDREI DRĂGĂNESCU AND TODD F. DUPONT

ABSTRACT. In this article we design and analyze multilevel preconditioners
for linear systems arising from regularized inverse problems. Using a scale-
independent distance function that measures spectral equivalence of operators,
it is shown that these preconditioners approximate the inverse of the opera-
tor to optimal order with respect to the spatial discretization parameter $h$.
As a consequence, the number of preconditioned conjugate gradient iterations
needed for solving the system will **decrease** when increasing the number of
levels, with the possibility of performing only one fine-level residual compu-
tation if $h$ is small enough. The results are based on the previously known
two-level preconditioners of Rieder (1997) (see also Hanke and Vogel (1999)),
and on applying Newton-like methods to the operator equation $X^{-1} - A = 0$.
We require that the associated forward problem has certain smoothing proper-
ties; however, only natural stability and approximation properties are assumed
for the discrete operators. The algorithm is applied to a reverse-time parabolic
equation, that is, the problem of finding the initial value leading to a given
final state. We also present some results on constructing restriction operators
with preassigned approximating properties that are of independent interest.

## 1. INTRODUCTION

We consider the equation

(1.1) $$\mathcal{K}u = f \ ,$$

where $\mathcal{K} : \mathcal{X} \to \mathcal{Y}$ is a compact linear operator between two Hilbert spaces.
Since (1.1) is generallly ill-posed, a generalized solution is given by $u = \mathcal{K}^\dagger f$, where

$$\mathcal{K}^\dagger : \mathcal{R}(\mathcal{K}) \oplus \mathcal{R}(\mathcal{K})^\perp \to \mathcal{X}$$

is the Moore-Penrose generalized inverse of $\mathcal{K}$. For $f \in \mathcal{D}(\mathcal{K}^\dagger)$, the solution $\mathcal{K}^\dagger f$
is defined as the element $u \in \mathcal{X}$ of smallest norm that minimizes the expression
$\|\mathcal{K}u - f\|^2$. If $\mathcal{R}(\mathcal{K})$ is not closed, then, in order to allow general perturbations in
the right-hand side of (1.1), one can approximate $\mathcal{K}^\dagger$ by a variety of regularization
operators that are defined on all of $\mathcal{Y}$ (see [16] for details). Perhaps the best known
is the Tikhonov regularization, which amounts to solving the identity-perturbed

normal equation

(1.2) 
$$H_\beta u = \beta^{-1} \mathcal{K}^* f \, ,$$

where

(1.3) 
$$H_\beta \stackrel{\text{def}}{=} I + \beta^{-1} H \quad \text{and} \quad H \stackrel{\text{def}}{=} \mathcal{K}^* \mathcal{K} \, ,$$

with $\mathcal{K}^*$ being the adjoint of $\mathcal{K}$. For convenience we have multiplied the normal equation by $\beta^{-1}$. It should be noted that solving (1.2) is equivalent to applying one step of Newton's method to the minimization of the quadratic functional

(1.4) 
$$\mathcal{J}_\beta(u) = \frac{\beta^{-1}}{2} \|\mathcal{K}u - f\|^2 + \frac{1}{2} \|u\|^2$$

starting from the initial guess $u = 0$, and that $H_\beta$ is the Hessian operator associated with the quadratic $\mathcal{J}_\beta$ via the $\mathcal{X}$-inner product. We add the subscript (resp. superscript) $h$ to a vector (resp. operator) to denote its discrete version, $h$ representing the discretization parameter. In this work, we are interested in fast solvers for computing the Tikhonov regularization of $\mathcal{K}^\dagger$ for operators $\mathcal{K}$ with certain smoothing properties.

If $\mathcal{K}^* \mathcal{K}$ is smoothing, and $\beta$ is small enough, then the spectrum $\sigma(H_\beta)$ has only a few eigenvalues that are significantly greater than 1 (and they correspond to smooth eigenvectors), while the remainder of the spectrum is clustered around the value 1. Therefore unpreconditioned Krylov-type methods are expected to be efficient in inverting $H_\beta^h$. For example, we can show that, if $\mathcal{K}$ is the solution operator for the heat equation mapping the initial value onto the final-time state, then the number of conjugate gradient (CG) iterations required to solve the regularized inverse problem down to machine precision is mesh-independent; moreover, the needed number of iterations grows only logarithmically as $\beta \to 0$ (see Chapter 7 in [16], and also [12, 1]). However, the number of iterations needed for convergence, even though independent of resolution, may be too large for practical use in the case of large-scale problems, where the application of $\mathcal{K}^h$ (i.e., the *direct problem*) requires, for example, solving a time-dependent three-dimensional partial differential equation. A standard strategy for solving large-scale optimization (or inverse) problems is to use reduced-order models, and perform optimization (or inversion) on the model rather than the direct (large) problem itself; this strategy usually restricts the accuracy and resolution of the solution process. In this work we show that multilevel techniques can be successfully employed in efficiently solving certain types of inverse problems without resorting to reduced models, thus producing highly resolved solutions.

Multilevel methods, especially multigrid methods (MGMs), have been used extensively to efficiently solve linear systems related to partial differential equations [2, 3, 4, 6, 8, 19]. The efficiency of MGMs relies on the availability of good *smoothers*, that is, iterators that are inexpensive to apply, and that remove high-frequency components from the approximation error. In the classical multigrid theory, originally developed for inverting elliptic operators, standard iterators such as Jacobi and Gauß-Seidel are natural choices for smoothers, this being related to the fact that the operator to be inverted has roughening properties (being of differential type). However, for operators with smoothing properties such as $H_\beta^h$, the aforementioned iterators remove the **low**-frequency components from the error, thus leaving an error that cannot be represented accurately on a coarser mesh. In [25], frequently

cited as one of the first studies of MGMs explicitly designed for ill-posed problems, a smoother was introduced based on the idea that on high-frequency eigenvectors, the operator $H_\beta^h$ acts essentially as the identity; therefore the projector on the high-frequency eigenspace acts as a smoother. A number of later studies [22, 23, 24] are dedicated to MGMs for unregularized inverse problems. "Unregularized" is meant as "not explicitly regularized"; in fact, regularization is achieved through discretization (by limiting the highest resolution that can be used) or by premature termination of the iterative process. The goal of MGMs for unregularized ill-posed problems is to design preconditioners that achieve a mesh-independent error reduction, in other words, whose approximation quality is mesh-independent. As mentioned before (see also [20, 30]), for the severely ill-posed problems under consideration, this aim is already achieved by unpreconditioned CG. In this article we show that, under certain circumstances, these types of problems can actually be solved at a lower cost than just a mesh-independent number of iterations; in fact, the number of iterations is shown to decrease with increasing number of levels. We first construct a two-level preconditioner that is closely related to the additive Schwarz preconditioner introduced in [30]; we denote this by TLAS (**T**wo-**L**evel **A**dditive **S**chwarz). The significant difference concerning TLAS between our work and [30, 20] is that the discrete operators we are considering here are not obtained by orthogonal projection from their continuous versions, but rather they are derived from natural discretizations. The only assumptions we make about the discretizations are related to their approximation and stability properties. Also, our two-level analysis is different from the one in [30, 20]; TLAS is shown here to approximate $(H_\beta^h)^{-1}$ to optimal order (see Section 4 for details). We use similar ideas to analyze a two-level preconditioner for an inverse semilinear parabolic problem in a forthcoming paper [13], the technical details therein being significantly more involved. By simply replacing the call to the coarse-level inverse operator in TLAS with a recursive call to TLAS, one could obtain a multilevel preconditioner which turns out to be of suboptimal quality. In particular, the number of preconditioned CG iterations would not decrease with increasing number of levels. However, the first Newton iterate of the operator-function $X \mapsto X^{-1} - H_\beta^h$ starting at this multilevel operator produces an optimal order multilevel preconditioner, denoted by MLAS (**M**ulti-**L**evel **A**dditive **S**chwarz). MLAS has a W-cycle structure, and we should point out that it escapes the "usual" paradigm of pre-smoothing followed by error-correction and post-smoothing; here smoothing is intertwined with error-correction, thus introducing an *inter-smoothing* step.

We would like to point out that, as in [30], the main results in this article apply in the regime when $h^p \ll \beta$, where $p$ is the convergence order of the discretization $\mathcal{K}^h$ (see also (2.5)). We would like to argue that this is the case for a class of problems that motivated our research. The small parameter $\delta$ dictating the choice of $h$ and $\beta$ for an inverse problem is related to the noise level in the data $f$; that is, one typically assumes that the actual "measured" data $f$ used in (1.1) satisfies

$$\|f - f_0\| < \delta \ ,$$

where $f_0$ is the "correct" data. The natural assumption is that $h$ should be chosen so that $h^p \sim \delta$. A priori rules for choosing $\beta$ usually take into account source conditions on $u^\dagger$ [16]: if a Hölder-type holds,

$$u^\dagger = \mathcal{K}^\dagger f_0 \in \mathcal{R}((\mathcal{K}^* \mathcal{K})^\nu) \ ,$$

then the choice $\beta \sim O(\delta^{\frac{2}{2\nu+1}})$ gives the optimal convergence rate of

$$\|u_\beta - u\| = O(\delta^{\frac{2\nu}{2\nu+1}}) \ , \quad \text{as } \delta \to 0 \ ,$$

where $u_\beta$ is the solution of (1.2). Therefore

$$\frac{h^p}{\beta} \sim \delta^{\frac{2\nu-1}{2\nu+1}} \ ,$$

which implies that, for $\nu > \frac{1}{2}$,

$$h^p \ll \beta, \quad \text{as } \delta \to 0.$$

We refer the reader also to [32, 29] for discussions regarding the optimal choice of the regularization parameter $\beta$.

One of our motivating problems is the question of recovering an early stage of an air pollutant in an atmospheric model, given later time measurements [1]. In this case the direct operator is given by $\mathcal{K} = \mathcal{S}(T)$, where $t \mapsto \mathcal{S}(t)$ is the solution operator of a linear advection-diffusion equation (see also Section 6), and $T > 0$ is a given time. If the coefficients of the parabolic equation are smooth, then $\mathcal{K}$ is infinitely smoothing. In general, Hölder-type conditions are too restrictive and should be replaced by logarithmic source conditions [21], the case in which an a priori choice of the regularization parameter leads to $\beta \sim O(\delta)$. However, in a practical and plausible scenario, where concentration measurements of the pollutant are sparsely placed in space, there is a time-lag $\tau$ between the moment of the original spill and the time when sensors determine that the pollution event has taken place. Thus it is natural to assume that the "initial" state $u^\dagger$ to be recovered has a history of length $\tau$, i.e., $u^\dagger = \mathcal{S}(\tau)u_0$. If $\tau > T$, that is, the "initial" state's history is longer than the difference $T$ between the measurements-collection time and the "initial" time, then $u^\dagger \in \mathcal{R}((\mathcal{K}^*\mathcal{K})^\nu)$ with $\nu > \frac{1}{2}$, under appropriate and reasonable assumptions on the parabolic equation (cf. Picard's criterion, see [16]). A similar situation is encountered in a data assimilation setting for, say, weather prediction. Here $\mathcal{S}(t)$ is the evolution operator for the velocity field of a fluid (air). Various direct or indirect measurements provide a part of the history of a process for which the entire state at a single time is unknown. Recovering a possible "initial state" at an artificially chosen time $t_0$ in the past enables "learning" the entire state at the current time, thus potentially improving predictions from this time on. Again, the initial state to be recovered naturally is $\mathcal{S}(\tau)u$, where $u$ is a state at an even earlier time, which again leads to a Hölder-type source condition with $\nu > \frac{1}{2}$, thus making the inverse problem amenable to the case when $h^p \ll \beta$.

This article closely follows [12] and is organized as follows: after describing the problem in Section 2, we define the spectral distance between operators with positive definite symmetric part in Section 3; this distance (a measure of spectral equivalence) introduces a framework that is convenient for analyzing the two-level (Section 4) and multilevel (Section 5) algorithms. Theorems 4.1 and 5.4 are the central results of this paper. We present an application of these methods to inverse problems of parabolic type in Section 6. The article concludes with a section on numerical results. Appendix B contains a brief notation summary.

## 2. Notation and problem formulation

Let $\mathcal{X} = \mathcal{Y} = L^2(\Omega)$, where $\Omega \subset \mathbb{R}^d$ is an open set ($d \geq 1$). Throughout this paper we shall denote by $H^m(\Omega), H_0^m(\Omega)$ ($m \in \mathbb{N}$) the standard Sobolev spaces,

and by $\|\cdot\| = \|\cdot\|_{L^2(\Omega)}$ and $\|\cdot\|_m$ (resp. $|\cdot|_m$) the corresponding norms (resp. semi-norms); furthermore $\langle\cdot,\cdot\rangle$ (resp. $\langle\cdot,\cdot\rangle_m$) will be the standard inner product in $L^2(\Omega)$ (resp. $H^m(\Omega)$). Let $\widetilde{H}^{-m}(\Omega)$ be the dual (with respect to the $L^2$-inner product) of $H^m(\Omega) \cap H_0^1(\Omega)$ for $m > 0$. For $T \in L(V_1, V_2)$ we denote the *operator norm* of $T$ by

$$\|T\|_{V_1,V_2} = \sup_{u \in V_1 \setminus \{0\}} \frac{\|Tu\|}{\|u\|} \quad \text{and} \quad \|T\|_V \overset{\text{def}}{=} \|T\|_{V,V}.$$

In the absence of any subscript, $\|T\|$ denotes $\|T\|_{L^2(\Omega)}$. We consider a set of approximating spaces $(\mathcal{V}_h)_{h \in I}$ with

(2.1) $$I = \{h_{\max}/2^i : i \in \mathbb{N}\}$$

that have the nesting property

(2.2) $$\mathcal{V}_{2h} \subset \mathcal{V}_h \subset H_0^1(\Omega), \quad \forall h \in I \setminus \{h_{\max}\}.$$

Furthermore, we denote by $\pi_h$ the $L^2$-orthogonal projection onto $\mathcal{V}_h$. Discretizations $\mathcal{K}^h \in L(\mathcal{V}_h)$ of $\mathcal{K}$ give rise to the following discrete quadratic functional to be minimized:

(2.3) $$\mathcal{J}_\beta^h(u_h) = \frac{\beta^{-1}}{2}\|\mathcal{K}^h u - f_h\|^2 + \frac{1}{2}\|u_h\|^2 ,$$

where $f_h = \pi_h(f) \in \mathcal{V}_h$. Throughout this paper it is assumed that the operators $\mathcal{K}$ and $\mathcal{K}^h$ together with their adjoints satisfy

**Condition 2.1.** There exists a number $p > 0$ (the approximation order) and constants $C_1 = C_1(p, \|\mathcal{K}\|, \Omega)$ and $C_2 = C_2(p, \Omega)$ such that for all $h \in I$ the following hold:

[a] stability:

(2.4) $$\|\mathcal{K}^h u\| \le C_1\|u\| \text{ and } \|(\mathcal{K}^h)^* u\| \le C_1\|u\|, \quad \forall u \in \mathcal{V}_h ;$$

[b] smoothed approximation:

(2.5) $$\|\mathcal{K}u - \mathcal{K}^h u\| \le C_1 h^p \|u\| \text{ and } \|\mathcal{K}^* u - (\mathcal{K}^h)^* u\| \le C_1 h^p \|u\|, \quad \forall u \in \mathcal{V}_h ;$$

[c] negative-index norm approximation of the identity by the projection:

(2.6) $$\|(I - \pi_h)u\|_{\widetilde{H}^{-p}(\Omega)} \le C_2 h^p \|u\|, \quad \forall u \in L^2(\Omega) ;$$

[d] smoothing:

(2.7) $$\|\mathcal{K}u\| \le C_1 \|u\|_{\widetilde{H}^{-p}} \text{ and } \|\mathcal{K}^* u\| \le C_1 \|u\|_{\widetilde{H}^{-p}}, \quad \forall u \in L^2(\Omega).$$

*Remark* 2.2. Condition 2.1 implies that there is a $C = C(p, \|\mathcal{K}\|, \Omega)$ such that

(2.8) $$\|Hu\|_{L^2} \le C \|u\|_{\widetilde{H}^{-p}}, \quad \forall u \in L^2(\Omega) ,$$

and

(2.9) $$\|\mathcal{K}(I - \pi_h)u\| \le C h^p \|u\| \text{ and } \|\mathcal{K}^*(I - \pi_h)u\| \le C h^p \|u\|, \quad \forall u \in L^2(\Omega).$$

In order to minimize $\mathcal{J}_\beta^h$ we need to invert the discrete version of $H_\beta$, namely

(2.10) $$H_\beta^h \overset{\text{def}}{=} I + \beta^{-1} H^h, \quad \text{with} \quad H^h \overset{\text{def}}{=} (\mathcal{K}^h)^* \mathcal{K}^h.$$

**Lemma 2.3.** *The following approximation property holds:*

(2.11) $$\|\pi_h(H^h - H)u\| \le C h^p \|u\|, \quad \forall u \in \mathcal{V}_h,$$

*for some constant $C = C(p, \|K\|, \Omega)$.*

*Proof.* For $u \in \mathcal{V}_h$ we have

$$
\begin{aligned}
\left| \langle \pi_h (H^h - H) u, u \rangle \right| &= \left| \|\mathcal{K}^h u\|^2 - \|\mathcal{K}u\|^2 \right| \\
&\leq \|\mathcal{K}^h u - \mathcal{K}u\| \cdot \left( \|\mathcal{K}^h u\| + \|\mathcal{K}u\| \right) \\
&\leq C(C + \|\mathcal{K}\|) \, h^p \|u\|^2 \ ,
\end{aligned}
$$

and (2.11) follows from the symmetry of $\pi_h(H^h - H) \in L(\mathcal{V}_h)$.  $\square$

It follows from the definition of $H_\beta$, that for all $u \in L^2(\Omega)$,

$$
(2.12) \qquad \|u\|^2 \leq \langle H_\beta u, u \rangle \leq \left( 1 + \frac{C}{\beta} \right) \|u\|^2 \ ,
$$

with $C = C(p, \|\mathcal{K}\|, \Omega)$, and similar estimates hold for the discrete Hessian. This implies that the quadratic $\mathcal{J}_\beta$ is positive definite and has a unique minimizer given by

$$
(2.13) \qquad u^{\min} = u_\beta^{\min} = \beta^{-1} (H_\beta)^{-1} \, \mathcal{K}^* f.
$$

Similarly the minimizer of the discrete quadratic is

$$
(2.14) \qquad u_h^{\min} = u_{h,\beta}^{\min} = \beta^{-1} (H_\beta^h)^{-1} \left( \mathcal{K}^h \right)^* f_h.
$$

We refer the reader to [15, 16] for results concerning convergence of $u_\beta^{\min}$ or to $\mathcal{K}^\dagger f$. The next result shows that $u_h^{\min}$ approximates $u^{\min}$ to optimal order in the $L^2$-norm and plays a role in the analysis of the multilevel method.

**Theorem 2.4.** *Assume that Condition 2.1 holds. Then there exists a constant $C = C(p, \|\mathcal{K}\|, \Omega)$ such that for $h \leq h_0(\beta, p, \|\mathcal{K}\|, \Omega)$ we have the following stability and error estimates:*

$$
(2.15) \qquad \|u_h^{\min}\| \leq C \left( \|u^{\min}\| + \beta^{-1} h^p \|f\| \right) \ ,
$$

$$
(2.16) \qquad \|u_h^{\min} - u^{\min}\| \leq C \frac{h^p}{\beta} \left( \|f\| + \|u^{\min}\| \right).
$$

*Proof.* Denote by $e_h = u_h^{\min} - u^{\min}$. We have

$$
H_\beta e_h = \beta^{-1} \left\{ \left( (\mathcal{K}^h)^* - \mathcal{K}^* \right) f_h - \mathcal{K}^* (I - \pi_h) f + (H - H^h) u_h^{\min} \right\}.
$$

Therefore

$$
\begin{aligned}
\beta \|e_h\|^2 \ &\overset{(2.12)}{\leq}\ \beta \langle H_\beta e_h, e_h \rangle \\
&= \left\langle \left( (\mathcal{K}^h)^* - \mathcal{K}^* \right) f_h - \mathcal{K}^* (I - \pi_h) f + (H - H^h) u_h^{\min}, e_h \right\rangle \\
&\overset{(2.5),\,(2.9)}{\leq}\ Ch^p \|f\| \, \|e_h\| + \left\langle (H - H^h) u_h^{\min}, e_h \right\rangle \\
&\overset{(I - \pi_h) e_h \perp \mathcal{V}_h}{=}\ Ch^p \|f\| \, \|e_h\| + \left\langle \pi_h (H - H^h) u_h^{\min}, \pi_h e_h \right\rangle \\
&\qquad + \left\langle H u_h^{\min}, (I - \pi_h) e_h \right\rangle \\
&\overset{(2.11)}{\leq}\ Ch^p \|e_h\| (\|f\| + \|u_h^{\min}\|) + \left\langle \mathcal{K} u_h^{\min}, \mathcal{K}(I - \pi_h) e_h \right\rangle \\
&\overset{(2.9)}{\leq}\ Ch^p \|e_h\| (\|f\| + \|u_h^{\min}\|) \ ,
\end{aligned}
$$

with $C = C(p, \|\mathcal{K}\|, \Omega)$. Hence it follows that

$$
(2.17) \qquad \|u_h^{\min} - u^{\min}\| \leq C \frac{h^p}{\beta} \left( \|f\| + \|u_h^{\min}\| \right).
$$

Since

$$\|u_h^{\min}\| \leq \|u^{\min}\| + \|u_h^{\min} - u^{\min}\| \leq \|u^{\min}\| + C\frac{h^p}{\beta}\left(\|f\| + \|u_h^{\min}\|\right) ,$$

we obtain (2.15) and (2.16), for $h$ small enough. $\qquad\qquad\qquad\qquad\qquad\square$

## 3. The spectral distance

In this section we define a scale-independent distance between operators with positive definite symmetric part, and we study its relevant properties. This *spectral distance* is a measure of spectral equivalence between two operators and introduces a convenient framework for the multilevel analysis in Section 5.

Throughout this section $(\mathcal{X}, \langle \cdot, \cdot \rangle)$ is a real, finite-dimensional Hilbert space. As usual, $\|u\| = \sqrt{\langle u, u \rangle}$ is the Hilbert-space norm of $u \in \mathcal{X}$. Denote by

$$L_+(\mathcal{X}) = \{T \in L(\mathcal{X}, \mathcal{X}) : \langle Tu, u \rangle > 0, \quad \forall u \in \mathcal{X} \setminus \{0\}\}$$

the set of operators with positive definite symmetric part. All operators in this section are assumed to be in $L_+(\mathcal{X})$. If $A \in L_+(\mathcal{X})$ is symmetric, we write $\|u\|_A \stackrel{\text{def}}{=} \sqrt{\langle Au, u \rangle} = \|A^{\frac{1}{2}}u\|$. Our object of study is the preconditioned Richardson iteration

(3.1) $$x_{n+1} = x_n + M(b - Hx_n)$$

leading to the solution $x^*$ of the equation

$$Hx = b ,$$

where $b \in \mathcal{X}$, $H \in L_+(\mathcal{X})$. Later we will specialize to $H$ being symmetric (we think of $H$ being Hessian $H_\beta^h$ from (2.10)), but we allow the preconditioner $M$ to be nonsymmetric, for reasons explained in Section 4.2. The results we prove for $H$ symmetric (essentially Theorem 3.12 and Corollary 3.13) apply to the reverse situation as well, with the preconditioner $M$ being symmetric and the operator $H$ nonsymmetric, a situation that has been studied beginning with [10, 14, 17]. Also, some of the techniques in this section are rooted in the aforementioned papers. It is well known that the error $e_n = x_n - x^*$ satisfies

$$e_n = (I - MH)^n e_0.$$

If $H$ is symmetric, then

$$e_n = H^{-\frac{1}{2}}(I - H^{\frac{1}{2}}MH^{\frac{1}{2}})^n H^{\frac{1}{2}}e_0$$

and

(3.2) $$\|e_n\|_H \leq \|I - H^{\frac{1}{2}}MH^{\frac{1}{2}}\|^n \cdot \|e_0\|_H .$$

Therefore the quantity $\|I - MH\|$ (or $\|I - H^{\frac{1}{2}}MH^{\frac{1}{2}}\|$, if $H$ is symmetric) is an upper bound for the convergence rate of (3.1). Another quality-measure for the preconditioner $M$, especially useful when neither $M$ nor $H$ are symmetric, is the spectral radius $\rho(I - MH)$. Although all the above quantities measure, in spirit, how far $M^{-1}$ is from $H$, none of them is a distance function in a strict mathematical sense. For technical reasons that will become clear in Section 5 we prefer to assess the quality of $M$ by using an actual (scale-free) distance function to measure how far $M^{-1}$ and $H$ are from each other.

We denote the complexification of $\mathcal{X}$ by

$$\mathcal{X}^{\mathbb{C}} = \mathcal{X} \otimes_{\mathbb{R}} \mathbb{C} = \{u + \mathbf{i}\,v \ : \ u, v \in \mathcal{X}\} ,$$

and we consider the natural extension of the inner product on $\mathcal{X}$ to a Hermitian product on the complex vector space $\mathcal{X}^{\mathbb{C}}$:

$$\langle u_1 + \mathbf{i}\, v_1, u_2 + \mathbf{i}\, v_2 \rangle = \langle u_1, u_2 \rangle + \langle v_1, v_2 \rangle + \mathbf{i}\, (\langle v_1, u_2 \rangle - \langle u_1, v_2 \rangle) \; ;$$

for $T \in L(\mathcal{X})$ define the *complexification* of $T$ to be $T^{\mathbb{C}} \in L(\mathcal{X}^{\mathbb{C}})$ (the space of $\mathbb{C}$-linear maps), where

$$T^{\mathbb{C}}(u + \mathbf{i}\, v) = T(u) + \mathbf{i}\, T(v).$$

We will drop the superscript $^{\mathbb{C}}$ whenever there is no potential for confusion. Furthermore, we denote by $\mathcal{B}_r(z_0)$ the open disc $\{z \in \mathbb{C} : |z - z_0| < r\}$.

**Definition 3.1.** Let $T_1, T_2 \in L_+(\mathcal{X})$. We define the *spectral distance* between $T_1$ and $T_2$ to be

$$
\begin{aligned}
d_{\mathcal{X}}(T_1, T_2) &= \sup_{w \in \mathcal{X}^{\mathbb{C}} \setminus \{0\}} \left| \ln \frac{\langle T_1^{\mathbb{C}} w, w \rangle}{\langle T_2^{\mathbb{C}} w, w \rangle} \right| \\
&= \sup_{(u,v) \in \mathcal{X} \times \mathcal{X} \setminus \{(0,0)\}} \left| \ln \frac{\langle T_1 u, u \rangle + \langle T_1 v, v \rangle + \mathbf{i}\, (\langle T_1 v, u \rangle - \langle T_1 u, v \rangle)}{\langle T_2 u, u \rangle + \langle T_2 v, v \rangle + \mathbf{i}\, (\langle T_2 v, u \rangle - \langle T_2 u, v \rangle)} \right| ,
\end{aligned}
$$

where ln is the branch of the logarithm corresponding to $\mathbb{C} \setminus (-\infty, 0]$.

We should point out that, for $T_1, T_2 \in L_+(\mathcal{X})$, the quotients $\langle T_1^{\mathbb{C}} w, w \rangle / \langle T_2^{\mathbb{C}} w, w \rangle$ do not lie in $(-\infty, 0]$. Also note that, if $T = T_s + T_a$ with $T_s$ (resp. $T_a$) being the symmetric (resp. antisymmetric) part of $T$, then

$$(3.3) \qquad \langle T^{\mathbb{C}}(u + \mathbf{i}\, v), u + \mathbf{i}\, v \rangle = \langle T_s u, u \rangle + \langle T_s v, v \rangle + 2\mathbf{i}\, \langle T_a v, u \rangle .$$

The polarization identity

$$
\begin{aligned}
\langle Tu, v \rangle &= \frac{1}{4}(\langle T(u+v), u+v \rangle - \langle T(u-v), u-v \rangle \\
&\quad + \mathbf{i}\, \langle T(u + \mathbf{i}\, v), u + \mathbf{i}\, v \rangle - \mathbf{i}\, \langle T(u - \mathbf{i}\, v), u - \mathbf{i}\, v \rangle)
\end{aligned}
$$

implies that $d_{\mathcal{X}}(T_1, T_2) = 0$ if and only if $T_1 = T_2$. We leave as an exercise to the reader the verification of the symmetry and triangle inequality (see also [12]).

The following elementary (but nontrivial) inequalities prove useful in evaluating the spectral distance.

**Lemma 3.2.** *If $\alpha \in (0,1)$ and $z \in \mathcal{B}_\alpha(1)$, then*

$$(3.4) \qquad \frac{\ln(1+\alpha)}{\alpha}|1 - z| \leq |\ln z| \leq \frac{|\ln(1-\alpha)|}{\alpha}|1 - z|.$$

*For $|\ln z| \leq \delta$ we have*

$$(3.5) \qquad \frac{1 - e^{-\delta}}{\delta}|\ln z| \leq |1 - z| \leq \frac{e^\delta - 1}{\delta}|\ln z|.$$

*Proof.* The modulus of the analytic function $f : \overline{\mathcal{B}_\alpha(1)} \to \mathbb{C}$ defined by $f(z) = \ln z/(1 - z)$ attains its extreme values on the boundary $\partial \mathcal{B}_\alpha(1)$, as $f$ has no zeros in $\mathcal{B}_\alpha(1)$. On the circle $\mathcal{C}_\alpha(1) = \{\zeta : |1 - \zeta| = \alpha\}$ we have $|f(z)| = \alpha^{-1}|\ln z|$; hence the extreme values of $|f(z)|$ are attained at the same points as the extremes of $|\ln z|^2$. The problem is thus reduced to showing that the maximum of $|\ln z|^2$ on the circle $\mathcal{C}_\alpha(1)$ is attained at $\zeta = 1 - \alpha$, and the minimum at $\zeta = 1 + \alpha$. We leave this calculus problem as an exercise (hint: Draw tangents from the origin to $\mathcal{C}_\alpha(1)$, and parametrize the two resulting arches using polar coordinates around 0. Use the

angle as the free variable. The functions to be analyzed are $g_{1,2}(\theta) = \ln(\rho_{1,2}(\theta))^2 + \theta^2$, with $\rho_{1,2}(\theta) = \cos\theta \pm \sqrt{\alpha^2 - \sin^2\theta}$). For (3.5) we proceed similarly.     $\square$

*Remark* 3.3. We should note the following short one-line proof of the fact that the global maximum of $|\ln z|$ on $\mathcal{C}_\alpha(1)$ is located at $z = 1 - \alpha$:

$$|\ln z| \quad = \quad \left| \sum_{n=1}^{\infty} (-1)^{n-1} \frac{(z-1)^n}{n} \right| \leq \sum_{n=1}^{\infty} \frac{\alpha^n}{n} = |\ln(1 - \alpha)|.$$

However, we found no such simple argument for showing that the global **minimum** is at $z = 1 + \alpha$.

The following result shows that, even if neither $M$ nor $H$ is symmetric, the spectral radius of $(I - MH)$ is controlled by the spectral distance between $M^{-1}$ and $H$.

**Lemma 3.4.** *Let $M, H \in L_+(\mathcal{X})$ such that $d_{\mathcal{X}}(M^{-1}, H) \leq \delta$. Then*

$$(3.6) \qquad\qquad \rho(I - MH) \leq \frac{e^\delta - 1}{\delta} d_{\mathcal{X}}(M^{-1}, H).$$

*In particular, if $d_{\mathcal{X}}(M^{-1}, H) < \ln 2$, then $\rho(I - MH) < 1$.*

*Proof.* Let $\lambda \in \sigma(I - MH)$. Then there exists a unit vector $u \in \mathcal{X}^{\mathbb{C}}$ such that $(I - MH)u = \lambda u$, which further implies that

$$\lambda = 1 - \frac{\langle Hu, u \rangle}{\langle M^{-1}u, u \rangle}.$$

Hence

$$\begin{aligned} \rho(I - MH) \quad &\leq \quad \sup\{|1 - z| \ : \ z = \langle Hu, u \rangle / \langle M^{-1}u, u \rangle \text{ for some } u \in \mathcal{X}^{\mathbb{C}} \setminus \{0\}\} \\ &\overset{(3.5)}{\leq} \quad \frac{e^\delta - 1}{\delta} d_{\mathcal{X}}(M^{-1}, H). \end{aligned}$$

$\square$

*Remark* 3.5. In general, the converse of (3.6) does not hold. For example, let $H = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, and $M = I$. Then $\rho(I - MH) = 0$, but $d_{\mathcal{X}}(M^{-1}, H) \neq 0$.

A concept related to the spectral distance is the *numerical range* or *field of values* of an operator $T$, which is defined in [18] as

$$W(T) = \left\{ \langle T^{\mathbb{C}}u, u \rangle : u \in \mathcal{X}^{\mathbb{C}}, \|u\| = 1 \right\} \ ;$$

the *numerical radius* is

$$w(T) = \sup\{|\lambda| : \lambda \in W(T)\}.$$

We recall two results from [18]:

**Theorem 3.6** (Theorem 1.3-1 in [18]). *For any $T \in L(\mathcal{X})$,*

$$(3.7) \qquad\qquad w(T) \leq \|T\| \leq 2w(T).$$

This is basically saying that $w(\cdot)$, which is a norm, is equivalent to the operator norm $\| \cdot \|$. However, for symmetric operators we have

$$\tag{3.8} \| T \| = w(T) \,,$$

since $W(T) = \mathrm{co}(\sigma(T))$ ($\mathrm{co}(A)$ denotes the *convex hull* of the set $A$). The next result is about mapping of the numerical radius under a power function:

**Theorem 3.7** (Theorem 2.1-1 in [18]). *For any $T \in L(\mathcal{X})$,*

$$\tag{3.9} w(T^n) \leq (w(T))^n \,, n = 1, 2, 3, \ldots.$$

By analogy we define the *joint numerical range* to be

$$\tag{3.10} W(T_1, T_2) = \left\{ \frac{\langle T_1 w, w \rangle}{\langle T_2 w, w \rangle} : w \in \mathcal{X}^{\mathbb{C}} \setminus \{0\} \right\}.$$

Note that $W(T, I) = W(T)$. With this definition the spectral distance becomes

$$\tag{3.11} d_{\mathcal{X}}(T_1, T_2) = \sup\{| \ln(\lambda)| : \lambda \in W(T_1, T_2)\}.$$

*Remark* 3.8. If $T_2$ is symmetric, then

$$\tag{3.12} W(T_1, T_2) = W(T_2^{-\frac{1}{2}} T_1 T_2^{-\frac{1}{2}}).$$

This is easily seen from the change of variable $w = T_2^{-\frac{1}{2}} v$ in (3.10).

**Lemma 3.9.** *If $T_1, T_2 \in L_+(\mathcal{X})$ are symmetric, then*

$$\tag{3.13} d_{\mathcal{X}}(T_1, T_2) = \sup_{u \in \mathcal{X} \setminus \{0\}} \left| \ln \frac{\langle T_1 u, u \rangle}{\langle T_2 u, u \rangle} \right|$$

*and*

$$\tag{3.14} \frac{\| I - T_1^{\frac{1}{2}} T_2^{-1} T_1^{\frac{1}{2}} \|}{d_{\mathcal{X}}(T_1, T_2)} \to 1 \quad \text{as } d_{\mathcal{X}}(T_1, T_2) \to 0.$$

*Proof.* Let $T_{12} = T_2^{-\frac{1}{2}} T_1 T_2^{-\frac{1}{2}}$. The symmetry of $T_{12}$ implies

$$W(T_1, T_2) \overset{(3.12)}{=} W(T_{12}) \overset{(3.3)}{=} \left\{ \frac{\langle T_{12} u, u \rangle + \langle T_{12} v, v \rangle}{\langle u, u \rangle + \langle v, v \rangle} : (u, v) \neq (0, 0) \right\}.$$

Since $W(T_{12})$ is convex (Theorem 1.1-2 in [18]) we have

$$\left\{ \frac{\langle T_{12} u, u \rangle + \langle T_{12} v, v \rangle}{\langle u, u \rangle + \langle v, v \rangle} : (u, v) \neq (0, 0) \right\} \subseteq \left\{ \frac{\langle T_{12} u, u \rangle}{\langle u, u \rangle} : u \neq 0 \right\}.$$

The reverse inclusion being evident, we get

$$W(T_1, T_2) = \left\{ \frac{\langle T_{12} u, u \rangle}{\langle u, u \rangle} : u \neq 0 \right\} \overset{u' = T_2^{-\frac{1}{2}} u}{=} \left\{ \frac{\langle T_1 u', u' \rangle}{\langle T_2 u', u' \rangle} : u' \neq 0 \right\},$$

and (3.13) follows by (3.11). Moreover, since all operators involved are symmetric,

$$W(T_{12}) = \mathrm{co}\left(\sigma(T_{12})\right) = [\lambda_{\min}, \lambda_{\max}] \,,$$

where $\lambda_{\min}$ (resp. $\lambda_{\max}$) is the smallest (resp. largest) eigenvalue of $T_{12}$. The fact that $d_{\mathcal{X}}(T_1, T_2) \to 0$ translates into $\lambda_{\min}, \lambda_{\max} \to 1$; hence

$$\frac{\| I - T_{12} \|}{d_{\mathcal{X}}(T_1, T_2)} = \frac{\max(|1 - \lambda_{\min}|, |1 - \lambda_{\max}|)}{\max(|\ln \lambda_{\min}|, |\ln \lambda_{\max}|)} \to 1. \qquad \square$$

The following lemma is a restating in terms of the spectral distance of equivalent results from [10, 14, 17].

**Lemma 3.10.** *For symmetric operators $T_1, T_2 \in L_+(\mathcal{X})$ we have*

$$(3.15) \qquad\qquad d_{\mathcal{X}}(T_1, T_2) = d_{\mathcal{X}}((T_1)^{-1}, (T_2)^{-1}).$$

*Proof.* As in the proof of Lemma 3.9 we have

$$
\begin{aligned}
W(T_1, T_2) \quad &= \quad \left\{ \frac{\langle T_{12} u, u \rangle}{\langle u, u \rangle} \; : \; u \neq 0 \right\} \overset{v = T_{12}^{\frac{1}{2}} u}{=} \left\{ \frac{\langle v, v \rangle}{\langle T_{12}^{-1} v, v \rangle} \; : \; v \neq 0 \right\} \\[2mm]
&= \quad \left\{ \frac{\langle v, v \rangle}{\left\langle T_2^{\frac{1}{2}} T_1^{-1} T_2^{\frac{1}{2}} v, v \right\rangle} \; : \; v \neq 0 \right\} \\[2mm]
&\overset{w = T_2^{\frac{1}{2}} v}{=} \left\{ \frac{\langle T_2^{-1} w, w \rangle}{\langle T_1^{-1} w, w \rangle} \; : \; w \neq 0 \right\} = W(T_2^{-1}, T_1^{-1}) \;,
\end{aligned}
$$

and the conclusion follows from (3.11). $\qquad\qquad\square$

The remainder of this section is devoted to showing that if $H$ is symmetric, and $H^{\frac{1}{2}} M H^{\frac{1}{2}}$ is close to $I$, then the spectral distance can replace $\| I - H^{\frac{1}{2}} M H^{\frac{1}{2}} \|$ in the estimate (3.2), at the expense of a slightly increased constant on the right-hand side (see (3.27)); this fact is not surprising, since the values of $d_{\mathcal{X}}(T_1, T_2)$ and $\| I - T_1^{-\frac{1}{2}} T_2 T_1^{-\frac{1}{2}} \|$ become asymptotically close as $d_{\mathcal{X}}(T_1, T_2) \to 0$. It will be convenient to regard the Richardson iteration (3.1) as an **iteration of preconditioners**. More precisely, we think of (3.1) as a sequence of **single** iterations with "improved" preconditioners:

$$(3.16) \qquad\qquad x_n = x_0 + M_n(b - Hx_0) \;,$$

with $M_0 = 0$ and $M_n$ recursively defined by $M_{n+1} = M + M_n - MHM_n$. A simple calculation shows that

$$(3.17) \qquad M_n = H^{-1} - (I - MH)^n H^{-1} = H^{-\frac{1}{2}} \left( I - \left( I - H^{\frac{1}{2}} M H^{\frac{1}{2}} \right)^n \right) H^{-\frac{1}{2}}.$$

In particular $M_1 = M$; moreover, $M_2 = 2M - MHM$ is the first Newton iterate with initial guess $M_1$ of the operator-function $X \mapsto X^{-1} - H$, as shown below.

*Remark* 3.11. Define $G : L_+(\mathcal{X}) \to L(\mathcal{X})$ by $G(X) = X^{-1} - H$. Since $G'(X)U = -X^{-1} U X^{-1}$ (see [9]), it follows that $(G'(X))^{-1} U = -XUX$; the Newton iteration for solving $G(X) = 0$ is

$$
X_{i+1} \quad = \quad X_i - (G'(X_i))^{-1}(X_i^{-1} - H) = 2X_i - X_i H X_i.
$$

Hence the iteration is $X_{i+1} = \mathcal{N}_H(X_i)$ with

$$(3.18) \qquad\qquad \mathcal{N}_H(X) \overset{\text{def}}{=} 2X - XHX.$$

We will be using the operator $\mathcal{N}_H$ in designing a multilevel preconditioner in Section 5 in the following way: given an initial (and unacceptable) guess $X$ at $H^{-1}$, the next best guess is $\mathcal{N}_H(X)$. By (3.17) the first Richardson iterate with $\mathcal{N}_H(X)$ as a preconditioner instead of $M$, namely the value $x_0 + \mathcal{N}_H(X)(b - Hx_0)$, is equal to $x_2$ from (3.1). Therefore applying $\mathcal{N}_H(X)$ amounts to performing two Richardson iterations with $M$ as a preconditioner.

**Theorem 3.12.** *Let $M, H \in L_+(\mathcal{X})$ and $M_n$ be defined as in (3.17). If $H$ is symmetric and $d_{\mathcal{X}}(M, H^{-1}) < \ln 2$, then $M_n \in L_+(\mathcal{X})$ and*

$$(3.19) \qquad d_{\mathcal{X}}(M_n, H^{-1}) \leq g_n\left(d_{\mathcal{X}}(M, H^{-1})\right) ,$$

*where $g_n(x) = |\ln(1 - (e^x - 1)^n)| = x^n + O(x^{n+1})$; hence*

$$(3.20) \qquad \lim_{n \to \infty} M_n = H^{-1} \quad \text{in the metric } d_{\mathcal{X}}.$$

*Proof.* Let $\alpha = d_{\mathcal{X}}(M, H^{-1})$. By (3.11) and the fact that

$$\max\{|e^w - 1| : w \in \mathcal{B}_\alpha(0)\} = e^\alpha - 1,$$

it follows that

$$W(H^{\frac{1}{2}} M H^{\frac{1}{2}}) = W(M, H^{-1}) \subset \mathcal{B}_{e^\alpha - 1}(1) ;$$

therefore

$$(3.21) \qquad w\left(I - H^{\frac{1}{2}} M H^{\frac{1}{2}}\right) \leq e^\alpha - 1.$$

Equation (3.17) implies that

$$(3.22) \qquad \left(I - H^{\frac{1}{2}} M_n H^{\frac{1}{2}}\right) = \left(I - H^{\frac{1}{2}} M H^{\frac{1}{2}}\right)^n ;$$

hence by (3.9) and (3.21) we get

$$(3.23) \qquad w\left(I - H^{\frac{1}{2}} M_n H^{\frac{1}{2}}\right) \leq \gamma \overset{\text{def}}{=} (e^\alpha - 1)^n.$$

Since $\max\{|\ln(z)| : z \in \mathcal{B}_\gamma(1)\} = |\ln(1 - \gamma)|$ (here we used $\gamma < 1$, which follows from $\alpha < \ln 2$ and also implies that $M_n$ is positive definite), we get

$$(3.24) \qquad d_{\mathcal{X}}(M_n, H^{-1}) \leq |\ln(1 - (e^\alpha - 1)^n)| ,$$

which concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

**Corollary 3.13.** *Under the conditions of Theorem 3.12 the sequence $x_n$ defined by the simple iteration (3.1) converges to the solution $x^*$, and we have the estimate*

$$(3.25) \qquad \|e_n\|_H \leq 2\, g\left(d_{\mathcal{X}}(M_n, H^{-1})\right) \cdot \|e_0\|_H ,$$

*where $g(x) = e^x - 1 = x + o(x)$.*

*Proof.* The relations (3.2), (3.23) and (3.7) imply

$$(3.26) \qquad \|e_n\|_H \leq 2\, w\left(I - H^{\frac{1}{2}} M_n H^{\frac{1}{2}}\right) \cdot \|e_0\|_H .$$

Similarly to the proof of Theorem 3.12 we have

$$w\left(I - H^{\frac{1}{2}} M_n H^{\frac{1}{2}}\right) \leq g\left(d_{\mathcal{X}}(M_n, H^{-1})\right) ,$$

with $g$ as in the hypothesis. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

*Remark* 3.14. In light of the asymptotic behavior of the functions $g_n$ and $g$ above, we rewrite (3.25) in the following way: for every $C > 2$ there exists $\delta(C) > 0$ such that, if $d_{\mathcal{X}}(M, H^{-1}) < \delta(C)$, then

$$(3.27) \qquad \|e_n\|_H \leq C\left(d_{\mathcal{X}}(M, H^{-1})\right)^n \cdot \|e_0\|_H .$$

## 4. Two-level preconditioners for the discrete Hessian

In this section we construct and analyze two-level preconditioners for the discrete Hessian $H_\beta^h$ defined in (2.10). Given that $H_\beta^h$ is symmetric, it is natural to seek preconditioners that are symmetric as well (Section 4.1). However, the results in Section 3 show that symmetry of the preconditioner is not essential for the Richardson iteration to converge. In Section 4.2 we investigate the possibility of using slightly nonsymmetric operators as preconditioners for $H_\beta^h$. Nonsymmetry may occur by allowing restriction operators other than exact $L^2$-projections. The main results in this section are Theorems 4.1 and 4.3, which essentially state that the spectral distance between the constructed preconditioners and the inverse of the Hessian is $O(h^p/\beta)$.

We should point out that the results presented in this section, although similar to the ones in [30, 20] are obtained in a different context. While the discrete operators in the aforementioned papers are obtained by orthogonal projections of their continuous versions on the finite-dimensional spaces under consideration, our discrete operators arise from natural finite element discretizations, and the only assumptions made are related to convergence and stability. Naturally the results in Section 4.1 have counterparts in [30, 20], but the techniques used here are different.

4.1. **A symmetric preconditioner.** For the remainder of this article we consider on $\mathcal{V}_h$ the Hilbert-space structure inherited from $L^2(\Omega)$. Let $\mathcal{W}_{2h}$ be the orthogonal complement of $\mathcal{V}_{2h}$ in $\mathcal{V}_h$, and let $\rho_{2h}$ be the orthogonal projection onto $\mathcal{W}_{2h}$, so that $\pi_{2h} + \rho_{2h} = I$. We shall refer to $\mathcal{V}_h$ as the **fine** space in relation to $\mathcal{V}_{2h}$, which will be the **coarse** space. We consider the splitting

$$H_\beta^h \;=\; \underbrace{\pi_{2h} H_\beta^h \pi_{2h}}_{A_1} + \underbrace{\rho_{2h} H_\beta^h \pi_{2h}}_{A_2} + \underbrace{\pi_{2h} H_\beta^h \rho_{2h}}_{A_3} + \underbrace{\rho_{2h} H_\beta^h \rho_{2h}}_{A_4}.$$

Standard Fourier analysis on uniform grids shows that $\pi_{2h}$ extracts the low-frequency components from a function in $\mathcal{V}_h$, while $\rho_{2h}$ extracts the high-frequency components. For $h \ll 1$ it is expected that $H^h$ inherits the smoothing character of $H = \mathcal{K}^*\mathcal{K}$, therefore annihilates high-frequency components, i.e., $H^h \rho_{2h} \approx 0$. Similarly, $\rho_{2h} H^h \approx 0$, since this operation amounts to extracting high-frequency components from a "smoothened" function. This suggests that the cross terms $A_2, A_3$ are negligible (since $A_2 \approx \rho_{2h}\pi_{2h} = 0$, $A_3 \approx \pi_{2h}\rho_{2h} = 0$), and that $A_4 \approx \rho_{2h}$ (see [11] for more details). On "smooth" components both $H^h$ and $H^{2h}$ approximate $H$ well, since all eigenvectors of $H$ corresponding to large eigenvalues are well represented on fine grids (the high-frequency eigenvectors of $H$ correspond to negligible eigenvalues); hence $A_1 \approx H_\beta^{2h}$. The resulting approximation is

$$H_\beta^h \;\approx\; M_\beta^h \;\overset{\text{def}}{=}\; H_\beta^{2h}\pi_{2h} + (I - \pi_{2h}).$$

Since $\pi_{2h}$ is a projection, the inverse $(M_\beta^h)^{-1}$ is given explicitly by

$$(4.1) \qquad L_\beta^h = L_\beta^h(\pi_{2h}) \;\overset{\text{def}}{=}\; (M_\beta^h)^{-1} = (H_\beta^{2h})^{-1}\pi_{2h} + (I - \pi_{2h}).$$

We propose $L_\beta^h \in L(\mathcal{V}_h)$ as a two-level preconditioner for $H_\beta^h$. Due to the similarity with the additive Schwarz construction from domain decomposition, and in accordance with [30], we denote it here as the TLAS (**T**wo-**L**evel **A**dditive **S**chwarz) preconditioner. We would like to remark that $M_\beta^h$ and $L_\beta^h$ are symmetric. It should be noted that the projector $I - \pi_{2h}$ removes high-frequency components from the

approximation error; therefore it acts as a smoother. However, unlike the case of the classical multigrid, neither is the solution a fixed point for this smoother, nor does it make sense to apply the smoother more than once.

The remainder of this section is devoted to proving the following theorem, which gives an estimate of the distance between $L_\beta^h$ and $(H_\beta^h)^{-1}$. We write $d_h$ instead of $d_{\mathcal{V}_h}$.

**Theorem 4.1.** *Assume that the operators $\mathcal{K}, \mathcal{K}^h$ satisfy Condition 2.1. Then there exists a constant $C = C(p, \|\mathcal{K}\|, \Omega)$ such that for $h \leq h_0$ with $h_0 = h_0(\beta, p, \|\mathcal{K}\|, \Omega)$ we have $M_\beta^h \in L_+(\mathcal{V}_h)$, and*

$$(4.2) \qquad\qquad d_h(H_\beta^h, M_\beta^h) \leq C \frac{h^p}{\beta}.$$

Before we proceed to the proof of the theorem we state the corollary which legitimizes the use of $L_\beta^h$ as a preconditioner for the Hessian.

**Corollary 4.2.** *Under the hypotheses of Theorem 4.1 there exists a constant $C = C(p, \|\mathcal{K}\|, \Omega)$ such that $L_\beta^h \in L_+(\mathcal{V}_h)$, and*

$$(4.3) \qquad\qquad d_h\left(L_\beta^h, \left(H_\beta^h\right)^{-1}\right) \leq C \frac{h^p}{\beta} \ ,$$

*for $h \leq h_0$ with $h_0 = h_0(\beta, p, \|\mathcal{K}\|, \Omega)$.*

*Proof.* The result follows directly from Theorem 4.1 and Lemma 3.10.    □

We would like to remark that, for fixed $\beta$, the quality of the preconditioner $L_\beta^h$ increases with $h \to 0$. This is different from standard multigrid preconditioning for elliptic problems, where the goal would be to show that the right-hand side of (4.3) is independent of $h$. It should not be surprising that (4.3) is such an optimistic estimate. This is related to the adverse eigenstructure of the eigenvalue-frequency correlation for the operator $H_\beta^h$. By increasing resolution we only add high-frequency eigenfunctions whose eigenvalues are close to 1. On the other hand, the low-frequency eigenvectors (with higher energy) are increasingly well approximating the continuous (smooth) eigenvectors of $H_\beta$, and therefore they are increasingly well approximated by corresponding eigenvectors of $H_\beta^{2h}$. Hence the fine-space operator is **increasingly well** approximated by the identity on the high-frequency part, and by the coarse-space operator on the low-frequency part. We return to the proof of Theorem 4.1:

*Proof.* We have the following relation on $\mathcal{V}_h$:

$$\beta(M_\beta^h - H_\beta^h) \;=\; \underbrace{\pi_{2h}\left(H^{2h} - H\right)\pi_{2h}}_{A} + \underbrace{\left(\pi_{2h}H\pi_{2h} - \pi_h H\pi_h\right)}_{B} + \underbrace{\left(\pi_h H\pi_h - H^h\right)}_{C}.$$

Note that all operators in the above sum are symmetric in $L(\mathcal{V}_h)$. We analyze the products $\langle Au, u\rangle, \langle Bu, u\rangle, \langle Cu, u\rangle$. For $u \in \mathcal{V}_h$,

$$|\langle Au, u\rangle| = \left|\left\langle \left(H^{2h} - \pi_{2h}H\right)\pi_{2h}u, \pi_{2h}u\right\rangle\right| \overset{(2.11)}{\leq} C(2h)^p \|\pi_{2h}u\|^2 \leq C(2h)^p \|u\|^2.$$

Similarly, we obtain $|\langle Cu, u\rangle| \leq Ch^p \|u\|^2$. For the middle term we have

$$
\begin{aligned}
|\langle Bu, u\rangle| \;&=\; |\langle(\pi_{2h}H\pi_{2h} - \pi_h H\pi_h)u, u\rangle| \overset{u\in\mathcal{V}_h}{=} |\langle H\pi_{2h}u, \pi_{2h}u\rangle - \langle Hu, u\rangle| \\
&=\; \left| \|\mathcal{K}\pi_{2h}u\|^2 - \|\mathcal{K}u\|^2 \right| = \|\mathcal{K}(I - \pi_{2h})u\| \cdot (\|\mathcal{K}\pi_{2h}u\| + \|\mathcal{K}u\|) \\
&\overset{(2.9)}{\leq}\; Ch^p \|u\|^2 \;,
\end{aligned}
$$

with $C = C(p, \|K\|, \Omega)$. Putting the above estimates together we get

$$
(4.4) \qquad \left| \frac{\left\langle M_\beta^h u, u\right\rangle}{\left\langle H_\beta^h u, u\right\rangle} - 1 \right| \;\leq\; C\frac{h^p}{\beta} \cdot \frac{\|u\|^2}{\beta^{-1}\left\langle H^h u, u\right\rangle + \|u\|^2} \;\leq\; C\frac{h^p}{\beta}.
$$

Assuming $C\beta^{-1}h_0^p = \alpha < 1$ and $0 < h \leq h_0$ we obtain $M_\beta^h \in L_+(\mathcal{V}_h)$, and

$$
\sup_{u\in\mathcal{V}_h\backslash\{0\}} \left| \ln \frac{\left\langle M_\beta^h u, u\right\rangle}{\left\langle H_\beta^h u, u\right\rangle} \right| \overset{(3.4)}{\leq} \frac{|\ln(1-\alpha)|}{\alpha} \sup_{u\in\mathcal{V}_h\backslash\{0\}} \left| \frac{\left\langle M_\beta^h u, u\right\rangle}{\left\langle H_\beta^h u, u\right\rangle} - 1 \right|
$$

$$
\leq \quad C\frac{|\ln(1-\alpha)|}{\alpha} \cdot \frac{h^p}{\beta}.
$$

The conclusion follows from the symmetry of $M_\beta^h$ and $H_\beta^h$ and from Lemma 3.9. $\quad\square$

### 4.2. Nonsymmetric preconditioners.

The application of the preconditioner $L_\beta^h(\pi_{2h})$ from Section 4.1 requires computing exact (or very close to exact) $L^2$-projections. In this section we prove a result similar to Theorem 4.1 that applies to a larger class of preconditioners, obtained by replacing in (4.1) the orthogonal projection $\pi_{2h}$ with a more general restriction operator $R_{2h}^h \in L(\mathcal{V}_h, \mathcal{V}_{2h})$ that verifies Condition 2.1[c]. A local restriction operator that satisfies a similar, but slightly weaker condition is defined in [5]. In Section 6.3 we show how to construct local restriction operators with the required properties.

**Theorem 4.3.** *Assume that the operators $\mathcal{K}, \mathcal{K}^h$ satisfy Condition 2.1, and that the restriction operator $R_{2h}^h \in L(\mathcal{V}_h, \mathcal{V}_{2h})$ satisfies for all $\beta \in (0,1]$ and $h \in I$,*

[e] *stability:*

$$
(4.5) \qquad \|R_{2h}^h u\| \leq C\|u\|, \quad \forall u \in \mathcal{V}_h ;
$$

[f] *approximation of the identity in the negative-index norm:*

$$
(4.6) \qquad \|(I - R_{2h}^h)u\|_{\widetilde{H}^{-p}(\Omega)} \leq Ch^p \|u\|, \quad \forall u \in \mathcal{V}_h ,
$$

*with $C = C(p, \|\mathcal{K}\|, \Omega)$. If we denote by*

$$
(4.7) \qquad L_\beta^h = L_\beta^h(R_{2h}^h) \overset{\text{def}}{=} (H_\beta^{2h})^{-1}R_{2h}^h + (I - R_{2h}^h) ,
$$

*then for $h \leq h_0$ with $h_0 = h_0(\beta, p, \|\mathcal{K}\|, \Omega)$ we have $L_\beta^h \in L_+(\mathcal{V}_h)$, and*

$$
(4.8) \qquad d_h\left((H_\beta^h)^{-1}, L_\beta^h\right) \leq C\frac{h^p}{\beta^2} ,
$$

*for some constant $C = C(p, \|\mathcal{K}\|, \Omega)$.*

Note that this result, though applying to a larger class of restriction operators, is weaker than Theorem 4.1, due to the extra power of $\beta$ in the denominator of the right-hand side of (4.8). We first prove

**Lemma 4.4.** *Under the hypotheses of Theorem 4.3 there exists a constant* $C = C(p, \|\mathcal{K}\|, \Omega)$ *such that, for all* $u \in \mathcal{V}_h$ *the following hold:*

   **(i)** *smoothing properties of the discrete Hessian:*

$$(4.9) \qquad \max\left(\|H^h(I - R^h_{2h})u\|,\ \|H^{2h}(\pi_{2h} - R^h_{2h}u)\|\right) \leq Ch^p\|u\|\ ;$$

   **(ii)** *smoothing properties of the inverse discrete Hessian:*

$$(4.10) \qquad \|((H^{2h}_\beta)^{-1} - I)(\pi_{2h} - R^h_{2h})u\| \leq C\beta^{-1}h^p\|u\|.$$

*Proof.* Throughout this proof $C$ will denote a generic constant depending only on $p, \|\mathcal{K}\|$ and $\Omega$.

   **(i)** Conditions (2.6) and (4.6) imply that for $u \in \mathcal{V}_h$,

$$(4.11) \qquad \left\|(\pi_{2h} - R^h_{2h})u\right\|_{\widetilde{H}^{-p}} \leq Ch^p\|u\|\,.$$

For $u \in \mathcal{V}_h$ we have

$$\left\|\left(H^{2h} - \pi_{2h}H\right)(\pi_{2h} - R^h_{2h})u\right\| \overset{(2.11)}{\leq} C(2h)^p\|(\pi_{2h} - R^h_{2h})u\| \overset{(4.5)}{\leq} Ch^p\|u\|\ ;$$

hence

$$\begin{aligned}
\left\|H^{2h}(\pi_{2h} - R^h_{2h})u\right\| &\leq \left\|\pi_{2h}H(\pi_{2h} - R^h_{2h})u\right\| + C(2h)^p\|u\| \\
&\overset{(2.8)}{\leq} C\left\|(\pi_{2h} - R^h_{2h})u\right\|_{\widetilde{H}^{-p}(\Omega)} + Ch^p\|u\| \\
&\overset{(4.11)}{\leq} Ch^p\|u\|\,.
\end{aligned}$$

The inequality

$$\left\|H^h(I - R^h_{2h})u\right\| \leq Ch^p\|u\|\,, \quad u \in \mathcal{V}_h\ ,$$

follows along the same lines, and (4.9) is proved.

   **(ii)** Since $\langle H^{2h}u, u\rangle \geq 0$, it follows that $\langle H^{2h}_\beta u, u\rangle \geq \|u\|^2$, $\forall u \in \mathcal{V}_{2h}$. Hence

$$(4.12) \qquad \|\left(H^{2h}_\beta\right)^{-1} u\| \leq \|u\|\,, \quad \text{for } u \in \mathcal{V}_{2h}.$$

For $u \in \mathcal{V}_h$,

$$\begin{aligned}
\|((H^{2h}_\beta)^{-1} - I)(\pi_{2h} - R^h_{2h})u\| &= \beta^{-1}\|(H^{2h}_\beta)^{-1}H^{2h}(\pi_{2h} - R^h_{2h})u\| \\
&\overset{(4.9),\ (4.12)}{\leq} C\beta^{-1}h^p.
\end{aligned}$$

$\square$

We now proceed to the proof of Theorem 4.3.

*Proof.* Recall that in Corollary 4.2 it was shown that

$$d_h\left(\left(M^h_\beta\right)^{-1}, \left(H^h_\beta\right)^{-1}\right) \leq C\beta^{-1}h^p\ ,$$

where

$$\left(M^h_\beta\right)^{-1} = \left(H^{2h}_\beta\right)^{-1}\pi_{2h} + (I - \pi_{2h}).$$

Since $\beta \leq 1$, it suffices to show that

$$(4.13) \qquad d_h\left(\left(M^h_\beta\right)^{-1}, L^h_\beta\right) \leq C\beta^{-2}h^p.$$

We rewrite (4.7) as

$$L_\beta^h = I + \left( \left( H_\beta^{2h} \right)^{-1} - I \right) R_{2h}^h \ ,$$

and we start by evaluating the antisymmetric part of $L_\beta^h$:

$$\begin{aligned}
&\left\langle L_\beta^h u, v \right\rangle - \left\langle L_\beta^h v, u \right\rangle \\
&= \ \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) R_{2h}^h u, v \right\rangle - \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) R_{2h}^h v, u \right\rangle \\
&= \ \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) R_{2h}^h u, \pi_{2h} v \right\rangle - \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) R_{2h}^h v, \pi_{2h} u \right\rangle \\
&= \ \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) (R_{2h}^h - \pi_{2h}) u, \pi_{2h} v \right\rangle \\
&\quad - \left\langle \left( \left( H_\beta^{2h} \right)^{-1} - I \right) (R_{2h}^h - \pi_{2h}) v, \pi_{2h} u \right\rangle .
\end{aligned}$$

For the last equality we have added and subtracted $\left\langle ((H_\beta^{2h})^{-1} - I) \pi_{2h} u, \pi_{2h} v \right\rangle$ and used the symmetry of $((H_\beta^{2h})^{-1} - I)$ on $\mathcal{V}_{2h}$. By Lemma 4.4 (ii) we get

$$(4.14) \qquad \left| \left\langle L_\beta^h u, v \right\rangle - \left\langle L_\beta^h v, u \right\rangle \right| \leq C \beta^{-1} h^p \left\| u \right\| \cdot \left\| v \right\| .$$

Note that $M_\beta^h = H_\beta^{2h} \pi_{2h} + (I - \pi_{2h}) = \beta^{-1} H^{2h} \pi_{2h} + I$; hence

$$(4.15) \qquad \left\langle u, u \right\rangle \leq \left\langle M_\beta^h u, u \right\rangle \leq C \beta^{-1} \left\langle u, u \right\rangle , \quad \text{for } u \in \mathcal{V}_h \ ,$$

with $C = C(p, \|K\|, \Omega)$. This implies

$$(4.16) \qquad C^{-1} \beta \left\langle u, u \right\rangle \leq \left\langle \left( M_\beta^h \right)^{-1} u, u \right\rangle \leq \left\langle u, u \right\rangle , \quad \text{for } u \in \mathcal{V}_h.$$

It follows that the projection of the set $W(L_\beta^h, (M_\beta^h)^{-1})$ onto the imaginary axis (see (3.3)) satisfies

$$\frac{\left| \left\langle L_\beta^h u, v \right\rangle - \left\langle L_\beta^h v, u \right\rangle \right|}{\left\langle (M_\beta^h)^{-1} u, u \right\rangle + \left\langle (M_\beta^h)^{-1} v, v \right\rangle} \quad \overset{(4.14),\ (4.16)}{\leq} \quad C \frac{h^p}{\beta^2} \cdot \frac{\left\| u \right\| \cdot \left\| v \right\|}{\left\| u \right\|^2 + \left\| v \right\|^2} \ ;$$

we obtain therefore

$$(4.17) \qquad \sup \left\{ \left| \Im w \right| \ : \ w \in W(L_\beta^h, (M_\beta^h)^{-1}) \right\} \leq C \beta^{-2} h^p \ ,$$

where $\Im w$ is the imaginary part of a complex number $w$. We now turn our attention to the projection of $W(L_\beta^h, (M_\beta^h)^{-1})$ onto the real axis. Since $(M_\beta^h)^{-1}$ is symmetric, this amounts to computing the joint numerical range of the symmetric part of $L_\beta^h$ and $(M_\beta^h)^{-1}$. By Lemma 3.9 we need to evaluate

$$(4.18) \qquad \frac{\left\langle L_\beta^h u, u \right\rangle}{\left\langle (M_\beta^h)^{-1} u, u \right\rangle} = 1 + \frac{\left\langle (L_\beta^h - (M_\beta^h)^{-1}) u, u \right\rangle}{\left\langle (M_\beta^h)^{-1} u, u \right\rangle}$$

for $u \in \mathcal{V}_h \setminus \{0\}$. A simple calculation shows that

$$L_\beta^h - (M_\beta^h)^{-1} \ = \ ((H_\beta^{2h})^{-1} - I)(R_{2h}^h - \pi_{2h}) \ ;$$

therefore, by Lemma 4.4 (ii) it follows that

$$(4.19) \qquad \left| \left\langle \left( L_\beta^h - \left( M_\beta^h \right)^{-1} \right) u, u \right\rangle \right| \leq C \beta^{-1} h^p \left\| u \right\|^2 \ ;$$

this, together with (4.18) and (4.16), implies

$$(4.20) \qquad \left\{ \Re w \ : \ w \in W(L_\beta^h, (M_\beta^h)^{-1}) \right\} \subset [1 - C \beta^{-2} h^p, \ 1 + C \beta^{-2} h^p] \ ,$$

where $\Re w$ is the real part of a complex number $w$. Finally from (4.17) and (4.20) it follows that

$$(4.21) \qquad\qquad W\left(L_\beta^h, \left(M_\beta^h\right)^{-1}\right) \subset \mathcal{B}_{C\beta^{-2}h^p}(1).$$

In particular this implies, possibly by further reducing $h_0$ by a factor of two, that $L_\beta^h \in L_+(\mathcal{V}_h)$. Lemma 3.2 implies that (4.13) holds for $h \leq h_0(\beta, p, \|\mathcal{K}\|, \Omega)$. $\qquad\square$

## 5. A multilevel preconditioner

In this section we define a multilevel preconditioner $K_\beta^h$ that is of comparable quality with the two-level preconditioner $L_\beta^h$ defined in Section 4. By analogy we will denote the preconditioner by MLAS (**M**ulti-**L**evel **A**dditive **S**chwarz).

5.1. **Design and work estimates.** Before constructing $K_\beta^h$ we want to point out that one natural way to extend $L_\beta^h$ to a multilevel preconditioner $G_\beta^h$ is by replacing $(H_\beta^{2h})^{-1}$ in (4.1) with the coarse-space preconditioner; thus we define the preconditioner recursively:

$$(5.1) \qquad\qquad G_\beta^h \overset{\text{def}}{=} G_\beta^{2h}\pi_{2h} + (I - \pi_{2h}).$$

At the coarsest level $h_0$ we use the conjugate gradient method as an "almost direct" solver. It is immediate that $G_\beta^h \in L_+(\mathcal{V}_h)$. This strategy yields a multilevel preconditioner with a $V$-cycle structure, which is very similar to the one defined in [30]. In Theorem 5.2 we prove that the distance between $G_\beta^h$ and $(H_\beta^h)^{-1}$ only depends on the coarsest level resolution, which is consistent with the results in [30]. This preconditioner is suboptimal, and this fact should not be surprising: the multilevel $G_\beta^h$ is in fact only a two-level operator, in the sense that we would obtain the same $G_\beta^h$ if we had used only the finest and the coarsest levels. For example, if we use three levels corresponding to resolutions $h$, $2h$, and $4h$, then

$$
\begin{aligned}
G_\beta^h &= G_\beta^{2h}\pi_{2h} + (I - \pi_{2h}) = \left(G_\beta^{4h}\pi_{4h} + (I - \pi_{4h})\right)\pi_{2h} + (I - \pi_{2h}) \\
&= G_\beta^{4h}\pi_{4h} + (I - \pi_{4h}), \quad \text{since } \pi_{4h}\pi_{2h} = \pi_{4h}.
\end{aligned}
$$

Also we should point out that, when applying $G_\beta^h$, residuals at intermediate levels are never computed. Thus the quality of the multilevel preconditioner $G_\beta^h$, while not degrading, also does not improve with increasing **fine-level** resolution. As a result, for example, the number of $G_\beta^h$-preconditioned conjugate gradient iterations will be constant if $h_0$ is kept fixed and $h \to 0$. Our goal is to design a multilevel preconditioner whose quality **improves** with increasing fine-level resolution. This would imply that asymptotically (as $h \to 0$) we may only need one preconditioned iteration at the finest level.

We adopt a different strategy in defining an improved preconditioner $K_\beta^h$: in addition to replacing $(H_\beta^{2h})^{-1}$ with the coarse-space preconditioner, like before, we perform one Newton iteration for the operator-function $X \to X^{-1} - H_\beta^h$, as explained in Remark 3.11. More precisely, if we denote by $\mathcal{G}^h : L(\mathcal{V}_{2h}) \to L(\mathcal{V}_h)$ the affine operator-function

$$(5.2) \qquad\qquad \mathcal{G}^h(T) = T\pi_{2h} + (I - \pi_{2h}) \, ,$$

then we define $K_\beta^h$ recursively using (3.18) by

$$(5.3) \qquad K_\beta^h \overset{\text{def}}{=} \mathcal{N}_{H_\beta^h}\left(\mathcal{G}^h(K_\beta^{2h})\right) = 2\mathcal{G}^h(K_\beta^{2h}) - \mathcal{G}^h(K_\beta^{2h}) \cdot H_\beta^h \cdot \mathcal{G}^h(K_\beta^{2h}).$$

With this notation the symmetric preconditioner from Section 4.1 is written as

$$L_\beta^h = \mathcal{G}^h((H_\beta^{2h})^{-1}).$$

We should note that $\mathcal{G}^h$ is symmetry-preserving; therefore $K_\beta^h$ is symmetric if $K_\beta^{h_0}$ is so. The application of $K_\beta^h$ requires one residual computation at the finest level (see Algorithm MLAS); moreover, performing one $K_\beta^h$-preconditioned Richardson iteration simply means to complete two $\mathcal{G}^h(K_\beta^{2h})$-preconditioned Richardson iterations, as explained in Section 3. This implies that there will be two calls to the coarser-space procedure, therefore $K_\beta^h$ has a W-cycle structure. We would like to remark that the calling sequence *pre-smoothing - restriction - error correction - interpolation - post-smoothing* from the classical multigrid is not appropriate for Algorithm MLAS. Here restriction applies to the right-hand side as a whole as opposed to the smoothed residual, as in the classical multigrid, and smoothing occurs again between the error-correction steps.

**Algorithm MLAS.** $(b \mapsto K_\beta^h b)$

    (1) if $h = h_0$ then return $K_\beta^{h_0} b$        // direct solve
    (2) else
    (3)     $b_c \leftarrow \pi_{2h} b$                 // restriction
    (4)     $u_1 \leftarrow b - b_c + K_\beta^{2h} b_c$       // smoothing and error correction
    (5)     $r \leftarrow b - H_\beta^h u_1$           // residual computation
    (6)     $r_c \leftarrow \pi_{2h} r$                // restriction
    (7)     $u_2 \leftarrow u_1 + K_\beta^{2h} r_c + r - r_c$    // smoothing and error correction
    (8)     return $u_2$

Let $h_i = 2^{-i} h_0$, $i = 0, 1, \ldots, l$. If we denote by $W(i)$ (resp. $R(i)$) the work needed to apply $K_\beta^{h_i}$ (resp. $H_\beta^{h_i}$), then

$$(5.4) \qquad\qquad W(i) \approx R(i) + 2W(i-1).$$

We have assumed that the cost of restriction and interpolation is negligible compared to that of a residual computation. Indeed this is the case when $H_\beta^{h_i}$ is represented by a dense matrix, and restriction by a sparse matrix, or when the direct problem (i.e., applying $\mathcal{K}$) is a space-time process and restriction operates only on the spatial variables (see Section 6). If $R(i-1) \le 2^{-g} R(i)$ for some $g > 0$, then (5.4) implies

$$(5.5) \qquad\qquad W(l) \le (1 - 2^{-g})^{-1} R(l) + 2^l W(0).$$

It is shown in [16] that the unpreconditioned conjugate gradient takes a level-independent number of iterations $N_{cg}$ to solve the exponentially ill-posed problems ($N_{cg}$ still depends on $\beta$). This results in $W(0) = C N_{cg} R(0)$, with $C = O(1)$ being a universal constant. Hence we have the estimate

$$(5.6) \qquad\qquad W(l) \le \left( (1 - 2^{-g})^{-1} + C 2^{l(1-g)} N_{cg} \right) R(l).$$

The number $g$ is related to the dimension of the problem to be solved. For example, if the forward operator $\mathcal{K}$ is a three-dimensional space-time operator, and we use a time-stepping procedure with the same convergence order as the spatial discretization, then $g = 4$. With $l = 3$ levels, (5.6) results in $W(3) \approx (16/15 + 0.002 \cdot C \cdot N_{cg}) R(3)$. In practice we have observed $N_{cg}$ ranging from 20 to

200 (see [1]). Hence for this numerical example the work for applying the three-level preconditioner is a small multiple of the work for a residual computation.

In practice, algorithm MLAS should be modified so that, at the finest level, no residual computation is performed inside the preconditioner; that is, at the finest level MLAS should return the value $u_1$ on line (4). Without this modification, $K_\beta^h$ would approximate $(H_\beta^h)^{-1}$ to an excessively high order, as seen in the discussion at the end of Section 5.2. However, we prefer to use the current definition of $K_\beta^h$ for the sake of keeping the presentation of the error estimates in the following section free of branching.

5.2. **Error estimates for the multilevel preconditioner.** The main goal of this section is to estimate the distance

$$(5.7) \qquad e_h = d_h(K_\beta^h, (H_\beta^h)^{-1}).$$

This quantity lies at the basis of the residual and error estimates, as pointed out in Section 3. We will also estimate $d_h(G_\beta^h, (H_\beta^h)^{-1})$.

**Lemma 5.1.** *Let* $T_1, T_2 \in L_+(\mathcal{V}_{2h})$ *be symmetric. Then*

$$(5.8) \qquad d_h\left(\mathcal{G}^h(T_1), \mathcal{G}^h(T_2)\right) \leq d_{2h}\left(T_1, T_2\right).$$

*Proof.* A simple calculation shows that for any $a, b > 0$ the function $g(x) = |\ln(a+x) - \ln(b+x)|$ is nonincreasing on $[0, \infty]$. Therefore

$$(5.9) \qquad \left|\ln \frac{a+x}{b+x}\right| \leq \left|\ln \frac{a}{b}\right|.$$

Since $(I - \pi_{2h})$ is positive semidefinite, we obtain for $u \in \mathcal{V}_h$ with $\pi_{2h} u \neq 0$,

$$\left|\ln \frac{\langle (T_1 \pi_{2h} + I - \pi_{2h})u, u \rangle}{\langle (T_2 \pi_{2h} + I - \pi_{2h})u, u \rangle}\right| \quad = \quad \left|\ln \frac{\langle T_1 \pi_{2h} u, \pi_{2h} u \rangle + \langle (I - \pi_{2h})u, u \rangle}{\langle T_2 \pi_{2h} u, \pi_{2h} u \rangle + \langle (I - \pi_{2h})u, u \rangle}\right|$$

$$\overset{(5.9)}{\leq} \quad \left|\ln \frac{\langle T_1 \pi_{2h} u, \pi_{2h} u \rangle}{\langle T_2 \pi_{2h} u, \pi_{2h} u \rangle}\right| \leq d_{2h}(T_1, T_2).$$

The conclusion follows after passing to the supremum over all $u \in \mathcal{V}_h$ with $\pi_{2h} u \neq 0$ in the inequality above. $\qquad \square$

**Theorem 5.2** (V-cycle estimates). *Under the hypotheses and in the notation of Theorem 4.1 there exists a constant* $C = C(p, \|\mathcal{K}\|, \Omega)$ *such that*

$$(5.10) \qquad d_h(G_\beta^h, (H_\beta^h)^{-1}) \leq C \frac{h_0^p}{\beta}.$$

*Proof.* Assume that $h = h_0 2^{-l}$ with $h_0$ small enough so that Theorem 4.1 applies. Then

$$d_h(G_\beta^h, (H_\beta^h)^{-1}) \quad \overset{\text{triangle ineq.}}{\leq} \quad d_h\left(G_\beta^h, L_\beta^h\right) + d_h(L_\beta^h, (H_\beta^h)^{-1})$$

$$\overset{\text{Cor. 4.2, (5.8)}}{\leq} \quad d_{2h}(G_\beta^{2h}, (H_\beta^{2h})^{-1}) + C\beta^{-1} h^p.$$

We apply the above recursively to obtain

$$(5.11) \qquad d_h(G_\beta^h, (H_\beta^h)^{-1}) \leq C\beta^{-1} h_0^p \sum_{i=0}^{l} 2^{-ip} < (1 - 2^{-p})^{-1} C\beta^{-1} h_0^p,$$

which proves (5.10). $\qquad \square$

**Lemma 5.3.** *Let $(e_i)_{i \geq 0}$ and $(a_i)_{i \geq 0}$ be positive numbers satisfying the recursive inequality*

$$(5.12) \qquad\qquad e_{i+1} \leq C(e_i + a_{i+1})^2$$

*and*

$$(5.13) \qquad\qquad a_{i+1} \leq a_i \leq f^{-1} a_{i+1}$$

*for some $0 < f < 1$. If $a_0 \leq \frac{f}{4C}$ and if $e_0 \leq 4Ca_0^2$, then*

$$(5.14) \qquad\qquad e_i \leq 4Ca_i^2, \quad \forall i > 0.$$

*Proof.* We proceed by induction: the base case is in the hypothesis, and if we assume $e_i \leq 4Ca_i^2$, then

$$
\begin{aligned}
e_{i+1} &\leq C(e_i + a_{i+1})^2 \leq C(4Ca_i^2 + a_{i+1})^2 \\
&\leq C(4Cf^{-1}a_i a_{i+1} + a_{i+1})^2 = C(4Cf^{-1}a_i + 1)^2 a_{i+1}^2 \\
&\leq C\underbrace{(4Cf^{-1}a_0 + 1)^2}_{\leq 4} a_{i+1}^2 \leq 4Ca_{i+1}^2 ,
\end{aligned}
$$

which concludes the proof. $\qquad\qquad\square$

**Theorem 5.4** (MLAS error estimates). *Assume that Condition 2.1 holds, and that*

$$(5.15) \qquad\qquad h_0^p \leq \frac{\min\left(0.1, 2^{-(p+3)}\right) \cdot \beta}{C} ,$$

*where $C$ is the constant from Theorem 4.1. If $e_{h_0} \leq 8\left(C\beta^{-1}h_0^p\right)^2$ (if we solve the system almost exactly on the coarsest grid this is automatically satisfied), then $K_\beta^h \in L_+(\mathcal{V}_h)$ and*

$$(5.16) \qquad\qquad d_h(K_\beta^h, (H_\beta^h)^{-1}) \leq 8\, C \frac{h^{2p}}{\beta^2} , \quad \forall h \in I.$$

*Proof.* Theorem 3.12 implies that, given operators $M, H \in L_+(\mathcal{V}_h)$ with $d_h(M, H^{-1}) < 0.4$, it follows that $M \in L_+(\mathcal{V}_h)$ and

$$(5.17) \qquad\qquad d_h(\mathcal{N}_H(M), H^{-1}) \leq 2\, d_h(M, H^{-1})^2 ,$$

because $|\ln(1 - (e^x - 1)^2)| \leq 2\,x^2$ for $x \in [0, 0.4]$. If $e_{2h} \leq 0.2$ and $C\beta^{-1}h^p \leq 0.1$, then

$$
\begin{aligned}
d_h(\mathcal{G}^h(K_\beta^{2h}), (H_\beta^h)^{-1}) &\leq d_h\left(\mathcal{G}^h(K_\beta^{2h}), L_\beta^h\right) + d_h(L_\beta^h, (H_\beta^h)^{-1}) \\
&= e_{2h} + C\beta^{-1}h^p \leq 0.3 ;
\end{aligned}
$$

hence by (5.17) we obtain

$$(5.18) \quad e_h \leq 2\left[d_h(\mathcal{G}^h(K_\beta^{2h}), (H_\beta^h)^{-1})\right]^2 \leq 2(e_{2h} + C\beta^{-1}h^p)^2 \leq 2 \cdot 0.3^2 < 0.2.$$

Assume that $h = 2^{-l}h_0$, and let $h_i = 2^{-i}h_0$. Denote by $e_i = e_{h_i}$ and by $a_i = C\beta^{-1}h_i^p$. An inductive argument implies that $e_i \leq 0.2$ for all $i$, provided that it holds for $e_0$, and that $C\beta^{-1}h_0^p \leq 0.1$. It follows that $(a_i)_{i \geq 0}$ satisfies (5.12) and (5.13) with $f = 2^{-p}$. The other hypotheses of the theorem are tailored to fit the corresponding hypotheses in Lemma 5.3, which implies that

$$(5.19) \qquad\qquad e_i \leq 8\left(C\beta^{-1}h_i^p\right)^2, \text{ for all } i \geq 0. \qquad\qquad\square$$

The previous result shows that the quality of the multilevel preconditioner improves with increasing resolution; moreover the estimate is optimal with respect to $h$. We should remark that a proper comparison of the multilevel preconditioner with the two-level preconditioner from Section 4 requires weighing $K_\beta^h = \mathcal{N}_{H_\beta^h}(K_\beta^{2h}\pi_{2h} + (I - \pi_{2h}))$ against $\mathcal{N}_{H_\beta^h}(L_\beta^h) = \mathcal{N}_{H_\beta^h}((H_\beta^{2h})^{-1}\pi_{2h} + (I - \pi_{2h}))$ (rather than $L_\beta^h$). Theorem 5.4, Corollary 4.2 and Theorem 3.12 show that the distances from $K_\beta^h$ and $\mathcal{N}_{H_\beta^h}(L_\beta^h)$ to $(H_\beta^h)^{-1}$ are asymptotically comparable, an assertion that is confirmed by the numerical results in Section 7.

In the regime of its applicability, namely when $h^p \ll \beta$, the estimate (5.16) seems to provide a solution mechanism that is better than needed: if the iterative process for finding $u_h^{\min}$ starts with an $O(1)$ close initial guess, then one $K_\beta^h$-preconditioned iteration will give an $O\left((\beta^{-1}h^p)^2\right)$ approximation to $u_h^{\min}$. This is better than the approximation of $u^{\min}$ by $u_h^{\min}$ as shown in (2.16), therefore unnecessary. As shown at the end of Section 5.1, $K_\beta^h$ should be implemented so that no finest-level residual computation should be performed inside the preconditioner. The multigrid preconditioner can be used in conjunction with a Krylov solver, and all finest-level Hessian-vector multiplications should be left to the iterations in the solver.

## 6. Application to the finite element Galerkin approximation for parabolic problems

In this section we identify a class of linear parabolic problems and their discretizations for which Condition 2.1 is satisfied; in particular it will follow that all results in Sections 4 and 5 apply. The verification amounts primarily to providing links to the corresponding results in the literature. In Section 6.2 we discuss a few results on regularization for inverse parabolic problems, and we devote Section 6.3 to the construction of local restriction operators.

6.1. **The inverse problem.** Let $\Omega \subset \mathbb{R}^d$ be an open set, with $d \geq 1$ an integer. We consider the following initial value problem:

$$
(6.1) \qquad \begin{cases} \partial_t u + A(t)u = 0 & \text{on } \Omega \times (0, \infty) \ , \\ u(x, t) = 0 & \text{on } \partial\Omega \times (0, \infty) \ , \\ u(x, 0) = u_0(x) & \text{for } x \in \Omega \ , \end{cases}
$$

where

$$
(6.2) \qquad A(t)u = -\sum_i \partial_i \left( \sum_j a_{ij}(x, t)\partial_j u + b_i(x, t)u \right) + c(x, t)u
$$

with $a_{ij}(x, t) = a_{ji}(x, t), b_i(x, t), c(x, t)$ being smooth functions with uniformly bounded derivatives of all orders on $\overline{\Omega} \times [0, \infty)$. We define the time-dependent bilinear form $a : (0, \infty) \times H_0^1 \times H_0^1 \to \mathbb{R}$ by

$$
(6.3) \qquad a(t; \phi, \psi) = \sum_{i,j} \langle a_{ij}\partial_j\phi, \partial_i\psi \rangle + \sum_i \langle b_i\phi, \partial_i\psi \rangle + \langle c\phi, \psi \rangle, \quad \text{for } \phi, \psi \in H_0^1.
$$

It is assumed that $a$ is coercive, i.e., there exists a constant $c_1 > 0$ independent of $t$ such that

$$
(6.4) \qquad a(t; \phi, \phi) \geq c_1 \|\phi\|_1^2, \quad \text{for } \phi \in H_0^1 \ ,
$$

and that the boundary $\partial\Omega$ is smooth enough for the following regularity condition to hold:

$$(6.5) \qquad \|\phi\|_2 \leq c_2 \|A(t)\phi\|, \quad \text{for } \phi \in H_0^1 \cap H^2.$$

This is verified for example in case $\Omega \subset \mathbb{R}^2$ is a convex polygonal domain and $a_{ij}(x,t) = \alpha(x,t)\delta_{ij}$ with $\alpha(x,t)$ a positive smooth function ($\delta$ is the Kronecker symbol), or if $\partial\Omega \in C^1$. For $u_0 \in L^2$ there exists a unique solution $u : (0,\infty) \to H_0^1$ to the weak formulation of (6.1), namely

$$(6.6) \qquad \begin{cases} \langle \partial_t u, \phi \rangle + a(t; u, \phi) = 0, & \text{for all } \phi \in H_0^1, \ t > 0, \\ \lim_{t \to 0} \|u(t) - u_0\| = 0. \end{cases}$$

(cf. [26] and [27]). We denote the time-$t$ solution operator by $\mathcal{S}(t)u_0 \overset{\text{def}}{=} u(\cdot, t)$.

Given a fixed $T > 0$ we define the operator $\mathcal{K} \in L(\mathcal{X})$ by

$$(6.7) \qquad \mathcal{K} \overset{\text{def}}{=} \mathcal{S}(T),$$

with $\mathcal{X} = L^2(\Omega)$. We restate Lemmas 2.3, 2.6 and equation (2.12) in [26] as

**Lemma 6.1.** *For $p = 0, 1, 2$ and $t > 0$,*

$$(6.8) \qquad \|\mathcal{S}(t)u_0\|_{H^p} \leq Ct^{-p/2}\|u_0\|, \quad \forall u_0 \in L^2(\Omega),$$

*where $C$ is independent of $t$ and $u_0$.*

For each fixed $t > 0$, the adjoint $(\mathcal{S}(t))^*$ is the time-$t$ solution operator of the equation obtained from (6.1) by switching the signs of the $b_i$'s, thus resulting in an equation of the same type as (6.1). Therefore Lemma 6.1 applies to $(\mathcal{S}(t))^*$ as well.

**Corollary 6.2.** *For $p = 0, 1, 2$ and $t > 0$,*

$$(6.9) \qquad \max(\|\mathcal{S}(t)u_0\|, \|(\mathcal{S}(t))^*u_0\|) \leq Ct^{-p/2}\|u_0\|_{\widetilde{H}^{-p}}, \quad \forall u_0 \in L^2(\Omega),$$

*where $C$ is independent of $t$ and $u_0$.*

*Proof.* For $u_0 \in L^2(\Omega)$ we have

$$\begin{aligned} \|\mathcal{S}(t)u_0\|^2 &= \langle u_0, (\mathcal{S}(t))^*\mathcal{S}(t)u_0 \rangle \leq \|u_0\|_{\widetilde{H}^{-p}} \cdot \|(\mathcal{S}(t))^*\mathcal{S}(t)u_0\|_{H^p} \\ &\overset{(6.8)}{\leq} Ct^{-p/2}\|u_0\|_{\widetilde{H}^{-p}} \cdot \|\mathcal{S}(t)u_0\|, \end{aligned}$$

which concludes the proof. $\qquad\square$

We conclude that the forward and adjoint operators $\mathcal{K}$ and $\mathcal{K}^*$ satisfy Condition 2.1[d]. We discretize the forward problem via the Galerkin method using continuous piecewise linear functions in space, and backward Euler in time. In order to verify Conditions 2.1[a] and [b] we resort to the literature on error estimates for fully discrete parabolic problems with irregular data [26, 27, 28]. For an extensive presentation of finite element methods for parabolic problems, see [31]. For simplicity we focus on the two-dimensional case ($d = 2$). Let $\mathcal{T}_{h_0}$ be a triangulation of the domain $\Omega$, and let $\mathcal{T}_{h/2}$ be defined inductively to be the Goursat refinement of $\mathcal{T}_h$ for all $h \in I = \{h_0/2^i : i \in \mathbb{N}\}$ (each triangle in $T \in \mathcal{T}_h$ is cut along the three lines obtained by joining the midpoints of its edges). Note that $(\mathcal{T}_h)_{h \in I}$ is a sequence of quasi-uniform triangulations. For $h \in I$ define

$$(6.10) \qquad \mathcal{V}_h = \{f \in \mathcal{C}(\overline{\Omega}) : \forall T \in \mathcal{T}_h \ f|_T \text{ linear, and } f|_{\partial\Omega} \equiv 0\}.$$

We have $\mathcal{V}_{h/2} \subset \mathcal{V}_h \subset H^1(\Omega)$. The following estimate holds for $v \in H_0^1 \cap H^2$:

$$(6.11) \qquad \inf_{\phi_h \in \mathcal{V}_h} \|v - \phi_h\|_j \le c_3 h^{2-j} \|v\|_2, \quad j \in \{0, 1\}$$

(cf. [7]). Let $t_m = mk$ with $m = 0, 1, \ldots, M_h$ so that $t_{M_h} = T$, with $k = k(h)$. The backward difference approximation $U^m$ to $u(t_m)$ is computed succesively by

$$(6.12) \qquad \begin{cases} \langle d_t U^{m+1}, \phi \rangle + a(t_{m+1}; U^{m+1}, \phi) = 0, \quad \forall \phi \in \mathcal{V}_h \ , \\ U^0 = \pi_h u_0 \ , \end{cases}$$

where $d_t U^{m+1} = k^{-1}(U^{m+1} - U^m)$. We will write for $u_0 \in L^2(\Omega)$,

$$(6.13) \qquad \mathcal{S}^h(t_m) u_0 \overset{\text{def}}{=} U^m, \ m = 1, 2, \ldots, M_h.$$

The following estimate holds with $C$ independent of $h, m$ (see [26, 27]):

$$(6.14) \qquad \|U^m - u(t_m)\| \le C \, t_m^{-1} \left( h^2 + k \right) \|u_0\|, \quad m = 1, 2, \ldots \ .$$

We will choose $k(h) = k_0 h^2$; therefore the discrete operator

$$(6.15) \qquad \mathcal{K}^h \overset{\text{def}}{=} \mathcal{S}^h(T) \in L(\mathcal{V}_h)$$

satisfies Conditions 2.1[**a**, **b**]. The adjoint $(\mathcal{S}^h(t_m))^*$ is the backward Euler Galerkin solution of the parabolic equation whose solution operator is $(\mathcal{S}(t))^*$ (under the assumption that at each time step the linear solve is performed exactly). Therefore $(\mathcal{K}^h)^*$ and $\mathcal{K}^*$ also satisfy Conditions 2.1[**a**, **b**]. Condition 2.1[**c**] is an easy consequence of the Bramble-Hilbert Lemma. Due to the fact that the time-stepping method is only first-order accurate in time, we would use four times less time steps for the coarser level; therefore a coarse-level residual computation is four times less costly than a fine-level residual computation (the exponent $g$ in Section 5.1 is $g = 4$). We have thus verified that all results in Sections 4 and 5 apply to the problem (1.1) with $\mathcal{K}, \mathcal{K}^h$ defined as in (6.7) and (6.15). In fact this particular model problem has been the driving force behind the development and analysis of the MLAS preconditioner.

6.2. **Some results on regularization for inverse parabolic problems.** In this section we present a few calculations that explain some of the qualitative behavior observed in the numerical solution of regularized inverse parabolic problems from Section 6.1. In Lemma 6.4 we show why the solution of the $L^2$-regularized inverse problem (1.1) with $\mathcal{K}$ defined by (6.7) has a square variation that is comparable to that of the "true" initial value if $\beta$ is chosen appropriately (see also Figure 1). This fact may seem surprising, given that $L^2$-regularization explicitly controls only size, and not derivatives. Example 6.5 is devoted to explaining another observed fact, namely that, when controlling the initial value for matching final-time measurements, recovery of **intermediate** states is of significantly better quality than that of the initial condition (see Figure 2 in Section 7). The immediate consequence is that, even if the recovered initial value may not be of acceptable quality (e.g. for localizing certain features), it still can be used successfully for improving predictions.

*Remark* 6.3. Let $\mathcal{X}, \mathcal{X}'$ be two Hilbert spaces such that $\mathcal{X}' \subseteq \mathcal{X}$ ($\mathcal{X}'$ does not have to be a Hilbert-subspace of $\mathcal{X}$). Furthermore, let $\mathcal{K} \in L(\mathcal{X})$ be such that

$\mathcal{K}(\mathcal{X}') \subseteq \mathcal{X}'$, and assume that $\mathcal{K}|_{\mathcal{X}'} \in L(\mathcal{X}')$ is compact. Then $\sigma(\mathcal{K}|_{\mathcal{X}'}) \setminus \{0\}$ only consists of a point-spectrum; therefore

$$(6.16) \qquad \sigma(\mathcal{K}|_{\mathcal{X}'}) \subseteq \sigma(\mathcal{K}) \cup \{0\}.$$

**Lemma 6.4.** *Assume $\mathcal{K}$ is given by (6.7) and that the "data" $f$ in (1.1) takes the form $f = \mathcal{K}u_0 + \varphi$, with $\varphi \in \mathcal{X}$; furthermore, assume that $u_0 \in H_0^1(\Omega)$. Then $H_\beta$ is continuously invertible in $H_0^1(\Omega)$ for all $\beta > 0$, and*

$$(6.17) \qquad \|u^{\min} - u_0\|_1 \quad \leq \quad \|(H_\beta)^{-1}\|_{H_0^1} \cdot (\|u_0\|_1 + CT^{-\frac{1}{2}} \beta^{-1}\|\varphi\|) ,$$

*where $C$ is the constant from (6.8).*

*Proof.* Let $\mathcal{X}' = H_0^1(\Omega)$. By Lemma 6.1 the operator $H = \mathcal{K}^* \cdot \mathcal{K}$ is compact both in $L(\mathcal{X})$ and $L(\mathcal{X}')$. Moreover, it is symmetric and positive definite in $\mathcal{X}$; hence $\sigma(H) \subset [0, \infty)$. Remark 6.3 implies that $\sigma(H|_{\mathcal{X}'}) \subset [0, \infty)$; therefore $H_\beta$ is invertible in $L(\mathcal{X}')$ for all $\beta > 0$. Since

$$(6.18) \qquad u^{\min} - u_0 = (H_\beta)^{-1} \left( -u_0 + \beta^{-1} \mathcal{K}^* \varphi \right) ,$$

the conclusion follows by (6.8). $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

It should be noted that, if $\mathcal{K}$ is the solution operator for the heat equation, then $H$ is symmetric with respect to the $H_0^1$ inner product as well; therefore $\|(H_\beta)^{-1}\|_{H_0^1} \leq 1$.

**Example 6.5.** Let $\mathcal{S}(t)$ be the solution operator for the one-dimensional heat equation on $\Omega = [0, 2\pi]$ with zero boundary conditions. We will show that the relative distance between the "exact" solution $\mathcal{S}(t)u_0$ and the recovered solution $\mathcal{S}(t)u^{\min}$ is decaying in time. For simplicity we restrict our attention to the case of unperturbed measurements, i.e. $f = \mathcal{S}(T)u_0$.

Let $\chi_k(x) = (2\pi)^{-1}e^{\mathbf{i}kx}$, $k \in \mathbb{Z}$, be the standard Fourier basis of $L^2(\Omega)$, and let $\hat{u} = \langle u, \chi_k \rangle$ be the Fourier coefficients of a function $u \in L^2(\Omega)$. We formulate our inverse problem in the space $\mathcal{X} = \{u \in L^2(\Omega) : \hat{u}_k + \hat{u}_{-k} = 0, \forall k \in \mathbb{Z}\}$ that is invariant under $\mathcal{S}(t)$. We prefer to express all quantities in terms of the basis $(\chi_k)_{k \in \mathbb{Z}}$ as opposed to an orthonormal basis of $\mathcal{X}$. With the formulation and notation from Lemma 6.4 we have

$$(6.19) \qquad (H_\beta)^{-1}u_0 = \sum_{k=-\infty}^{\infty} \frac{\hat{u}_k}{1 + \beta^{-1} e^{-2Tk^2}} \chi_k ,$$

where $u_0 = \sum_{k=-\infty}^{\infty} \hat{u}_k \chi_k$. For simplicity we assume that $\hat{u}_k \neq 0$ if $k \neq 0$ (by construction $\hat{u}_0 = 0$), and that $\beta \ll 1$. By (6.18),

$$(6.20) \qquad \frac{\|\mathcal{S}(t)(u_0 - u^{\min})\|^2}{\|\mathcal{S}(t)u_0\|^2} = \frac{\sum_{k=-\infty}^{\infty}(1 + \beta^{-1}e^{-2Tk^2})^{-2} e^{-2tk^2}|\hat{u}_k|^2}{\sum_{k=-\infty}^{\infty} e^{-2tk^2}|\hat{u}_k|^2}.$$

Since for moderate and large values of $k$ the eigenvalues $\mu_{k,\beta} = (1 + \beta^{-1}e^{-2Tk^2})^{-1}$ of $(H_\beta)^{-1}$ are close to 1, the error formula (6.18) implies that the corresponding components of $u_0$ are not well approximated by those of $u^{\min}$; this results in a fairly large relative difference between $\mathcal{S}(t)u^{\min}$ and $\mathcal{S}(t)u_0$ for small $t$. A slight increase in $t$ strongly reduces the size of the high-frequency components both in $\mathcal{S}(t)u^{\min}$ and $\mathcal{S}(t)u_0$, and the eigenvalues of $(H_\beta)^{-1}$ associated with the remaining low-frequency components are of size $O(\beta)$. As a result, the expression on the right-hand side of (6.20) will decrease over a short time interval from $O(1)$ to $O(\beta^2)$. A rough

2026 ANDREI DRĂGĂNESCU AND TODD F. DUPONT

estimate, obtained by separating the first nonzero components from the rest in the numerator of the expression in (6.20), gives

$$
\frac{\|\mathcal{S}(t)(u_0 - u^{\min})\|^2}{\|\mathcal{S}(t)u_0\|^2} \leq (\beta e^{2T})^2 + \frac{\sum_{|k|\geq 2} e^{-2t(k^2-1)}|\hat{u}_k|^2}{|\hat{u}_1|^2 + |\hat{u}_{-1}|^2 + \sum_{|k|\geq 2} e^{-2t(k^2-1)}|\hat{u}_k|^2}
$$
$$
\leq (\beta e^{2T})^2 + \frac{e^{-6t}\|u_0\|^2}{|\hat{u}_1|^2 + |\hat{u}_{-1}|^2 + e^{-6t}\|u_0\|^2} \ ,
$$

which clearly show the decay to be exponentially in $t$.

6.3. **Construction of local restriction operators.** In this section we show a way to construct local restriction operators that satisfy conditions $[\mathbf{e}, \mathbf{f}]$ in Theorem 4.3. Local operators, whether defined in an explicit or implicit way, are more attractive for computational purposes than orthogonal projections, which require inverting the mass matrix on the coarse level. The local restriction operator inherited from the classical multigrid, defined as the adjoint of the natural interpolation with respect to a mesh-dependent inner product approximating $\langle \cdot, \cdot \rangle_{L^2}$, typically satisfies the approximation property (4.6) only up to order $p = 1$ (see Example 6.7). A local, explicit, symmetric restriction operator that approximates the identity operator in negative-index norm up to higher order is defined in [5]. More precisely, for continuous piecewise linear elements the inequality (4.6) is verified for any $p \in [0, 3/2)$, and a similar result holds for quadratics. However, our estimates in Section 4 require, e.g. for linear elements, that (4.6) holds for $p = 2$. The main result of this section, Theorem 6.6, shows how to construct restriction operators with needed approximation properties.

For simplicity we restrict our exposition to the two-dimensional case. Consider the quasi-uniform triangulations $(\mathcal{T}_h)_{h \in I}$ defined in Section 6.1. For a triangle $T \in \mathcal{T}_h$ denote by $\mathcal{P}^s(T)$ the space of polynomials of total degree $\leq s - 1$ restricted to $T$, and let $s \geq 2$ be fixed. Let $\mathcal{V}_h = \mathcal{V}_h^{(s)}$ the Lagrange finite element spaces

(6.21) $\qquad \mathcal{V}_h = \{u \in \mathcal{C}^0(\overline{\Omega}) : u|_T \in \mathcal{P}^s(T) \text{ for each } T \in \mathcal{T}_h, \ u|_{\partial\Omega} \equiv 0\}.$

Denote by $(\Phi_i^h)_{1 \leq i \leq N_h}$ the Lagrange nodal basis of $\mathcal{V}_h$.

**Theorem 6.6.** *Let $R_{2h}^h \in L(\mathcal{V}_h, \mathcal{V}_{2h})$ be a restriction operator that satisfies the stability condition (4.5) and is local in the following sense: for each $i = 1, \ldots, N_h$ there exists a set $\mathcal{N}_i^h$ which is a union of triangles in $\mathcal{T}_h$, such that*

(a) $\mathrm{supp}(R_{2h}^h \Phi_i^h) \cup \mathrm{supp}(\Phi_i^h) \subset \mathcal{N}_i^h$;
(b) *no triangle $T \in \mathcal{T}_h$ lies in more than $M$ of the $\mathcal{N}_i^h$'s, with $M$ independent of $h$;*
(c) $(I - R_{2h}^h)\Phi_i^h \perp \mathcal{P}^s(\mathcal{N}_i^h)$;
(d) $\mathrm{diam}(\mathcal{N}_i^h) \leq Kh$ *with $K$ independent of $h$.*

*Then there exists a constant $C = C(K, M, s)$, independent of $h$, such that*

(6.22) $\qquad \|(I - R_{2h}^h)u\|_{\widetilde{H}^{-s}(\Omega)} \leq Ch^s\|u\|, \text{ for } u \in \mathcal{V}_h.$

*Proof.* Let $u = \sum_{i=1}^{N_h} u_i \Phi_i^h$ be a function in $\mathcal{V}_h$ and $v \in H_0^s(\Omega)$. Then

$$
\langle (I - R_{2h}^h)u, v \rangle = \sum_{i=1}^{N_h} u_i \int_{\mathcal{N}_i^h} (I - R_{2h}^h)\Phi_i^h(x)\, v(x)dx.
$$

For any $q \in \mathcal{P}^s(\mathcal{N}_i^h)$,

$$
\left| \int_{\mathcal{N}_i^h} (I - R_{2h}^h)\Phi_i^h(x)\, v(x)dx \right| \overset{(c)}{=} \left| \int_{\mathcal{N}_i^h} (I - R_{2h}^h)\Phi_i^h(x)\, (v(x) - q(x))dx \right|
$$

$$
\leq \quad \|(I - R_{2h}^h)\Phi_i^h\|_{L^2(\mathcal{N}_i^h)} \cdot \|v - q\|_{L^2(\mathcal{N}_i^h)}.
$$

The Bramble-Hilbert Lemma and (d) imply

$$
\left| \int_{\mathcal{N}_i^h} (I - R_{2h}^h)\Phi_i^h(x)\, v(x)dx \right| \overset{(4.5)}{\leq} \quad C\|\Phi_i^h\| \cdot \inf_{q \in \mathcal{P}^s(\mathcal{N}_i^h)} \|v - q\|_{L^2(\mathcal{N}_i^h)}
$$

$$
\leq \quad C(Kh)^s \|\Phi_i^h\| \cdot |v|_{H^s(\mathcal{N}_i^h)}.
$$

Therefore

$$
\left| \langle (I - R_{2h}^h)u, v \rangle \right| \leq \quad C(Kh)^s \sum_{i=1}^{N_h} |u_i| \cdot \|\Phi_i^h\| \cdot |v|_{H^s(\mathcal{N}_i^h)}
$$

$$
\leq \quad CK^s h^s \left( \sum_{i=1}^{N_h} u_i^2 \|\Phi_i^h\|^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^{N_h} |v|_{H^s(\mathcal{N}_i^h)}^2 \right)^{\frac{1}{2}}
$$

$$
\overset{(b)}{\leq} \quad (CK^s M^{\frac{1}{2}})h^s \left( \sum_{i=1}^{n} u_i^2 \|\Phi_i^h\|^2 \right)^{\frac{1}{2}} \cdot |v|_{H^s(\Omega)}
$$

$$
\leq \quad C_1 h^s \|u\| \cdot |v|_{H^s(\Omega)} \ ;
$$

for the last inequality we used the quasi-uniformity of the triangulation. It follows that

$$
\|(I - R_{2h}^h)u\|_{\widetilde{H}^{-s}(\Omega)} = \sup_{v \in H_0^s(\Omega)} \frac{\left| \langle (I - R_{2h}^h)u, v \rangle \right|}{\|v\|_{H^s(\Omega)}} \leq C_1 h^s \|u\|.
$$

$$\square$$

**Example 6.7** (Standard restriction operator on a uniform grid). Let $\Omega$ be the unit square $(0,1) \times (0,1) \subset \mathbb{R}^2$ and $(\mathcal{T}_{2^{-k}})_{k=0,1,\dots}$ the uniform three-line meshes[1] with each of the small triangles in $\mathcal{T}_{2^{-k}}$ having side length $h_k = 2^{-k}$. Obviously $\mathcal{T}_{2^{-(k+1)}}$ is the Goursat refinement of $\mathcal{T}_{2^{-k}}$. Let $s = 2$ in (6.21) and define our spaces to consist of continuous, piecewise linear, doubly-periodic functions:

$$(6.23) \qquad \mathcal{V}_{h_k}^{\text{per}} = \{u \in \mathcal{V}_{h_k} \ : \ u(x_1, 0) = u(x_1, 1),\ u(0, x_2) = u(1, x_2)\}.$$

We extend the functions in $\mathcal{V}_{h_k}$ to $\mathbb{R}^2$ in an obvious way. Let $\{P_i^k \ : \ i = 1, \dots, 2^k\}$ be the vertices of $\mathcal{T}_{2^{-k}}$ **not** lying on $\{(x_1, x_2) : x_1 = 1 \text{ or } x_2 = 1\}$, and denote by $(\Phi_i^k)_{i=1,\dots,2^k}$ the corresponding nodal basis functions. We introduce, as in [7], the mesh-dependent inner products on $\mathcal{V}_{h_k}^{\text{per}}$:

$$(6.24) \qquad \langle u, v \rangle_k \overset{\text{def}}{=} \frac{h_k^2}{2} \sum_{i=1}^{2^k} u(P_i^k)v(P_i^k), \quad \text{for } u, v \in \mathcal{V}_{h_k}^{\text{per}}.$$

---

[1] The *three-line mesh* is obtained by dividing the square into equally sized squares with sides parallel to the coordinate axes, and by further cutting each little square along its slope-one diagonal.

It should be noted that $\langle \cdot, \cdot \rangle_k$ is obtained by applying the second-order correct cubature rule

$$(6.25) \qquad \int_{\triangle P_1 P_2 P_3} u(x)\, dx \;\approx\; \text{area}(\triangle P_1 P_2 P_3) \sum_{i=1}^{3} u(P_i)$$

to the $L^2$-inner product; hence

$$(6.26) \qquad \langle u, 1 \rangle_k = \langle u, 1 \rangle = \int_\Omega u(x)\, dx, \quad \text{for } u \in \mathcal{V}_{h_k}^{\text{per}}.$$

The standard restriction operator $R_{k-1}^k : \mathcal{V}_{h_k}^{\text{per}} \to \mathcal{V}_{h_{k-1}}^{\text{per}}$ is defined by the equation

$$(6.27) \qquad \left\langle R_{k-1}^k u, v \right\rangle_{k-1} = \langle u, v \rangle_k, \quad \forall u \in \mathcal{V}_{h_k}^{\text{per}},\ v \in \mathcal{V}_{h_{k-1}}^{\text{per}}.$$

A simple calculation shows that

$$(6.28) \qquad (R_{k-1}^k u)(P_i^k) = \frac{1}{4} u(P_i^k) + \frac{1}{8} \sum_{\alpha=1}^{6} u(P_{\iota_\alpha}^k) \,,$$

where $P_{\iota_1}^k, \ldots, P_{\iota_6}^k$ are the vertices directly connected to $P_i^k$ in the graph of $\mathcal{T}_{h_k}$ (periodicity assures that every vertex has exactly six neighbors in the graph). This shows that

$$(6.29) \qquad \text{supp}(R_{k-1}^k \Phi_i^k) = \bigcup_{\alpha=1}^{6} \text{supp}(\Phi_{\iota_\alpha}^k) \stackrel{\text{def}}{=} \mathcal{N}_i^k.$$

Obviously $\mathcal{N}_i^k$ satisfies conditions (a), (b) and (d) from Theorem 6.6 (diam$(\mathcal{N}_i^k) \le 8\sqrt{2}\, h_k$). It remains to verify (c). By construction, if $u \in \mathcal{V}_{h_k}^{\text{per}}$, then

$$(6.30) \qquad \left\langle (I - R_{k-1}^k)u, 1 \right\rangle = \left\langle (I - R_{k-1}^k)u, 1 \right\rangle_k = \langle u, 1 \rangle_k - \left\langle R_{k-1}^k u, 1 \right\rangle_k = 0.$$

Orthogonality on linear functions follows by symmetry: fix $i$ and define the reflection with respect to $P_i^k$ to be $x \xmapsto{r_i^k} P_i^k - x$. Then $\mathcal{N}_i^k$ is $r_i^k$-invariant $(r_i^k(\mathcal{N}_i^k) = \mathcal{N}_i^k)$ and $\Phi_i^k$ is $r_i^k$-symmetric (i.e. $\Phi_i^k \circ r_i^k = \Phi_i^k$), and so is $(I - R_{k-1}^k)\Phi_i^k$. If we denote by $p_1$ and $p_2$ the projections onto the two coordinate axes, then the functions

$$x \mapsto (I - R_{k-1}^k)\Phi_i^k(x) \cdot p_\alpha(x - P_i^k), \quad \alpha = 1, 2,$$

are $r_i^k$-antisymmetric (if $T : A \to A$ is a bijection and $f : A \to \mathbb{C}$, then $f$ is called $T$-antisymmetric if $f \circ T + f = 0$); hence

$$(6.31) \qquad \int_{\mathcal{N}_i^k} (I - R_{k-1}^k)\Phi_i^k(x) \cdot p_\alpha(x - P_i^k)\, dx = 0.$$

The equalities (6.30) and (6.31) can be rewritten as

$$(6.32) \qquad (I - R_{k-1}^k)\Phi_i^k \perp \mathcal{P}^2(\mathcal{N}_i^h).$$

Theorem 6.6 implies that

$$(6.33) \qquad \|(I - R_{k-1}^k)u\|_{\widetilde{H}^{-2}(\Omega)} \le Ch^2 \|u\| \quad \text{for } u \in \mathcal{V}_{h_k}^{\text{per}}.$$

For unstructured triangular meshes the restriction $R_{k-1}^k$ still satisfies

$$(I - R_{k-1}^k)\Phi_i^k \perp 1 \,,$$

but typically $(I - R^k_{k-1})\Phi^k_i$ is not orthogonal on all linear functions; hence in general we can only conclude that

$$(6.34) \qquad \|(I - R^k_{k-1})u\|_{\widetilde{H}^{-1}(\Omega)} \le Ch\|u\| \quad \text{for } u \in \mathcal{V}^{\text{per}}_{h_k}.$$

*Remark* 6.8. We note that a result similar to Theorem 6.6 holds for tensor-product finite elements, and that the same argument as in Example 6.7 can be used to show that the standard restriction operator used for continuous, piecewise linear tensor-product finite elements in $\mathbb{R}^d$ also satisfies the optimal negative-index norm estimates (6.33), provided the mesh is uniform and the functions are periodic.

**Example 6.9.** We return to the framework of Section 6.1, that is, the case of an unstructured quasi-uniform triangular grid on a polygonal domain $\Omega \subset \mathbb{R}^2$, with $\mathcal{V}_h$ consisting of continuous piecewise linear functions ($s = 2$). We define a local restriction operator $R^h_{2h}$ by explicitly enforcing condition (c) in Theorem 6.6 in the following way: for each fine nodal basis function $\Phi^h_M$ choose a coarse-mesh triangle $\Delta A_1 A_2 A_3$ with all vertices in the interior of $\Omega$, such that, in the fine mesh, $M$ is adjacent or equal to one of $A_1$, $A_2$ or $A_3$. (E.g., if $M$ is a coarse node, choose $A_1 = M$; if $M$ is the midpoint of an interior edge, let $A_1, A_2$ be the vertices of that edge; if $M$ is the midpoint of an edge connecting a coarse-mesh vertex $A_1$ to the boundary, then again let $A_2, A_3$ be two neighbors of $A_1$ that are adjacent to each other.) Let $(\Phi^{2h}_i)_{i=1,2,3}$ be the nodal variables in $\mathcal{V}_{2h}$ that are associated with $A_{1,2,3}$ respectively. We define

$$(6.35) \qquad R^h_{2h}\Phi^h_M = \sum_{j=1}^3 \alpha_j \Phi^{2h}_j \,,$$

with $(\alpha_j)_{j=1,2,3}$ chosen so that

$$(6.36) \qquad (I - R^h_{2h})\Phi^h_M \perp \mathcal{L}, \quad \forall \mathcal{L} \in \mathcal{P}^2.$$

Let $(L_i)_{i=1,2,3}$ be the basis of $\mathcal{P}^2$ obtained by linearly extending $\Phi^{2h}_i|_{\Delta A_1 A_2 A_3}$ to $\mathbb{R}^2$. The system (6.36) translates into

$$(6.37) \qquad \sum_{j=1}^3 \left\langle \Phi^{2h}_j, L_i \right\rangle \alpha_j = \left\langle \Phi^h_M, L_i \right\rangle, \quad i = 1, 2, 3.$$

Under additional regularity assumptions on the meshes $(\mathcal{T}_h)_{h \in I}$ we can show that the system (6.37) is diagonally dominant, therefore nonsingular, and that the resulting restriction operator is uniformly bounded, thus verifying condition [e] in Theorem 4.3. More precisely we can show that for every node $M$ in the fine mesh, the system (6.37) has a solution $(\alpha_j)_{j=1,2,3}$ with $|\alpha_j| \le C$, where $C$ is independent of $h$. This implies that the restriction operators are uniformly bounded. For details see Appendix A. The other hypotheses in Theorem 6.6 are satisfied by construction; therefore

$$(6.38) \qquad \|(I - R^h_{2h})u\|_{\widetilde{H}^{-2}(\Omega)} \le Ch^2 \|u\|, \quad \forall u \in \mathcal{V}_h.$$
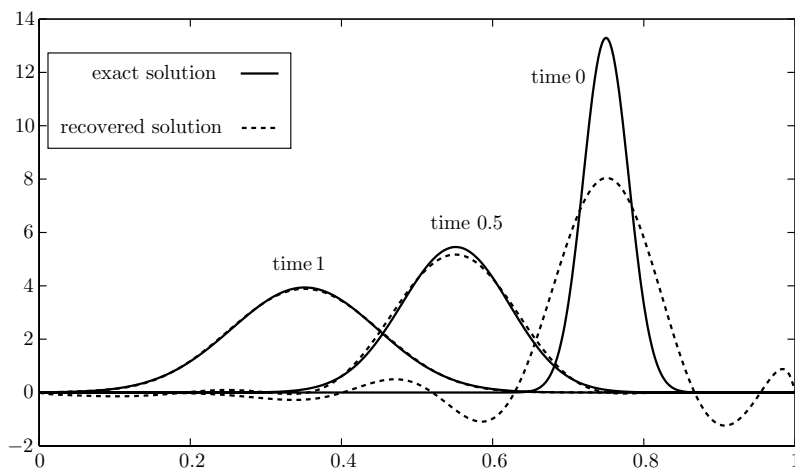
FIGURE 1. Evolution of exact solution versus recovered system for $\beta = 10^{-3}$; we used 400 intervals.

## 7. A NUMERICAL EXAMPLE

We present a simple numerical example that illustrates the application of the multilevel preconditioner MLAS to the inverse problem defined in Section 6.1.

**Data:** In this example, $\Omega = (0, 1)$, and the coefficients are $a_{11}(x, t) \equiv 4 \cdot 10^{-3}$, $b_1(x, t) \equiv 0.4$, $c(x, t) \equiv 0.05$. The exact initial value $u_0$ to be recovered is a Gaussian with $\mu = 0.75$ and $\sigma = 0.03$. The final time is $T = 1$. We consider six uniform grids on $(0, 1)$ with the base grid having $h_0 = 1/200$, and for sequent grids $h_k = 2^{-k}h_0$, $i = 1, 2, 3, 4, 5$. We approximate the solution by its Galerkin projection onto the space of continuous piecewise linear functions, and we use backward Euler for time discretization. At each level the time step is uniform, with the base level time step of $dt_0 = 1/100$. Due to the fact that the approximation in time is only first order, we refine the time step by $dt_{k+1} = dt_k/4$. For each of the levels $k = 0, 1, 2, 3, 4$ we solve the inverse problem using CG preconditioned by MLAS with $1 \leq i \leq 6 - k$ levels of refinement (one level of refinement simply means unpreconditioned CG); that is, with each of the base cases we refine the problem until the maximal resolution of $h_5 = 1/6400$ is reached. For each of $\beta = 10^{-3}, 10^{-4}, 10^{-6}$ we test the V-cycle and the W-cycle preconditioner (MLAS) from Section 5. MLAS is modified so that, at the finest level, only one recursive call to the coarse space solver is performed at each iteration; therefore no residual is computed at the finest level **in the preconditioner**. This way comparison between the W-cycle and the V-cycle preconditioner is fair (in fact the two are identical if only two levels are used). In Figure 1 we show the time-evolution of the system starting from the "exact" initial value versus the recovered initial value at times $0, 0.5$ and $1$ for $\beta = 10^{-3}$.

**Targets:** Due to the fact that the Hessian is never formed, it is difficult to evaluate directly the spectral distance between the inverse of the Hessian and the tested preconditioners. We evaluate the quality of the preconditioners indirectly by looking at the number of preconditioned CG iterations needed to obtain a residual whose

norm is $rtol = 10^{-12}$ times smaller than the norm of the right-hand side. We use the zero-function as the initial guess. Our goal is to verify the following consequences of Theorems 5.2 and 5.4:

**1.** V-cycle: for a fixed base level the number of iterations $N_{\text{it}}$ is independent of the number of refinement levels (provided the base level is sufficiently fine); $N_{\text{it}}$ decreases with increasing base-level resolution;

**2.** W-cycle: given a sufficiently fine base level, $N_{\text{it}}$ decreases with increasing number of refinement levels; in fact, $N_{\text{it}}$ depends only on the finest resolution (if at least two levels are used).
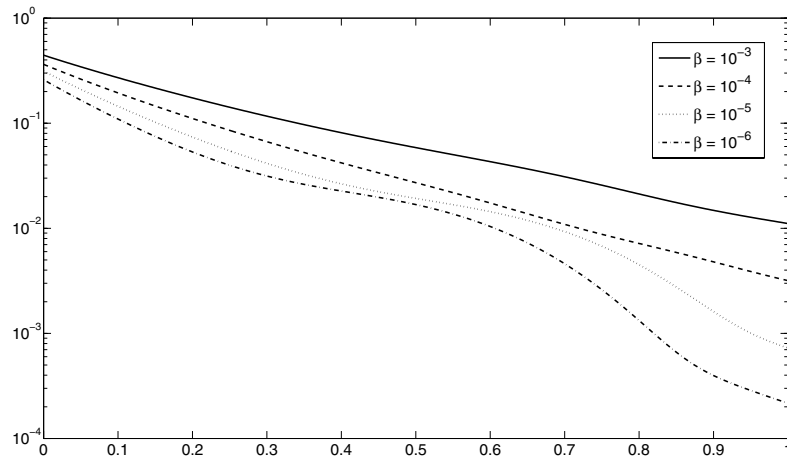
Furthermore, we verify that unpreconditioned CG solves the problem in a number of iterations that is independent of the mesh size, and increases slowly with $\beta \to 0$. We consider two ways of measuring the success of the solution process. An absolute measuring unit is the number of flops needed for a direct (forward) solve. We call the cost of the inverse solve measured in forward solves the *efficiency factor*, denoted by I/F:

$$\text{I/F} = \frac{\text{cost of inverse solve}}{\text{cost of forward solve}}.$$

For unpreconditioned CG each iteration requires a matrix-vector multiplication, which costs approximately two forward solves. Therefore the CG solve with $N_{\text{it}}$ iterations costs about $2N_{\text{it}}$ forward solves; the whole solution process includes also a gradient computation; therefore, the efficiency factor for the unpreconditioned CG is $\approx 2(N_{\text{it}} + 1)$. In Tables 1–6 we show for each case the number of iterations, and in parantheses the cost of the solution process for the inverse problem computed in "forward solves". An alternative way of measuring success, which inherently takes into account the difficulty of the problem, is by using the work of unpreconditioned CG as a measuring unit (since it is quasi-proportional with the space-time size of the problem). E.g., if the measuring time $T$ is large, then nearly all information about the initial value we are trying to recover is lost, and essentially we are inverting the identity operator. CG captures this fact by solving the problem in a small number of iterations.

**Results and conclusions:** A visual inspection of the recovered solution in Figure 1 (dotted line at time 0) shows that the recovered $u^{\text{min}}$ is a smooth curve, as predicted by Lemma 6.4. At the same time it shows that, despite the recovery of the initial value itself not being very accurate ($u^{\text{min}}$ exhibits low-frequency oscillations), $\mathcal{S}(0.5)u^{\text{min}}$ is very close to $\mathcal{S}(0.5)u_0$; note that no measurements were taken at time $t = 0.5$. The plot in Figure 2 shows an exponential-like decay of the relative error $\|\mathcal{S}(t)(u^{\text{min}} - u_0)\|/\|\mathcal{S}(t)u_0\|$, as expressed in Example 6.5.

The left column in Tables 1–6 shows the number $N$ of intervals used for discretization; the column headers indicate the number of refinement levels used. E.g., for $N = 400$ and three levels of refinements the finest grid will have $2^{3-1}N = 1600$ intervals. The entries in the first column support the assertion that unpreconditioned CG solves the inverse problem in a mesh-independent number of iterations. The two targets mentioned above are verified in the runs with $\beta = 10^{-3}$ (Tables 1 and 2), and $\beta = 10^{-4}$ (Tables 3 and 4). Most importantly, we notice that the number of MLAS-preconditioned CG iterations depends only on the finest-level resolution, independent of the number of levels. In particular, this confirms the conclusion

FIGURE 2. Time evolution of relative error $\|\mathcal{S}(t)(u^{\min} - u_0)\|/\|\mathcal{S}(t)u_0\|$

of Theorem 5.4, thus showing MLAS to have the same approximation quality as TLAS. A stronger decrease in regularization parameter, namely $\beta = 10^{-6}$, causes the base case with $N = 200$ to become unacceptably coarse. Specifically we see the V-cycle preconditioned CG taking a fairly large number of iterations, thus making this solution process as expensive as the unpreconditioned CG. Also as a result of the base case being too coarse we see the W-cycle preconditioner losing its positive-definiteness. Another interesting fact is noticed when using $N = 400$ as a base case. Although the V-cycle preconditioner stagnates after $\approx 14$ iterations (Table 5), the W-cycle algorithm shows a certain ability of self-correction, thus overcoming the suboptimal choice of base level. The correct base level for $\beta = 10^{-6}$ seems to be the one with $N = 800$, a point at which predicted behavior sets in. In terms of efficiency at a given resolution, the I/F factor readings suggest that MLAS is most efficient when maximizing the number of levels, constrained by having the base level acceptably coarse. Using this strategy we found MLAS-preconditioned CG to work up to four times faster than unpreconditioned CG (e.g. Table 6, $N = 800$ with 4 levels has an I/F of 17.7 compared to a predicted I/F of 72 for unpreconditioned CG).

TABLE 1. Iteration count (I/F) for the V-cycle; $\beta = 10^{-3}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 15 | (32.3) | 11 | (61.1) | 12 | (31.2) | 12 | (26.4) | 12 | (26.1) | 12 | (26) |
| 400 | 16 | (34.1) | 9 | (48) | 9 | (25.7) | 10 | (22.4) | 10 | (22) | | |
| 800 | 16 | (34) | 7 | (38) | 8 | (20.9) | 8 | (18.4) | | | | |
| 1600 | 16 | (34) | 6 | (32) | 6 | (16.3) | | | | | | |
| 3200 | 17 | (36) | 5 | (26.7) | | | | | | | | |

TABLE 2. Iteration count (I/F) for the W-cycle; $\beta = 10^{-3}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 15 | (32.3) | 11 | (61.1) | 9 | (29.6) | 7 | (19.4) | 6 | (16.2) | 5 | (13.7) |
| 400 | 16 | (34.1) | 9 | (48) | 7 | (22.8) | 6 | (16.8) | 5 | (13.8) | | |
| 800 | 16 | (34) | 7 | (38) | 6 | (19.8) | 5 | (14.4) | | | | |
| 1600 | 16 | (34) | 6 | (32) | 5 | (16.9) | | | | | | |
| 3200 | 17 | (36) | 5 | (26.7) | | | | | | | | |

TABLE 3. Iteration count (I/F) for the V-cycle; $\beta = 10^{-4}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 21 | (44.4) | 13 | (82.2) | 13 | (35.1) | 13 | (28.9) | 13 | (28.1) | 13 | (28) |
| 400 | 20 | (42.1) | 10 | (61.3) | 11 | (29.5) | 11 | (24.7) | 11 | (24.1) | | |
| 800 | 21 | (44) | 8 | (47.5) | 8 | (21.8) | 8 | (18.5) | | | | |
| 1600 | 21 | (44) | 6 | (36.2) | 6 | (16.9) | | | | | | |
| 3200 | 21 | (44) | 5 | (31.7) | | | | | | | | |

TABLE 4. Iteration count (I/F) for the W-cycle; $\beta = 10^{-4}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 21 | (44.4) | 13 | (82.2) | 15 | (49.6) | 10 | (27.5) | 8 | (21) | 6 | (16.1) |
| 400 | 20 | (42.1) | 10 | (61.3) | 7 | (25.5) | 7 | (19.5) | 5 | (13.9) | | |
| 800 | 21 | (44) | 8 | (47.5) | 6 | (21) | 6 | (17) | | | | |
| 1600 | 21 | (44) | 6 | (36.2) | 5 | (18) | | | | | | |
| 3200 | 21 | (44) | 5 | (31.7) | | | | | | | | |

TABLE 5. Iteration count (I/F) for the V-cycle; $\beta = 10^{-6}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 27 | (56.7) | 19 | (169) | 23 | (65.6) | 25 | (54.3) | 22 | (46.3) |
| 400 | 32 | (66.2) | 15 | (117.4) | 14* | | 13* | | – | |
| 800 | 34 | (70) | 9 | (73.2) | 10 | (29.5) | 10 | (22.9) | | |
| 1600 | 34 | (70) | 7 | (57.2) | 7 | (21.2) | | | | |
| 3200 | 35 | (72) | 6 | (45.7) | | | | | | |

\* stagnated at a low residual (corresponding to $rtol = 10^{-11}$) before converging.

TABLE 6. Iteration count (I/F) for the W-cycle; $\beta = 10^{-6}$

| N | 1 | | 2 | | 3 | | 4 | | 5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| 200 | 27 | (56.7) | 19 | (169) | –† | | –† | | –† | |
| 400 | 32 | (66.2) | 15 | (117.4) | 10 | (39.3) | 9 | (25.6) | 7 | (18.8) |
| 800 | 34 | (70) | 9 | (73.2) | 8 | (30.6) | 6 | (17.7) | | |
| 1600 | 34 | (70) | 7 | (57.2) | 6 | (23.7) | | | | |
| 3200 | 35 | (72) | 6 | (45.7) | | | | | | |

† multilevel preconditioner not positive definite.

APPENDIX A. SUFFICIENT CONDITIONS FOR WELL DEFINEDNESS
OF CERTAIN LOCAL RESTRICTION OPERATORS

In this section we refer to the notation and context of Example 6.9. We introduce a set of conditions on the meshes $(\mathcal{T}_h)_{h \in I}$ that prove sufficient for the systems (6.37) to have uniformly bounded solutions.

For each vertex $A$ of $\mathcal{T}_h$ we denote its *vicinity* by

$$\mathcal{K}_A^h = \bigcup \{T \in \mathcal{T}_h \; : \; A \text{ is a vertex of } T\}.$$

Furthermore, let $d_A^h = \sup\{r > 0 \; : \; \mathcal{B}_r(A) \subseteq \mathcal{K}_A^h\}$ and $D_A^h = \inf\{r > 0 \; : \; \mathcal{B}_r(A) \supseteq \mathcal{K}_A^h\}$. If we assume that all triangles in the meshes $(\mathcal{T}_h)_{h \in I}$ have angles $\leq \pi/2$, then the vicinity of each vertex is convex. The next condition states that the vicinity is somewhat balanced around the vertex.

**Condition A.1.** There are two constants $c_1, c_2 \geq 1$, independent of $h$, such that:
  (i) for any vertex $A$ of $\mathcal{T}_h$, $D_A^h \leq c_1 d_A^h$;
  (ii) for any neighboring vertices $A, B$ of $\mathcal{T}_h$, we have $c_2^{-1} d_A^h \leq d_B^h \leq c_2 d_A^h$.

It is easy to verify that (i) implies (ii) in Condition A.1, with $c_2 = c_1$; however, in general, $c_2 \leq c_1$, and here it is advantageous to keep track of the individual constants. E.g., for a *regular* mesh (we call a mesh regular if all triangles are equilateral) $c_1 = 2/\sqrt{3}$, and $c_2 = 1$, and for the usual three-line mesh (obtained from a rectangular grid with square elements by cutting each square along its slope-one diagonal) we have $c_1 = 2$, $c_2 = 1$. We wish to point out that, in case the triangulations are obtained by Goursat refinement, then the constants $c_1, c_2$ in Condition A.1 propagate from $\mathcal{T}_{h_0}$ to finer meshes; in other words, finding the optimal constants in Condition A.1 only requires looking at the coarsest grid. For a point $A$ and $\rho > 0$ denote by $\psi_{A,\rho}$ the "cone-hat" function

$$\psi_\rho^A(x) = \begin{cases} 1 - \|x - A\|/\rho & \text{for } x \in \mathcal{B}_\rho(A), \\ 0 & \text{otherwise.} \end{cases}$$

It follows that a nodal basis function $\Phi_A$ satisfies

(A.1)                                 $\psi_{d_A^h}^A \leq \Phi_A^h \leq \psi_{D_A^h}^A.$

The convexity of the $\mathcal{K}_{A_i}^{2h}$ implies that $L_i \geq 0$ on $\text{supp}(\Phi_i^{2h})$; therefore

(A.2)                 $a_{ii} = \left\langle \Phi_i^{2h}, L_i \right\rangle \geq \int_{\mathcal{B}_{d_i}(A_i)} \psi_{d_i}^{A_i} L_i = \frac{\pi d_i^2}{3}\,,$

where $d_i = d_{A_i}^{2h}$, $D_i = D_{A_i}^{2h}$, $i = 1, 2, 3$. For $i \neq j$,

(A.3)                 $a_{ij} = \left\langle \Phi_j^{2h}, L_i \right\rangle \overset{(A.1)}{\leq} \int_{L_i \geq 0} \psi_{D_j}^{A_j} L_i + \int_{L_i \leq 0} \psi_{d_j}^{A_j} L_i.$

Let $\delta_i$ be the distance from $A_i$ to the opposite side in $\Delta A_1 A_2 A_3$. Then $\delta_i \geq d_i \geq c_2^{-1} d_j$, and

(A.4)      $\int_{L_i \geq 0} \psi_{D_j}^{A_j} L_i = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} d\theta \int_0^{D_j} \left(1 - \frac{\rho}{D_j}\right) \cdot \frac{\rho^2 \cos \theta}{\delta_i} \, d\rho = \frac{D_j^3}{6\,\delta_i}.$

Similarly

(A.5)                 $\int_{L_i \leq 0} \psi_{d_j}^{A_j} L_i = -\int_{L_i \geq 0} \psi_{d_j}^{A_j} L_i = \frac{d_j^3}{6\,\delta_i}.$

Hence

(A.6)
$$a_{ij} \leq \frac{D_j^3 - d_j^3}{6\,\delta_i} \leq c_2 \frac{c_1^3 - 1}{6}\, d_j^2 .$$

Since $-a_{ij}$ satisfies the same inequality we conclude that (A.6) holds with $|a_{ij}|$ in place of $a_{ij}$. We should note that the estimate (A.6) can be quite inaccurate even for "tame" meshes such as the three-line mesh mentioned previously, where symmetry implies that $a_{ij} = 0$ for $i \neq j$, while (A.6) is insufficient to prove that $|a_{ij}| \leq a_{ii}$ (because $7/6 > \pi/3$). If we divide the $i^{\text{th}}$ equation of (6.37) by $a_{ii}$, then the resulting equivalent system will be represented by a matrix of the form $I - B$, $B = (b_{ij})_{1 \leq i,j \leq 3}$, with $b_{ij} = (\delta_{ij} - 1)a_{ij}/a_{ii}$. If

(A.7)
$$\max_{i=1}^{3} \sum_{j \neq i} |a_{ij}| \leq (1 - \gamma)a_{ii} ,$$

for some $\gamma \in (0, 1)$, then $\|B\|_\infty \leq 1 - \gamma$; hence

(A.8)
$$\|(I - B)^{-1}\|_\infty \leq \sum_{k=0}^{\infty} \|B\|_\infty^k = \gamma^{-1}.$$

We conclude our discussion with

**Lemma A.2.** *If the meshes $(\mathcal{T}_h)_{h \in I}$ satisfy Condition A.1 such that*

(A.9)
$$\gamma \stackrel{\text{def}}{=} 1 - \frac{c_2(c_1^3 - 1)}{\pi} > 0 ,$$

*then there exists a constant $C$ independent of $h$ such that for each triangle-vertex of $\mathcal{T}_h$ and choice of coarse-mesh triangle $\Delta A_1 A_2 A_3 \in \mathcal{T}_{2h}$, the solution $(\alpha_j)_{j=1,2,3}$ of the system (6.37) satisfies $|\alpha_j| \leq C$, $j = 1, 2, 3$.*

*Proof.* Quasi-uniformity of the meshes implies that for each vertex $M$ of $\mathcal{T}_h$,

$$c\,h \leq d_M^h \leq C\,h ,$$

for some constants $c, C$ independent of $h$. Therefore, due to the proximity of $M$ to the chosen triangle $\Delta A_1 A_2 A_3 \in \mathcal{T}_{2h}$, the right-hand side of the modified system (6.37) satisfies $|\langle \Phi_M^h, L_i \rangle|/a_{ii} \leq C'$. In light of the estimates (A.2) and (A.6), condition (A.9) implies (A.7). By (A.8) we obtain $|\alpha_j| \leq C'/\gamma$, $j = 1, 2, 3$. $\square$

The condition (A.9) is not necessary for the system (6.37) to be nonsingular, as seen from the three-line mesh example. However, it is not very restrictive either, given that for a regular mesh the value of $\gamma$ is fairly large ($\gamma \approx 0.82$). Hence any mesh that is locally "alike" to a regular mesh will likely satisfy (A.9).

Appendix B. Notation summary

For quick reference we provide a list of the notation used throughout the article, and the place where each is first used or defined, where applicable:

| | | |
|---|---|---|
| $\mathcal{K}$ | a compact operator from $L^2(\Omega)$ to $L^2(\Omega)$ | Section 1 |
| $\|\mathcal{K}\|$ | the operator norm of $\mathcal{K}$ (from $L^2(\Omega)$ to $L^2(\Omega)$) | Section 2 |
| $H$ | $K^*K$ | (1.3) |
| $H_\beta$ | $I + \beta^{-1}H$ | (1.3) |
| $\mathcal{V}_h$ | approximation space | Section 2 |
| $h_0$ | coarsest $h$ that can be chosen as base case | |
| $\mathcal{K}^h$ | discretization of $\mathcal{K}$ in $\mathcal{V}_h$ | Section 2 |
| $H^h$ | $(K^h)^*K^h$ | (2.10) |
| $H_\beta^h$ | $I + \beta^{-1}H^h$ | (2.10) |
| $\pi_h$ | $L^2$-projection onto $\mathcal{V}_h$ | Section 2 |
| $d_{\mathcal{X}}(T_1, T_2)$ | spectral distance between $T_1$ and $T_2$ | Def. 3.1 |
| $\mathcal{N}_H(X)$ | $2H - XHX$ | (3.18) |
| $d_h$ | $d_{\mathcal{V}_h}$ | Section 4 |
| $R_{2h}^h$ | restriction operator | Section 4.2 |
| $L_\beta^h$ | two-level preconditioner | (4.1), (4.7) |
| $\mathcal{G}^h(T)$ | $T\pi_{2h} + (I - \pi_{2h})$ | (5.2) |
| $K_\beta^h$ | MLAS preconditioner | (5.3) |
| $e_h$ | $d_h(K_\beta^h, (H_\beta^h)^{-1})$ | (5.7) |
| $\mathcal{S}(t)$ | solution operator for a parabolic PDE | Section 6 |

References

1. Volkan Akçelik, George Biros, Andrei Drăgănescu, Omar Ghattas, Judith C. Hill, and Bart G. van Bloemen Waanders, *Dynamic data driven inversion for terascale simulations: Real-time identification of airborne contaminants*, Proceedings of SC2005 (Seattle, WA), IEEE/ACM, November 2005.

2. Randolph E. Bank and Todd Dupont, *An optimal order process for solving finite element equations*, Math. Comp. **36** (1981), no. 153, 35–51. MR82b:65113

3. James H. Bramble, *Multigrid methods*, Pitman Research Notes in Mathematics Series, vol. 294, Longman Scientific & Technical, Harlow, 1993. MR1247694 (95b:65002)

4. James H. Bramble and Joseph E. Pasciak, *New estimates for multilevel algorithms including the V-cycle*, Math. Comp. **60** (1993), no. 202, 447–471. MR1176705 (94a:65064)

5. James H. Bramble, Joseph E. Pasciak, and Panayot S. Vassilevski, *Computational scales of Sobolev norms with application to preconditioning*, Math. Comp. **69** (2000), no. 230, 463–480. MR1651742 (2000k:65088)

6. James H. Bramble, Joseph E. Pasciak, and Jinchao Xu, *Parallel multilevel preconditioners*, Math. Comp. **55** (1990), no. 191, 1–22. MR1023042 (90k:65170)

7. Susanne C. Brenner and L. Ridgway Scott, *The mathematical theory of finite element methods*, second ed., Texts in Applied Mathematics, vol. 15, Springer-Verlag, New York, 2002. MR2003a:65103

8. William L. Briggs, Van Emden Henson, and Steve F. McCormick, *A multigrid tutorial*, second ed., Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2000. MR2001h:65002

9. Henri Cartan, *Calcul différentiel*, Hermann, Paris, 1967. MR36:6243

10. E. G. D'Jakanov, *On an iterative method for the solution of a system of finite difference equations*, Dokl. Akad. Nauk. SSSR **138** (1961), 522–525.

11. Andrei Drăgănescu, *A fast multigrid method for inverting linear parabolic problems*, Tech. Report TR-2004-01, University of Chicago, Department of Computer Science, 2004, presented at the Eighth Copper Mountain Conference on Iterative Methods.

12. ———, *Two investigations in numerical analysis: Monotonicity preserving finite element methods, and multigrid methods for inverse parabolic problems*, Ph.D. thesis, University of Chicago, August 2004.

13. Andrei Drăgănescu and Todd F. Dupont, *A multigrid algorithm for backwards reaction-diffusion equations*, in preparation.

14. Todd Dupont, *A factorization procedure for the solution of elliptic difference equations*, SIAM J. Numer. Anal. **5** (1968), 753–782. MR0246528 (39:7832)

15. Heinz W. Engl, *Necessary and sufficient conditions for convergence of regularization methods for solving linear operator equations of the first kind*, Numer. Funct. Anal. Optim. **3** (1981), no. 2, 201–222. MR627122 (82j:47018)

16. Heinz W. Engl, Martin Hanke, and Andreas Neubauer, *Regularization of inverse problems*, Mathematics and its Applications, vol. 375, Kluwer Academic Publishers Group, Dordrecht, 1996. MR97k:65145

17. James E. Gunn, *The solution of elliptic difference equations by semi-explicit iterative techniques*, J. Soc. Indust. Appl. Math. Ser. B Numer. Anal. **2** (1965), 24–45. MR0179962 (31:4199)

18. Karl E. Gustafson and Duggirala K. M. Rao, *Numerical range*, Universitext, Springer-Verlag, New York, 1997, The field of values of linear operators and matrices. MR98b:47008

19. Wolfgang Hackbusch, *Multigrid methods and applications*, Springer Series in Computational Mathematics, vol. 4, Springer-Verlag, Berlin, 1985. MR814495 (87e:65082)

20. Martin Hanke and Curtis R. Vogel, *Two-level preconditioners for regularized inverse problems. I. Theory*, Numer. Math. **83** (1999), no. 3, 385–402. MR2001h:65069

21. Thorsten Hohage, *Regularization of exponentially ill-posed problems*, Numer. Funct. Anal. Optim. **21** (2000), no. 3-4, 439–464. MR1769885 (2001e:65095)

22. Barbara Kaltenbacher, *On the regularizing properties of a full multigrid method for ill-posed problems*, Inverse Problems **17** (2001), no. 4, 767–788. MR1861481 (2002h:65094)

23. ———, *V-cycle convergence of some multigrid methods for ill-posed problems*, Math. Comp. **72** (2003), no. 244, 1711–1730 (electronic). MR1986801 (2004d:65069)

24. Barbara Kaltenbacher and Josef Schicho, *A multi-grid method with a priori and a posteriori level choice for the regularization of nonlinear ill-posed problems*, Numer. Math. **93** (2002), no. 1, 77–107. MR1938323 (2003h:65076)

25. J. Thomas King, *Multilevel algorithms for ill-posed problems*, Numer. Math. **61** (1992), no. 3, 311–334. MR1151773 (92k:65090)

26. Mitchell Luskin and Rolf Rannacher, *On the smoothing property of the Galerkin method for parabolic equations*, SIAM J. Numer. Anal. **19** (1982), no. 1, 93–113. MR83c:65245

27. ———, *On the smoothing property of the Crank-Nicolson scheme*, Applicable Anal. **14** (1982/83), no. 2, 117–135. MR83m:65072

28. Rolf Rannacher, *Finite element solution of diffusion problems with irregular data*, Numer. Math. **43** (1984), no. 2, 309–327. MR85c:65145

29. Teresa Regińska, *Regularization of discrete ill-posed problems*, BIT **44** (2004), no. 1, 119–133. MR2057365 (2005h:65092)

30. Andreas Rieder, *A wavelet multilevel method for ill-posed problems stabilized by Tikhonov regularization*, Numer. Math. **75** (1997), no. 4, 501–522. MR97k:65299

31. Vidar Thomée, *Galerkin finite element methods for parabolic problems*, Springer Series in Computational Mathematics, vol. 25, Springer-Verlag, Berlin, 1997. MR98m:65007

32. Curtis R. Vogel, *Computational methods for inverse problems*, Frontiers in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. MR2003i:65004

Sandia National Laboratories,[2] Albuquerque, New Mexico 87125
*Current address*: Department of Mathematics and Statistics, University of Maryland, Baltimore County, Baltimore, Maryland 21250
*E-mail address*: `draga@math.umbc.edu`

Department of Computer Science, University of Chicago, Chicago, Illinois 60637
*E-mail address*: `t-dupont@uchicago.edu`

---