

# SPECTRAL ANALYSIS OF IMPLICIT $s$ -STAGE BLOCK RUNGE-KUTTA PRECONDITIONERS\*

MARTIN J. GANDER<sup>†</sup> AND MICHAL OUTRATA<sup>‡</sup>

**Abstract.** We analyze the recently introduced family of preconditioners in [23] for the stage equations of implicit Runge-Kutta methods for  $s$ -stage methods. We simplify the formulas for the eigenvalues and eigenvectors of the preconditioned systems for a general  $s$ -stage method and use these to obtain convergence rate estimates for preconditioned GMRES for some common choices of the implicit Runge-Kutta methods. This analysis is based on understanding the inherent matrix structure of these problems and exploiting it to qualitatively predict and explain the main observed features of the GMRES convergence behavior, using tools from approximation and potential theory based on Schwarz-Christoffel maps for curves and close, connected domains in the complex plane. We illustrate our analysis with numerical experiments showing very close correspondence of the estimates and the observed behavior, suggesting the analysis reliably captures the essence of these preconditioners.

**Key words.** implicit Runge-Kutta methods, stage equations, preconditioned GMRES, convergence estimates, Schwarz-Christoffel maps, potential theory

**MSC codes.** 65L06, 65F10, 65E05

**1. Introduction.** Runge-Kutta methods are a well-established family of one-step solvers for systems of ordinary differential equations (ODEs; see [31, 30] for an overview and further references). For implicit methods (IRK), their efficiency depends on the efficiency of a solver for the so-called *stage equations* – in general a system of  $ms$  non-linear equations, where  $m$  is the number of scalar ODEs in the system and  $s$  is the number of stages of the Runge-Kutta method. An important application arises from the space discretization of time-dependent partial differential equations (PDEs), resulting in a system of ODEs with *very* large  $m$ . If the spatial operator is *linear*, then the stage equations also form a system of linear algebraic equations and are often solved by an iterative solver, e.g., a Krylov method. In [23], the authors introduced a family of preconditioners for GMRES for the stage equations, numerically showing that these preconditioners give an *outstanding* performance, especially under refinement of the spatial mesh, i.e., as  $m$  grows. Recently, there have also been other contributions in the direction of preconditioning the *fully implicit* Runge-Kutta stage equations for PDEs, see [27, 26] but also [20, 19] and [3], introducing new ideas in terms of implementation as well as formulation and testing these numerically on a variety of test problems.

We focus on the setting considered in [23], expand the 2-stage method analysis given in [10], and consider the general  $s$ -stage case, giving a theoretical background for the performance and spectral properties observed. Using the classical ideal GMRES bound we use the structural properties of the stage equations to obtain computable expressions for the spectrum. These then justify the use of estimates based on conformal mapping theory (see [5]) of the ideal GMRES bound and ultimately lead to descriptive estimates for GMRES convergence properties for the preconditioned systems.

---

\*Submitted to the editors DATE.

**Funding:** This work was partially supported by the SNF grant number 178752 and by the FCS Swiss Excellence PhD Fellowship program of the Swiss Federation (ESKAS No. 2019.0384).

<sup>†</sup>Section de Mathématiques, Université de Genève

<sup>‡</sup>Section de Mathématiques, Université de Genève

First, we recall some important preliminaries in Section 2 so that we can deliver the analysis, based on the spectral analysis of the preconditioned system, in Section 3. We support the analysis by considering more involved examples in Section 4.

**2. Model problem and preliminaries.** The analysis in this paper applies to any spatial discretizations of  $\partial_t u = \mathcal{L}u + f$  with a diffusive elliptic operator  $\mathcal{L}$  that leads to a symmetric definite problem (the main assumptions being (3.6) in Section 3). However, in order to facilitate the understanding and put the emphasis on the preconditioners and their performance we choose for its exposition the simplest concrete problem and its discretization – the heat equation. We thus consider the heat equation on the unit square and a time interval  $(0, T_{\text{end}})$ , i.e.,

$$(2.1) \quad \begin{aligned} \frac{\partial}{\partial t} u &= \Delta u + f \quad \text{in } \Omega \times (0, T_{\text{end}}), \\ u &= g \quad \text{on } \partial\Omega \times (0, T_{\text{end}}) \quad \text{and} \quad u = u_0 \quad \text{in } \Omega \times \{0\}, \end{aligned}$$

where  $\Delta$  is the Laplace operator,  $f, g, u_0$  are given functions and  $\Omega$  is the unit square  $\Omega := (0, 1) \times (0, 1)$ . As in [10] we discretize in space using a finite difference scheme on an equidistant grid with  $N + 1$  rows and columns, and with mesh size  $h = 1/N$ . The values at the interior grid points become unknown functions of time, which are governed by the system of ODEs

$$(2.2) \quad \frac{\partial}{\partial t} u_i(t) = \frac{u_{i-N}(t) + u_{i-1}(t) - 4u_i(t) + u_{i+1}(t) + u_{i+N}(t)}{h^2} + b_i^{(ST)}(t),$$

for  $i = N + 1, \dots, N(N - 1) - 1$ , where  $b_i^{(ST)}(t)$  collects the known values from the source terms, given by  $g$  and  $f$ , at the given point. Combining the unknowns in each grid column into one vector denoted by  $\mathbf{u}_k(t)$ , i.e.,

$$\mathbf{u}_k(t) := [u_{Nk+2} \quad u_{Nk+3} \quad \cdots \quad u_{N(k+1)-1}]^T(t), \quad \mathbf{u}(t) := [\mathbf{u}_1^T(t) \quad \cdots \quad \mathbf{u}_{N-1}^T(t)]^T,$$

and also analogously for  $\mathbf{b}_k(t)$  and  $\mathbf{b}(t)$ , we rewrite (2.2) as

$$(2.3) \quad \frac{\partial}{\partial t} \mathbf{u}(t) = \frac{1}{h^2} L \mathbf{u}(t) + \mathbf{b}^{(ST)}(t),$$

with

$$(2.4) \quad L = \begin{bmatrix} T & I & & \\ I & \ddots & \ddots & \\ & \ddots & \ddots & I \\ & & I & T \end{bmatrix}, \quad T = \begin{bmatrix} -4 & 1 & & \\ 1 & \ddots & \ddots & \\ & \ddots & \ddots & 1 \\ & & 1 & -4 \end{bmatrix}, \quad I = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix},$$

where  $L$  is of dimension  $n := (N - 1)^2$  and the blocks  $T, I$  are of dimension  $N - 1$ . We discretize  $[0, T_{\text{end}}]$  with  $M_{T_{\text{end}}} + 1$  equidistant time points with time step  $\tau = T_{\text{end}}/M_{T_{\text{end}}}$ , i.e.,

$$\{0 = t_0 < \cdots < t_{M_{T_{\text{end}}}} = T_{\text{end}}\}, \quad \tau = \frac{T_{\text{end}}}{M_{T_{\text{end}}}} \quad \text{and} \quad t_m = \tau \cdot m, \quad m = 0, \dots, M_{T_{\text{end}}}.$$

72 Having a *Butcher tableau*

$$73 \quad (2.5) \quad \begin{array}{c|c} \mathbf{c} & A \\ \hline & \mathbf{b} \end{array} := \begin{array}{c|ccc} c_1 & a_{1,1} & \dots & a_{1,s} \\ \vdots & \vdots & \ddots & \vdots \\ c_s & a_{s,1} & \dots & a_{s,s} \\ \hline & b_1 & \dots & b_s \end{array},$$

74 the corresponding IRK method applied to (2.3) at the  $m$ -th time step gives the ap-  
75 proximation  $\mathbf{u}^m \approx \mathbf{u}(t_m)$  as

$$76 \quad (2.6) \quad \mathbf{u}^m = \mathbf{u}^{m-1} + \tau \sum_{i=1}^s b_i \mathbf{k}_i^m,$$

77 where the vectors  $\mathbf{k}_1^m, \dots, \mathbf{k}_s^m \in \mathbb{R}^n$  are the solutions of the linear system

$$78 \quad (2.7) \quad \left( \begin{bmatrix} I & & \\ & \ddots & \\ & & I \end{bmatrix} - \frac{\tau}{h^2} \begin{bmatrix} a_{1,1}L & \dots & a_{1,s}L \\ \vdots & \ddots & \vdots \\ a_{s,1}L & \dots & a_{s,s}L \end{bmatrix} \right) \mathbf{k}^m = \begin{bmatrix} \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{ST})}(t_{m-1} + c_1 \tau) \\ \vdots \\ \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{ST})}(t_{m-1} + c_s \tau) \end{bmatrix},$$

79 with

$$80 \quad \mathbf{k}^m := [\mathbf{k}_1^m \quad \dots \quad \mathbf{k}_s^m]^T \in \mathbb{R}^{ns}.$$

81 Using the Kronecker product formulation (denoted by  $\otimes$ ; see [29] and references  
82 therein), (2.7) becomes

$$83 \quad (2.8) \quad \underbrace{\left( I_s \otimes I_n - \frac{\tau}{h^2} (A \otimes L) \right)}_{=: M} \mathbf{k}^m = \begin{bmatrix} \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{ST})}(t_{m-1} + c_1 \tau) \\ \vdots \\ \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{ST})}(t_{m-1} + c_s \tau) \end{bmatrix}.$$

84 We note that (2.8) can be reformulated into a *matrix equation*, which is in general  
85 better suited for using a Krylov solver (see [22]). Here we focus on the analysis of  
86 the results in [23] and thus we do not address this any further but a study of the  
87 preconditioners from [23] in the matrix equations setting seems worthwhile. Having  
88  $p \leq 2s$  as the order of convergence of the IRK method we assume that it is balanced  
89 with the spatial discretization error, i.e., that  $h^2 = C_e \tau^p$  for some  $C_e > 0$ .

90 The problem (2.8) with the sparse system matrix  $M$  can be very large for  $h$  (and  
91  $\tau$ ) small, suggesting an iterative solver such as GMRES, BiCG or GCR should be  
92 used, which in turn requires a preconditioner to attain efficiency. In [23], the authors  
93 introduce the block preconditioners

$$94 \quad (2.9) \quad \begin{aligned} P^d &= I_s \otimes I_n - \frac{\tau}{h^2} \text{diag}(A) \otimes L, \\ P^u &= I_s \otimes I_n - \frac{\tau}{h^2} D_A U_A \otimes L \quad \text{and} \quad P^l = I_s \otimes I_n - \frac{\tau}{h^2} L_A D_A \otimes L, \end{aligned}$$

95 where  $L_A, D_A, U_A$  are the LDU factors of the Butcher tableau matrix  $A$ . In addition,  
96 the authors also consider the block triangular preconditioners

$$97 \quad (2.10) \quad P^{\text{GSL}} = I_s \otimes I_n - \frac{\tau}{h^2} A_L \otimes L \quad \text{and} \quad P^{\text{GSU}} = I_s \otimes I_n - \frac{\tau}{h^2} A_U \otimes L,$$

98 where  $GSL/GSU$  stands for *Gauss-Seidel lower/upper*, and  $A_{L,U}$  is the lower/upper  
 99 triangular part of  $A$ , i.e.,

$$100 \quad (A_L)_{ij} = \begin{cases} a_{ij} & \text{if } i \geq j \\ 0 & \text{otherwise} \end{cases}, \quad (A_U)_{ij} = \begin{cases} a_{ij} & \text{if } i \leq j \\ 0 & \text{otherwise} \end{cases}.$$

101 Some of these –  $P^d$  and  $P^{GSL}$  – were considered already in [28]. Notice that if  $a_{ii} > 0$   
 102 for all  $i = 1, \dots, s$ , then the preconditioners are invertible as  $L$  is symmetric, negative-  
 103 definite. More general conditions for non-singularity of the preconditioners can be also  
 104 derived analogously to [27, Lemma 1].

105 Using GMRES for a linear system  $C\mathbf{x} = \mathbf{f}$  with  $C$  being diagonalizable, i.e.,  
 106  $C = \Lambda S^{-1}$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ , a standard convergence bound for the residuals  
 107  $\mathbf{r}_\ell$  reads

$$108 \quad (2.11) \quad \frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \kappa(S) \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{1 \leq i \leq d} |\varphi(\lambda_i)|,$$

109 where  $\kappa(S)$  is the 2-norm condition number of the matrix  $S$ , see, e.g., [18, Section  
 110 5.7.2]. We highlight some aspects of the bound (2.11) that are often used to study  
 111 GMRES convergence behavior.

112 *Remark 2.1.* As indicated above, the spectral information of the system matrix  
 113 in GMRES (in our case of the preconditioned system) does not generally govern the  
 114 convergence (see [12], [11] and [1] and also [18, Chapter 2 and 5.7] and the references  
 115 therein). If the system matrix is normal, i.e., it is diagonalizable with  $S$  unitary,  
 116 then the spectral information is enough to evaluate the ideal GMRES bound (2.11).  
 117 However, if  $C$  is non-normal, then a convincing argument needs to be put forward to  
 118 validate linking spectral information with the convergence behavior of GMRES as the  
 119 authors in [18, p. 303, Remark 1] point out.

120 Moreover, particular knowledge of the interaction of  $S$  and the initial residual  $\mathbf{r}_0$   
 121 can lead to a *qualitative* and *quantitative* improvement on (2.11), see, e.g., [17]. How-  
 122 ever, studying GMRES behavior with the bound (2.11), this interaction is completely  
 123 lost.

124 In cases where (2.11) is justifiable, the next step is usually to bound from above  
 125 the mixed<sup>1</sup> min-max problem in the right-hand side of (2.11) by replacing the discrete  
 126 set over which we take the maximum, let us denote it by  $\sigma^{\text{discr}}$ , by a non-discrete one,  
 127 which we denote by  $\sigma^{\text{non-discr}}$ , so that we have  $\sigma^{\text{discr}} \subset \sigma^{\text{non-discr}}$ . We highlight two  
 128 important aspects of this step:

- 129 (a) It is *functional* only if we can further bound or evaluate the solution of the  
 130 min-max problem over  $\sigma^{\text{non-discr}}$  and obtain a reasonably fast convergence  
 131 estimate.
- 132 (b) It is *appropriate* only if<sup>2</sup>  $\partial_{\mathbb{C}} \sigma^{\text{non-discr}}$  is reasonably uniformly covered by  
 133  $\sigma^{\text{discr}}$ .<sup>3</sup> In case of clusters, we should consider having  $\sigma^{\text{non-discr}}$  as a union

<sup>1</sup>Mixed in the sense that the minimum is over a non-discrete set while the maximum is over a discrete one.

<sup>2</sup>We denote the boundary of a set  $S \subset \mathbb{C}$  in  $\mathbb{C}$  by  $\partial_{\mathbb{C}} S$ .

<sup>3</sup>Intuitively, we could expect that the bound will be appropriate only if  $\sigma^{\text{discr}}$  covers the entirety of  $\sigma^{\text{non-discr}}$  but because polynomials of complex variables are harmonic we can conclude that the maximum of the modulus of a polynomial over the set  $\sigma^{\text{non-discr}}$  is attained along  $\partial_{\mathbb{C}} \sigma^{\text{non-discr}}$  and therefore only the relation of  $\partial_{\mathbb{C}} \sigma^{\text{non-discr}}$  and  $\sigma^{\text{discr}}$  is important for the GMRES bound, see [5, Section 2].

of separate non-discrete sets  $\sigma_i^{\text{non-discr}}$  each of which captures one of the clusters, i.e., is covered by one of the clusters reasonably uniformly. For example, in (2.11) we can replace the spectrum  $\sigma^{\text{discr}} = \{\lambda_1, \dots, \lambda_d\}$  by a disc containing all of the eigenvalues  $\sigma^{\text{non-discr}} = \{z \in \mathbb{C} \mid |z - c| \leq \eta\}$ . Assuming  $|c| > \eta$ , a crude but sometimes useful approximation of the original bound is available,

$$(2.12) \quad \frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \kappa(S) \left( \frac{\eta}{|c|} \right)^k,$$

see [25, Section 6.11.2, Corollary 6.33 and Lemma 6.26 and below]. Here,  $\sigma^{\text{non-discr}} = \{z \in \mathbb{C} \mid |z - c| \leq \eta\}$  was clearly chosen with the *functionality* aspect in mind as we know the polynomial that realizes the bound (see [25, Lemma 6.26]) and it gives a good convergence bound as long as  $\eta \not\approx |c|$ . However, it is usually far from being *appropriate* if the eigenvalues don't spread uniformly over the circle bounding the disc. One notable exception is the case of tightly clustered eigenvalues around a single point  $c$  – in this case the clustering usually makes this bound appropriate as we can choose  $\eta$  *very* small. We emphasize that the adjectives *functional* and *appropriate* make sense only if the original bound (2.11) was itself descriptive of the GMRES convergence bound, i.e., only if the system matrix is either close to normal or the initial residual is restricted to a subspace on which the system matrix is not too far from being normal.

**3. Analysis of the block preconditioners.** We start by transforming the calculations into the eigenbasis of the spatial operator. Denoting the eigenpairs of  $L$  by  $(\lambda_k, \mathbf{v}_k)$ , we organize the eigenvectors into an  $n$ -by- $n$  matrix  $V$  and define the block transformation matrix  $Q$ ,

$$(3.1) \quad V := [\mathbf{v}_1, \dots, \mathbf{v}_n], \quad \text{and} \quad Q := \begin{bmatrix} V & & \\ & \ddots & \\ & & V \end{bmatrix} \in \mathbb{R}^{sn \times sn}.$$

Transforming  $M$  blockwise into the  $V$  basis gives  $\tilde{M} := QMQ^T$ ,

$$(3.2) \quad \tilde{M} = \begin{bmatrix} I & & \\ & \ddots & \\ & & I \end{bmatrix} - \frac{\tau}{h^2} \begin{bmatrix} a_{1,1}\Lambda & \dots & a_{1,s}\Lambda \\ \vdots & \ddots & \vdots \\ a_{s,1}\Lambda & \dots & a_{s,s}\Lambda \end{bmatrix},$$

with  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . With the preconditioners proposed in (2.9-2.10) we write the spectrum of the preconditioned system as

$$\text{sp}(MP^{-1}) = \text{sp}(Q^T MP^{-1}Q) = \text{sp}(Q^T MQQ^T P^{-1}Q) = \text{sp}(\tilde{M}\tilde{P}^{-1}),$$

where  $\tilde{P} := Q^T PQ$  stands for one of the right-preconditioners  $P^{\text{d,GSU,u}}$  and an analogous formulation follows also for the left-preconditioners  $P^{\text{GSL,l}}$ . As the preconditioners are defined blockwise as scalar multiplications of  $L$  and  $I$ , their blockwise transformation into the eigenbasis of  $L$  is a straight-forward calculation – replacing  $L$  with  $\Lambda$  (and keeping  $I$ ). Next, such matrices – block matrices with each block being a square, diagonal matrix – can be permuted into classical block-diagonal matrices as the following lemma shows.

LEMMA 3.1 (see [10, Lemma 1]). *Let  $C \in \mathbb{R}^{ns \times ns}$  be a real matrix with block*

169 structure such that every block is a square diagonal matrix, i.e.,

$$170 \quad (3.3) \quad C = \begin{bmatrix} \Lambda_{11} & \dots & \Lambda_{1s} \\ \vdots & \ddots & \vdots \\ \Lambda_{s1} & \dots & \Lambda_{ss} \end{bmatrix}, \quad \text{with} \quad \Lambda_{ij} = \text{diag} \left( \lambda_1^{(ij)}, \dots, \lambda_n^{(ij)} \right) \quad \forall ij.$$

171 Then there exists a permutation matrix  $\Pi \in \mathbb{R}^{ns \times ns}$  such that

$$172 \quad (3.4) \quad \Pi^T C \Pi = \begin{bmatrix} C_1 & & \\ & \ddots & \\ & & C_n \end{bmatrix} \quad \text{with} \quad C_\ell = \begin{bmatrix} \lambda_\ell^{(11)} & \dots & \lambda_\ell^{(1s)} \\ \vdots & \ddots & \vdots \\ \lambda_\ell^{(s1)} & \dots & \lambda_\ell^{(ss)} \end{bmatrix} \in \mathbb{R}^{s \times s},$$

173 for any  $\ell = 1, \dots, n$ .

174 Hence,  $C$  is diagonalizable if and only if  $C_\ell$  is diagonalizable for all  $\ell = 1, \dots, n$ ,  
 175 and if  $C_\ell = V_\ell^{-1} D_\ell V_\ell$  is the eigendecomposition of  $C_\ell$  with  $D_\ell = \text{diag}(\mu_\ell^{(1)}, \dots, \mu_\ell^{(s)})$ ,  
 176 then

$$177 \quad \text{sp}(C) = \bigcup_{\ell=1}^n \bigcup_{i=1}^s \mu_\ell^{(i)},$$

178 and if  $(\mu, \mathbf{v})$  is an eigenpair of some  $C_\ell$ , then  $(\mu, \Pi^T (\mathbf{v} \otimes \mathbf{e}_\ell))$  is an eigenpair of  $C$ .  
 179 As a result, if  $C$  is diagonalizable with  $C = V^{-1} D V$ , then

$$180 \quad \kappa(V) = \frac{\max_{\ell=1, \dots, n} \sigma_1^{(\ell)}}{\max_{\ell=1, \dots, n} \sigma_s^{(\ell)}},$$

181 where  $\kappa(\cdot)$  is the 2-norm condition number and the matrices  $V_\ell$  have the singular  
 182 values  $\sigma_1^{(\ell)} \geq \dots \geq \sigma_s^{(\ell)} \geq 0$ .

183 *Remark 3.2.* We note that an analogous lemma to Lemma 3.1 can also be for-  
 184 mulated for non-normal matrices (replacing  $Q^T$  by  $Q^{-1}$ ). Considering the Jordan  
 185 canonical (or the Schur decomposition form) of  $C_\ell$ , Lemma 3.1 can be reformulated  
 186 to obtain a block upper bi-diagonal (or block upper-triangular) matrix.

187 We take  $W$  as the matrix of eigenvectors of  $L$ , and in order to shorten the notation  
 188 we set

$$189 \quad (3.5) \quad \theta_k := \frac{\tau}{h^2} \lambda_k \quad \text{and} \quad \Theta := \frac{\tau}{h^2} \Lambda,$$

190 as these quantities always appear together in the computations, and we use  $p$  as the  
 191 order of the Runge-Kutta scheme (see [31, Section II.1, Definition 1.2]). Assuming  
 192 the time and space discretization errors are kept in balance, i.e., there exists a  $C$  so  
 193 that  $h^2 = C\tau^p$ , a direct calculation (see [21, Appendix B.8, pages 228–229]) leads us  
 194 to the following limit behavior of  $\theta_k$  as  $\tau, h \rightarrow 0$ :

$$195 \quad (3.6) \quad \underbrace{(\theta_n, \theta_1) \rightarrow \left(-\frac{8}{C_e}, 0\right), \quad (\theta_1^{-1}, \theta_n^{-1}) \rightarrow \left(-\infty, -\frac{C_e}{8}\right)}_{(\text{LIM})_{p=1}}, \quad \underbrace{(\theta_n, \theta_1) \rightarrow (-\infty, 0), \quad (\theta_1^{-1}, \theta_n^{-1}) \rightarrow (-\infty, 0)}_{(\text{LIM})_{p>1}}.$$

196 Next we define the  $s$ -by- $s$  matrices

$$197 \quad M_k := \begin{bmatrix} 1 - a_{11}\theta_k & -a_{12}\theta_k & \cdots & -a_{1s}\theta_k \\ -a_{21}\theta_k & 1 - a_{22}\theta_k & & \vdots \\ \vdots & & \ddots & \vdots \\ -a_{s1}\theta_k & \cdots & \cdots & 1 - a_{ss}\theta_k \end{bmatrix} \quad \text{and} \quad P_k^* := \begin{bmatrix} 1 - \alpha_{11}\theta_k & -\alpha_{12}\theta_k & \cdots & -\alpha_{1s}\theta_k \\ -\alpha_{21}\theta_k & 1 - \alpha_{22}\theta_k & & \vdots \\ \vdots & & \ddots & \vdots \\ -\alpha_{s1}\theta_k & \cdots & \cdots & 1 - \alpha_{ss}\theta_k \end{bmatrix},$$

198 where  $\alpha_{ij}$  are the entries of the replacement for  $A$  in  $M$ , e.g., taking  $\star = d$  we have  
 199  $\alpha_{ij} = a_{ij}$  for  $i = j$  and  $\alpha_{ij} = 0$  otherwise, while taking  $\star = u$  we have  $\alpha_{ij} = (D_A U_A)_{ij}$   
 200 where  $A = L_A D_A U_A$  is the LDU factorization of  $A$  and so on. Using Lemma 3.1, we  
 201 obtain the following result.

202 **PROPOSITION 3.3.** *Take  $M$  as in (2.8) and a preconditioner  $P$  from (2.9, 2.10).  
 203 Assuming  $P$  is invertible, the spectrum of  $MP^{-1}$  (or  $P^{-1}M$ ) is given as the union of  
 204 the spectra of the matrices  $X_k$  given by*

$$205 \quad (3.7) \quad X_k^\star := M_k (P_k^\star)^{-1} \quad (\text{or } (P_k^\star)^{-1} M_k),$$

206 for  $k = 1, \dots, n$ . If all  $X_k^\star$  are diagonalizable with

$$207 \quad (3.8) \quad (S_k^\star)^{-1} X_k^\star S_k^\star = \text{diag}(\xi_1^{(k)}, \dots, \xi_s^{(k)}),$$

208 then the condition number of the matrix of the eigenvectors of the preconditioned  
 209 system is given by

$$210 \quad \kappa(W) \cdot \max_{k=1, \dots, n} \kappa(S_k^\star).$$

211 If the  $\theta_k$  have multiplicity at most  $m$ , then the eigenvalues of the preconditioned system  
 212 have algebraic multiplicity at most  $ms$ . In particular, the preconditioned system can  
 213 be non-diagonalizable but the longest Jordan vector chain has length at most  $ms$ .

214 *Proof.* Transforming  $MP^{-1}$  (or  $P^{-1}M$ ) into the basis of  $Q$  we use Lemma 3.1  
 215 for the matrix  $\tilde{M}\tilde{P}^{-1}$  (see (3.2)) and obtain the result.  $\square$

216 Now we are ready to generalize the results shown in [10] for  $s = 2$  to a general  $s$ -stage  
 217 method.

218 **COROLLARY 3.4** ([21, Proposition 7.5]). *Under the assumptions of Proposi-  
 219 tion 3.3, we have for the right-preconditioner  $P^d$  the formula*

$$220 \quad (3.9) \quad X_k^d = \begin{bmatrix} 1 & -\frac{a_{12}\theta_k}{1-a_{22}\theta_k} & \cdots & -\frac{a_{1s}\theta_k}{1-a_{ss}\theta_k} \\ -\frac{a_{21}\theta_k}{1-a_{11}\theta_k} & 1 & & \vdots \\ \vdots & & \ddots & \vdots \\ -\frac{a_{1s}\theta_k}{1-a_{11}\theta_k} & \cdots & \cdots & 1 \end{bmatrix},$$

221 with the characteristic polynomial

$$222 \quad p_k^{(s)}(\lambda) = (1 - \lambda)^s + \beta_{s-2}(1 - \lambda)^{s-2} + \beta_{s-3}(1 - \lambda)^{s-3} + \dots + \beta_1(1 - \lambda) + \beta_0,$$

223 where  $\beta_j$  are continuous functions of  $\theta_k$  and  $a_{ii}$  for  $i = 1, \dots, s$ . Hence, the eigenvalues  
 224 become  $1 - \mu$ , where  $\mu$  is a root of the parametrized polynomial

$$225 \quad \tilde{p}_k^{(s)}(t) = t^s + \beta_{s-2}t^{s-2} + \beta_{s-3}t^{s-3} + \dots + \beta_1t + \beta_0.$$

COROLLARY 3.5 ([21, Proposition 7.6]). *Under the assumptions of Proposition 3.3, the block upper-triangular preconditioners  $P^{\text{GSU},u}$  give*  
(3.10)

$$X_k^{\text{GSU},u} = \begin{bmatrix} 1 & 0 & \dots & \dots & \dots & 0 \\ \star & & & & & \\ \vdots & (M_k(P_k^{\text{GSU},u})^{-1})_{2:s,2:s} & & & & \\ \vdots & & & & & \\ \star & & & & & \end{bmatrix}, \quad X_k^{\text{GSL},l} = \begin{bmatrix} 1 & \star & \dots & \dots & \dots & \star \\ 0 & & & & & \\ \vdots & ((P_k^{\text{GSL},l})^{-1}M_k)_{2:s,2:s} & & & & \\ \vdots & & & & & \\ 0 & & & & & \end{bmatrix},$$

and hence have one eigenvalue equal to one for each  $k$ . The entries replaced by  $\star$  above do not affect the spectrum, only the eigenbasis.

These results suggest 1 as a natural “central point” of the spectrum of the preconditioned system, generalizing the observations made for  $s = 2$ . We note that using these results we get both quantitative and qualitative insight into the spectra shown in [23, Figure 4.1 – 4.4], e.g., we see that for  $s = 3$  the eigeninformation of  $M(P^u)^{-1}$  and  $(P^l)^{-1}M$  can still be obtained explicitly (see also [21, Section 7.4]) and on the other hand for  $s \geq 6$  there is no hope for these in general – but any bound on the eigeninformation of  $L$  can be used to obtain a bound on the eigeninformation of the preconditioned system by calculating with  $X_k$ , see [10, Section 4].

We show the spectra of the preconditioned systems and the corresponding GMRES convergence behavior in Figure 1 and 2, demonstrating observations and results from above. Notably, the bounds leave something to be desired, especially for  $P^d$  where they are not descriptive at all. Moreover, increasing  $s$  seems to noticeably affect the quality of the preconditioners – see also [23] for further numerical tests with various  $s$  and  $h$ . These numerical examples (as well as the ones in [3, 10]) are, as far as we can tell, representative of the general experience with these preconditioners. We highlight several key features illustrated in Figures 1 and 2 that remained true in all of our experiments:

1. For  $s$  small, we have observed the staircase-like convergence behavior visible in the left upper-most plot in Figure 2 (and also in the first row of Figure 5), where GMRES makes very little progress for a number of iterations, then improves notably in one iteration and repeats this cycle going forward. This behavior was most pronounced for the preconditioner  $P^d$ , and for  $s = 2$  was described and explained in [10, Figure 2 and below].
2. We have usually not observed the desired *superlinear* convergence behavior, except for a speed-up after an initial stagnation (or slower speed convergence) phase.
3. In the vast majority of cases, the number of GMRES iterations to reach a certain tolerance grows only very moderately under mesh refinement and for  $P^u, P^l$  it remains almost constant.
4. In all of the experiments the spectra had the characteristic arc-like structure that we see in Figure 1.

Our goal is to explain all these features here as well as to investigate other bounds or estimates that would be more descriptive of the convergence behavior. This insight is of clear interest on its own but can be also used to further improve the used methods, e.g., looking at *numerical optimization* of the Butcher tableau in the spirit of [10, Section 4]. We also note that the above results translate in a straight-forward fashion to the *transformed system* after we multiply (2.8) with  $(A^{-1} \otimes I_n)$  from the



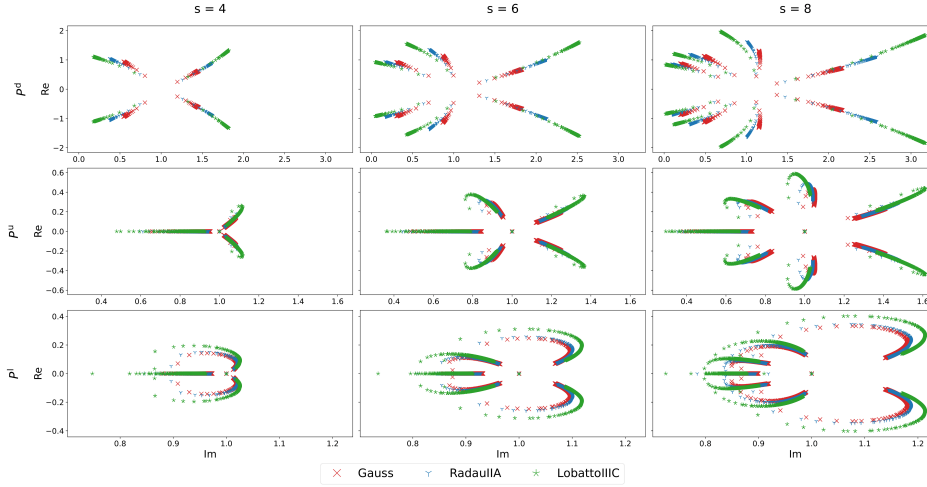


FIG. 1. The spectra of the preconditioned systems  $M (P^{u,d})^{-1}$  and  $(P^l)^{-1} M$  for  $s = 4, 6, 8$  and for three classical choices of fully implicit Runge-Kutta schemes - Gauss, RadauIIA and LobattoIIIC. The spectra seemingly assemble in  $s$  “branches” in the first row and into  $s - 1$  “branches” in the other two with a central point at  $1 + 0i$ . We set  $N = 50$ .

left, obtaining

$$\underbrace{\left( A^{-1} \otimes I_n - \frac{\tau}{h^2} I_s \otimes L \right)}_{=: M^{\text{transf}}} \mathbf{k}^m = (A^{-1} \otimes I_n) \begin{bmatrix} \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{BC})}(t_{m-1} + c_i \tau) \\ \vdots \\ \frac{1}{h^2} L \mathbf{u}^{m-1} + \mathbf{b}^{(\text{BC})}(t_{m-1} + c_i \tau) \end{bmatrix},$$

and getting analogously the preconditioners,

$$\begin{aligned} R^d &= \text{diag}(A^{-1}) \otimes I_n - \frac{\tau}{h^2} I_s \otimes L, \\ R^l &= (D_{A^{-1}} U_{A^{-1}}) \otimes I_n - \frac{\tau}{h^2} I_s \otimes L \quad \text{and} \quad R^u = (L_{A^{-1}} D_{A^{-1}}) \otimes I_n - \frac{\tau}{h^2} I_s \otimes L, \\ R^{\text{GSL}} &= (A^{-1})_L \otimes I_n - \frac{\tau}{h^2} I_s \otimes L \quad \text{and} \quad R^{\text{GSU}} = (A^{-1})_U \otimes I_n - \frac{\tau}{h^2} I_s \otimes L, \end{aligned}$$

where  $A^{-1}$  has the LDU factorization  $A^{-1} = L_{A^{-1}} D_{A^{-1}} U_{A^{-1}}$  and  $(A^{-1})_{L,U}$  are defined analogously to (2.10). These preconditioners were proposed in [20] and then used further in [19] but also [27, 26]. For a *general* Butcher tableau, it is not possible to say whether the preconditioned transformed system gives a better performance than the original one. However, in [27, 26] the authors propose different preconditioners and our analysis adapted to their framework is going to be considered elsewhere. Also, we note that the extension of the above analysis for FEM discretization is a straightforward task – more details on both of these topics can be found in [21, Sections 7.6 and 7.7].

**3.1. Spectral analysis.** Next we turn to the spectral analysis, keeping in mind its limitation in the sense of Remark 2.1. For block-diagonal problems we obtain

$$(3.11) \quad \frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{j=1, \dots, n} \|\varphi(X_j)\|,$$

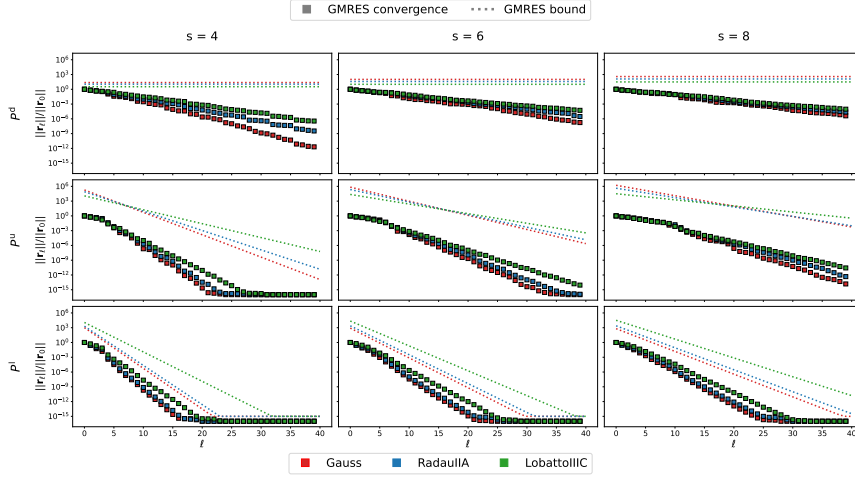


FIG. 2. The preconditioned GMRES convergence behavior for the preconditioned systems  $M(P^{u,d})^{-1}$  and  $(P^l)^{-1}M$  for  $s = 4, 6, 8$  and three classical choices of fully implicit Runge-Kutta schemes - Gauss, RadauIIA and LobattoIIIC - together with the GMRES bound (2.12) with  $c = 1$  (we set the values to 1 if  $\eta \geq 1$ ). We set  $N = 50$ .

which was studied in [9], where the authors showed that the extremal polynomials (i.e., the polynomial realizing the above bound) satisfies the equioscillation property but only every  $s$  iterations, where  $s$  is the size of the diagonal blocks. Relabeling the blocks in (3.11) we get

$$\frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{j=1, \dots, n} \|\varphi(X_j)\| = \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{\theta_j \in \text{sp}(\frac{\tau}{h^2}L)} \|\varphi(X_{\theta_j})\|.$$

Assuming each  $X_{\theta_j}$  is diagonalizable as in Proposition 3.3, we notice that  $\{\theta_j\}$  covers reasonably well the intervals  $I_{h,\tau,\dots}$  as  $h \rightarrow 0$  (see (3.6)) and, in the spirit of Section 2, the natural bound of (3.11) becomes

$$\frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{\theta \in I_{h,\tau,\dots}} \|\varphi(X_\theta)\|.$$

First, let us assume there is a uniform bound  $\kappa(S_\theta) \leq \kappa_S$  for all  $\theta \in I_{h,\tau,\dots}$ , which experimentally seems to be the case (see [21]) and can be confirmed analytically for  $s = 2, 3$  (see [10]) – this is an important and non-trivial assumption and a proper justification is an open problem. Next, we notice that the matrices  $X_\theta$  depend *smoothly*<sup>4</sup> on  $\theta$  and as a result so do their eigenproperties. In particular, the eigenvalues  $\xi_\theta^{(i)}$  of  $X_\theta$  will – by definition – form an *algebraic curve*<sup>5</sup> with  $s$  arcs (sometimes also called *branches*) some of which can be degenerate, e.g., reduced to just a point (incidentally,

<sup>4</sup>That is, for our model problem of the negative-definite Laplacian. However in most cases of interest this assumption is also satisfied, partially due to the stability assumptions/conditions coming from the Runge-Kutta scheme.

<sup>5</sup>We say that  $\Gamma$  is an algebraic curve provided there exists a bi-variate polynomial  $p(\theta, \xi)$  such that  $\Gamma = \{(\theta, \xi) | p(\theta, \xi) = 0\}$ . Locally, this can also be viewed through the lens of perturbation theory, see [14, Chapter 2 Section 1.1].

this is the case for at least one arc of the algebraic curve for any of the triangular preconditioners due to Corollary 3.5). Denoting the algebraic curve for the given Butcher tableau  $A$  and a choice of preconditioner  $P^*$  by  $\Gamma$ , we obtain

$$(3.12) \quad \frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \leq \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{\theta \in I_{h,\tau,\dots}} \kappa(S_\theta) \max_{i=1,\dots,s} \left| \varphi\left(\xi_\theta^{(i)}\right) \right| \leq \kappa_S \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{\xi \in \Gamma} |\varphi(\xi)|.$$

Notice that if we replace in (3.12) the interval  $I_{h,\tau,\dots}$  with its limit  $I_{\text{lim}}$  as  $h, \tau \rightarrow 0$  (see (3.6)), we obtain a bound for all mesh sizes. Noticing that, in our case, the preconditioned system matrix has a limit as  $\theta$  tends to either of the endpoints of  $I_{\text{lim}}$ , it follows that the arcs of the corresponding algebraic curve correspond to the eigenvalues of these limit matrices. Hence, the effect of mesh refinement becomes sampling more points along  $\Gamma$  and stretching it towards these fixed endpoints (and possibly in increasing  $\kappa_S$ ). This suggests that from a certain mesh size onward, the mesh refinement will have little effect on  $\Gamma$  and hence will not affect the min-max part of (3.12), shedding some light on why these preconditioners are quite robust under mesh refinement.

*Remark 3.6.* Note that the numerical experiments in [23, 3] as well as in [21] and in Section 4 clearly show that the spectra of the preconditioned systems cover reasonably well an algebraic curve. For two-stage methods, this behavior has been observed, proved and used to obtain descriptive GMRES bounds in [10]. Moreover, for any algebraic curve  $\Gamma$  we have  $\Gamma = \partial_{\mathbb{C}}\Gamma$ , which is convenient from the point of view of choosing  $\sigma^{\text{non-discr}}$ , see Remark 2.1 and below.

We also emphasize that, in general, these preconditioners do not cluster eigenvalues (that is, any more than the  $\theta \in I_{h,\tau,\dots}$  already are) but rather place them along a particular algebraic curve  $\Gamma \subset \mathbb{C}$ . Hence, if the conditioning of the eigenbasis is not very bad, we can reasonably expect linear convergence as opposed to superlinear, which can often be linked with clusters and numbers of outliers, in the sense of [18, Section 5.6.4].

Remark 3.6 also explains that the bound (2.12) is unlikely to be very descriptive or even usable. Indeed, the algebraic curves can reach into the left half-plane  $\{\text{Re}(z) < 0\}$  (making the bound useless due to 0 being included in the bounding circle) or, in the more favorable case, the arcs of the algebraic curve are *extremely* unlikely to align with the circle so that the bound have some resemblance of being what we earlier called *appropriate*. Naturally, the bound on the right-hand side of (3.12) is constructed to remedy that but the key question becomes if this bound is also *functional*, namely if we can (approximately) evaluate it.

To this end, we follow the excellent paper [5] on this topic and start by looking at the *asymptotic* convergence factor (justified by Remark 3.6 above). Considering (3.12) we are led to look at the so-called *logarithmic capacity* of  $\Gamma$ , denoted by  $\text{cap}(\Gamma)$ , which can be viewed as a measure of a compact set without isolated points in  $\mathbb{C}$ ; see [24, 13, 5] for the definition and further reading, but also [2] for progress on the calculation of logarithmic capacities. Importantly,  $\text{cap}(\Gamma)$  is known to asymptotically correspond to the maximal modulus of the *extremal polynomials* (sometimes also called Chebyshev polynomials) associated with  $\Gamma$ , namely

$$(3.13) \quad \left( \min_{\deg(\varphi) \leq \ell} \max_{z \in \Gamma} |\varphi(z)| \right)^{1/\ell} \rightarrow \text{cap}(\Gamma), \quad \text{as } \ell \rightarrow +\infty,$$

where the quantity on the left-hand side relates to the quantities we have seen in the

GMRES bounds. There are two important caveats to using  $\text{cap}(\Gamma)$ . The first one, which has been also highlighted as a caveat for using the analysis in [5] overall, is the fact that (3.13) only provides some information about the *limit behavior* as  $\ell \rightarrow +\infty$ , whereas we are interested in the behavior for relatively small values of  $\ell$ , say  $\ell \leq 50$  or 100. To large extent this issue is addressed by Remark 3.6 that states that we expect a linear convergence throughout the iteration. The second one is the fact that (3.13) describes the limit scaling of the maximal modulus over *all polynomials* – it lacks the crucial scaling  $\varphi(0) = 1$  of Krylov methods. This issue can be fixed by re-scaling (see [5, Section 2]), shifting our attention from the logarithmic capacity to *Green's functions associated with*  $\Gamma$ , as long as  $\Gamma$  is compact and without any isolated points.

Things simplify considerably if we assume that  $\Gamma$  is connected as then the normalized quantity

$$\left( \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{z \in \Gamma} |\varphi(z)| \right)^{1/\ell}$$

can be evaluated directly using conformal maps, in particular the Schwarz-Christoffel maps. Without going into the details (the interested reader can find these in [5, Sections 2 and 3]), we obtain the *asymptotic convergence factor estimate*  $\rho_{\text{est}}$  as

$$(3.14) \quad \rho_{\text{est}} := \lim_{\ell \rightarrow +\infty} \left( \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \max_{z \in \Gamma} |\varphi(z)| \right)^{1/\ell} = \frac{1}{|\Phi(0)|},$$

where  $\Phi(z)$  is the Schwarz-Christoffel map that maps the exterior of  $\Gamma$  to the exterior of the unit circle. In [5, Section 3, Theorem 2 and below], the authors put this as

“... if  $\Gamma$  is connected, the estimated asymptotic convergence factor for a matrix iteration depends on how far the origin is from  $\Gamma$  – provided that this distance is measured by level curves associated with the exterior conformal map.”

We would like to emphasize the word *estimate* when talking about  $\rho_{\text{est}}$  because we truly do not get a bound anymore – in fact we get an *underestimate* as highlighted also in [5, Section 5, equation (STEP1) and also Table 1]. However, we expect this estimate to be descriptive as explained above.

For not too complicated connected, compact sets the map  $\Phi$  and its value at the origin can be calculated using the Schwarz-Christoffel MATLAB toolbox [4], but we immediately notice that in Figure 1 the set of eigenvalues along  $\Gamma$  is not connected and the actual algebraic curve  $\Gamma$  itself is also not available in an easy form, i.e., neither of these can be directly given as an input to the SC toolbox. We take the natural next step and approximate  $\Gamma$  by its linear interpolation based on the available eigenvalues  $\xi_{\theta}^{(i)}$ . The linear interpolation gives us a good approximation of the arcs of  $\Gamma$  and we use the point  $1 + 0i$  as the natural point to join them (also by linear interpolation) and denote the resulting set  $\Gamma_h$ . Recalling the limit behavior in (3.6), we also see that  $\Gamma_h$  will tend towards  $\Gamma$  as  $h \rightarrow 0$  for our model problem.

The calculation of  $\xi_{\theta}^{(i)}$  is independent for each  $k = 1, \dots, n$  but for large  $n$  the SC toolbox can suffer numerically when calculating with  $\Gamma_h$  that is densely populated by

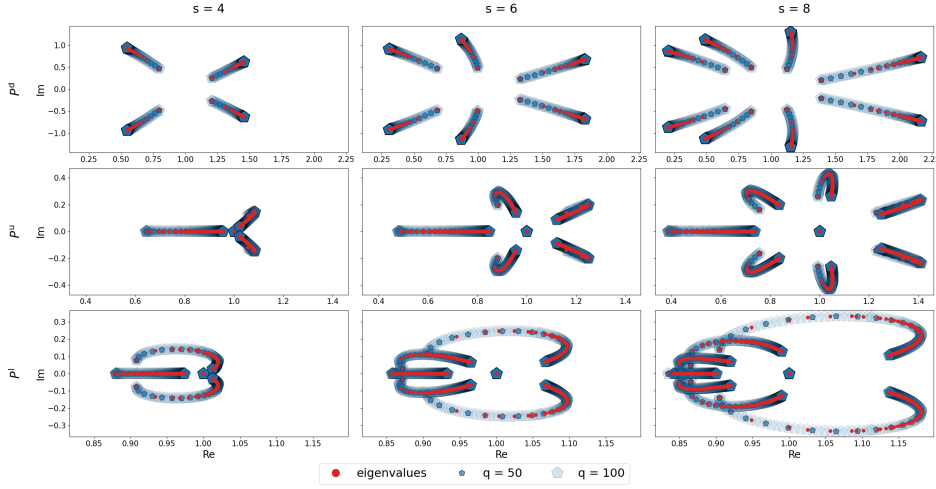


FIG. 3. The eigenvalues of the matrices  $X_{\theta_k}$  (red) and  $X_{\vartheta_k}$  (blue, for different values of  $q$ ), using the preconditioner  $P^d$ . Joining these together with line segments would yield the curves  $\Gamma_h$  (red) and  $\Gamma_q$  (blue).

the interpolation points – both in the sense of large computational complexity as well as in the sense of numerical issues (called *over-crowding*, see [4] but also [6, Section 2.6]). Moreover, we usually have only rough estimates on the extremal eigenvalues  $\theta_{\min}$  and  $\theta_{\max}$  of  $L$  rather than its full spectrum. To this end, we recall the idea in [10, Section 4] and instead of calculating  $\Gamma_h$  we use the information about  $\theta_{\min, \max}$  and artificially sample a fixed number of “fake” points  $\vartheta_k$  between them, say  $q$  of them. Then we replace  $\theta_k$  by  $\vartheta_k$  in the definition of  $\Gamma_h$ , obtaining  $\Gamma_q$  – an approximation of  $\Gamma_h$  (and a further approximation of  $\Gamma$ ) based on the linear interpolation given by the eigenvalues of the matrices  $X_{\vartheta_k}$ . We illustrate these points in Figure 3.

Another key point is that using the SC toolbox<sup>6</sup> – namely the functions `extermmap` and `evalinv` – has difficulties (as far as we understand it) when the arcs of  $\Gamma_q$  intersect, e.g., as is the case for  $s = 8$  and the preconditioner  $P^l$ , see Figure 1. Intuitively, this makes sense as the exterior of  $\Gamma_q$  then has multiple components, making the original set-up more complicated (a theoretical treatment of such problems could be approached based on [8]). We address this issue by taking the “envelope” of the arcs – if two arcs intersect, we follow the one staying outwards, e.g., in the case of  $s = 6$  (or  $s = 8$ ) and the preconditioner  $P^l$  we would exclude a portion of the densely populated end of the arc (two arcs) closer to the real axis as these portions lie “inward” relative to the arcs with the larger imaginary part, see Figure 1 and Figure 4 ahead. Finally, we illustrate the calculated Schwarz-Christoffel maps – or rather their contours – in Figure 4 together with the used inputs  $\Gamma_q$  (with the exception of  $s = 6, 8$  and the preconditioner  $P^l$ , where we used the “envelopes”) and also the asymptotic convergence factor estimate  $\rho_{\text{est}}$  in Figure 5. First, we see that the results in Figure 5 fully support the arguments in Remark 3.6 for considering  $\rho_{\text{est}}$  as *the* descriptive quantity for the convergence factor. Including an estimate for  $\kappa_S$  then gives also an estimate for GMRES convergence – not just its rate, see Section 4. Second, we note that for  $s = 8$  and the preconditioner  $P^u$ , the arcs turned so that the right-

<sup>6</sup>In our case,  $\Gamma_q$  qualifies as a degenerate polygon acceptable by the toolbox.

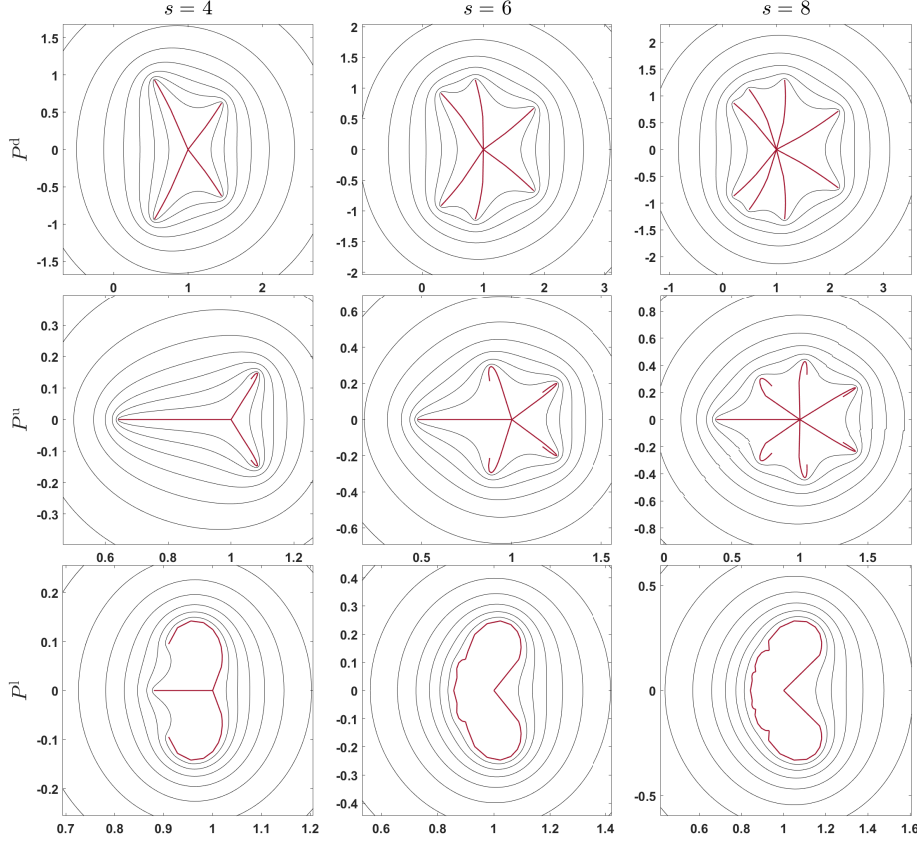


FIG. 4. In red: the curves  $\Gamma_q$  (first plots 1 to 7) and their “envelopes” (plots 8 and 9) for the Gauss Butcher tableau, taking  $q = 15$ . In black: the contours of the corresponding Schwarz-Christoffel map of the exterior of these curves (or envelopes) mapped to the exterior of the unit circle, see `extermap` in [4].

most arcs almost intersect themselves. This causes problems for the toolbox, which during the calculations raises a flag stating that the calculated map did not converge as expected. Although the predicted  $\rho_{\text{est}}$  seems accurate, we see in Figure 4 that contours have ripples, confirming that the calculated results should be taken with caution. This can be fixed by a similar “envelope-like” approach we described for  $s = 6, 8$  and the preconditioner  $P^l$ , see Section 4, obtaining a further approximation. Although there are a few similar caveats concerning the implementation of the above ideas, we have always found that a simple solution (such as considering the envelope or pruning the fake points in order to alleviate the crowding) can be used to fix them and still give an *appropriate* insight into the GMRES convergence factor. As long as  $\kappa_S$  does not completely dominate the ideal GMRES bound (2.11) this then translates to descriptive GMRES convergence estimates, see Section 4.

The above analysis also gives insight into the staircase-like behavior, which has been observed and explained for  $s = 2$  and the preconditioner  $P^d$  in [10] working with the minimal residual polynomial  $\varphi_\ell^{\text{MR}}$  (sometimes also called the GMRES polynomial; see [18, Section 5.7.1]). The arguments used in [10] remain valid as long as the

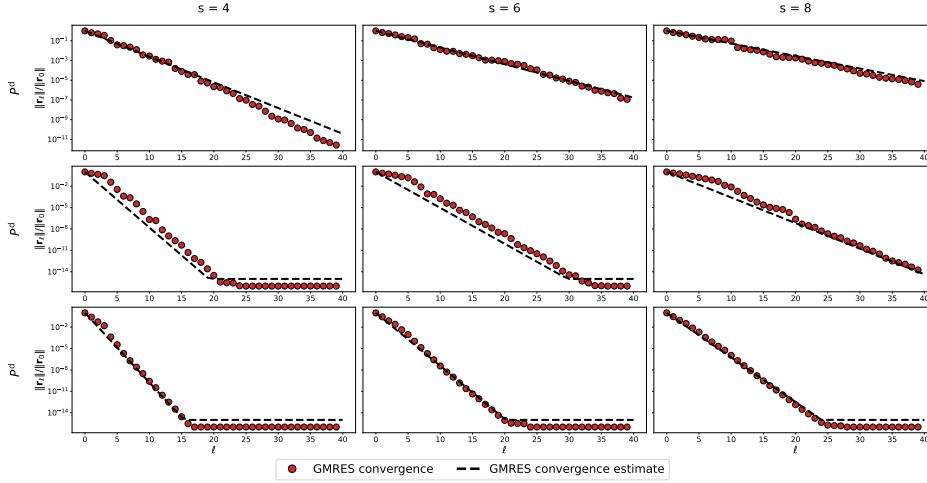


FIG. 5. The convergence behavior of preconditioned GMRES, using the Gauss Butcher tableau, together with the convergence factor estimates based on  $\rho_{\text{est}}$ .

branches are not very close to each other<sup>7</sup> – as long as the branches are far apart, the maximum of the polynomial  $\varphi_\ell^{\text{MR}}$  will decrease significantly more at the steps  $\ell = s \cdot j$  for  $j = 1, 2, \dots$  because only then each branch can get some attention. If the branches become close, then we do not expect this extra jump because keeping the absolute value of the polynomial small along one of the branches naturally translates into keeping the absolute value of the polynomial also small enough along another one. This is most pronounced in the first  $s$  iterations of GMRES, as we can see in Figure 5, where the convergence curves begin with a slower convergence phase – *precisely  $s$  steps* – for  $P^d$  and  $P^u$ , in contrast to the ones of  $P^l$ , where the arcs intersect and are, in general, closer to each other. We illustrate this further in Figure 6 for the preconditioner  $P^d$  for  $s = 4, 8$  by looking at the polynomial  $\varphi_\ell^{\text{MR}}$  and its roots (called harmonic Ritz values). We see that in the first row (4 branches, far apart) the possibility of “placing” one root along each of the branches was much more crucial (resulted in a more significant decrease of the modulus of the polynomial over the spectrum of the preconditioned system) than for the second row (8 branches with two complex conjugate pairs of branches that are close to each other). We note that an example of explanation (and prediction) of a *complete* staircase behavior of GMRES can be found in [5, Figure 9 and below].

Having analyzed the model problem, we want to emphasize that the approach relied on two assumptions – (a) the spectrum of  $L$  covers (reasonably) uniformly a real interval  $I_{\tau, h, \dots}$  and (b) the condition numbers  $\kappa(S_\theta)$  stay bounded for  $\theta \in I_{\tau, h, \dots}$ . Importantly, in *many* problems (a) is not satisfied even though the spectrum of  $L$  still shows the crucial “one-dimensionality”, i.e., the eigenvalues of  $L$  densely populate a curve  $\Psi \subset \mathbb{C}$ . To demonstrate, we consider a model problem of 1D advection-diffusion,

$$(3.15) \quad \partial_t u = (\partial_x - \kappa \partial_{xx}) u + f \quad \text{in } \mathbb{R} \times (0, T_{\text{end}}),$$

<sup>7</sup>In [10], the branches are two line segments parallel to the imaginary axis that are, moreover, reasonably well separated along the real line, i.e., a natural case of being “not very close to each other”.

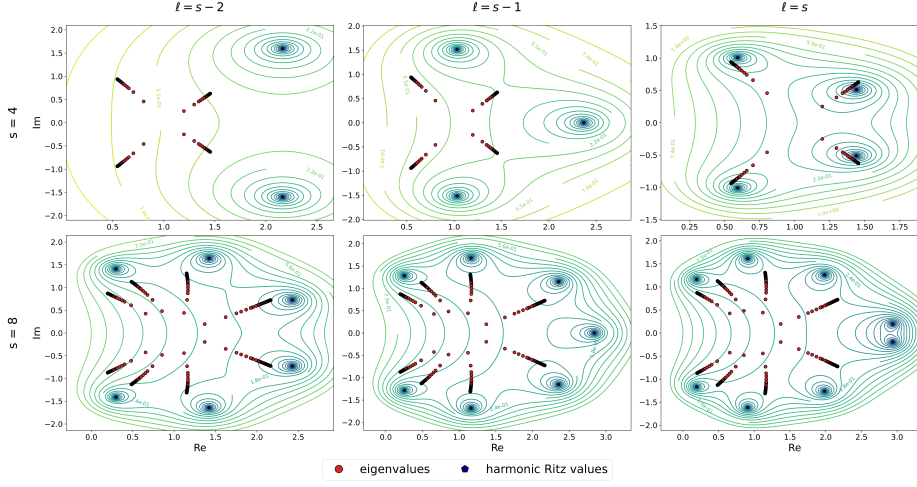


FIG. 6. The level curves of the GMRES polynomials  $\varphi_l^{MR}$  for the preconditioned system  $(P^d)^{-1}M$  together with the spectrum of this system as well as the roots of  $\varphi_l^{MR}$  (so-called harmonic Ritz values). We set  $N = 50$ .

which we discretize in space with a centered finite difference scheme with a mesh size  $h$ , obtaining an infinite tri-diagonal matrix  $L$  with the stencil<sup>8</sup>

$$\begin{bmatrix} \ddots & h - \kappa & & \\ -h - \kappa & 2\kappa & h - \kappa & \\ & -h - \kappa & \ddots & \end{bmatrix},$$

which can be in real calculations replaced by a finite matrix with, e.g., the periodic boundary conditions. To remain concise we focus only on the bound here and postpone an example with GMRES convergence graphs to Section 4. The advection-diffusion problem is suitable as the eigenpairs can be calculated explicitly,

$$(3.16) \quad \lambda_k = 2\kappa - 2\kappa \cos(k\pi h) + i \cdot 2h \sin(k\pi h) \quad \mathbf{v}_k = [\exp(ik\pi jh)]_{j \in \mathbb{Z}} \quad \text{for any } k,$$

and hence we see that  $\theta_k$  densely populate the ellipse  $\Psi$  centered at  $2\kappa\tau/h^2$  with semi-axis parallel to the real and imaginary axis and with width  $4\kappa\tau/h^2$  and height  $2\tau/h$ . First, we note that both Corollary 3.4 and 3.5 still hold. Importantly, we can sample  $\vartheta_k$  from  $\Psi$  and proceed in completely analogous manner, only now having  $X_{\vartheta_k} \in \mathbb{C}^{s \times s}$ . This seems to suggest that the symmetry of the branches of  $\Gamma_q$  wrt to the real axis is lost. However, as long as we sample  $\vartheta_k$  *symmetrically* wrt to the real axis the branch symmetry is preserved as we show next.

PROPOSITION 3.7. Let  $\vartheta \in \mathbb{C}$  have positive imaginary part. Taking  $M_\vartheta, P_\vartheta^\star$  and  $X_\vartheta^\star$  as in Proposition 3.3 for any  $\star \in \{d, \text{GSU}, u, \text{GSL}, l\}$  we get

$$X_\vartheta^\star \mathbf{v}_\vartheta = \xi_\vartheta \mathbf{v}_\vartheta \implies X_{\bar{\vartheta}}^\star \bar{\mathbf{v}}_\vartheta = \bar{\xi}_\vartheta \bar{\mathbf{v}}_\vartheta,$$

where  $\bar{\cdot}$  stands for the entry-wise complex conjugation. In particular, the eigenvalues of  $X_\vartheta^\star$  are complex conjugate to those of  $X_{\bar{\vartheta}}^\star$ .

<sup>8</sup>We keep the notation consistent with Section 2 and hence  $L$  has the  $1/h^2$  scaling in front.



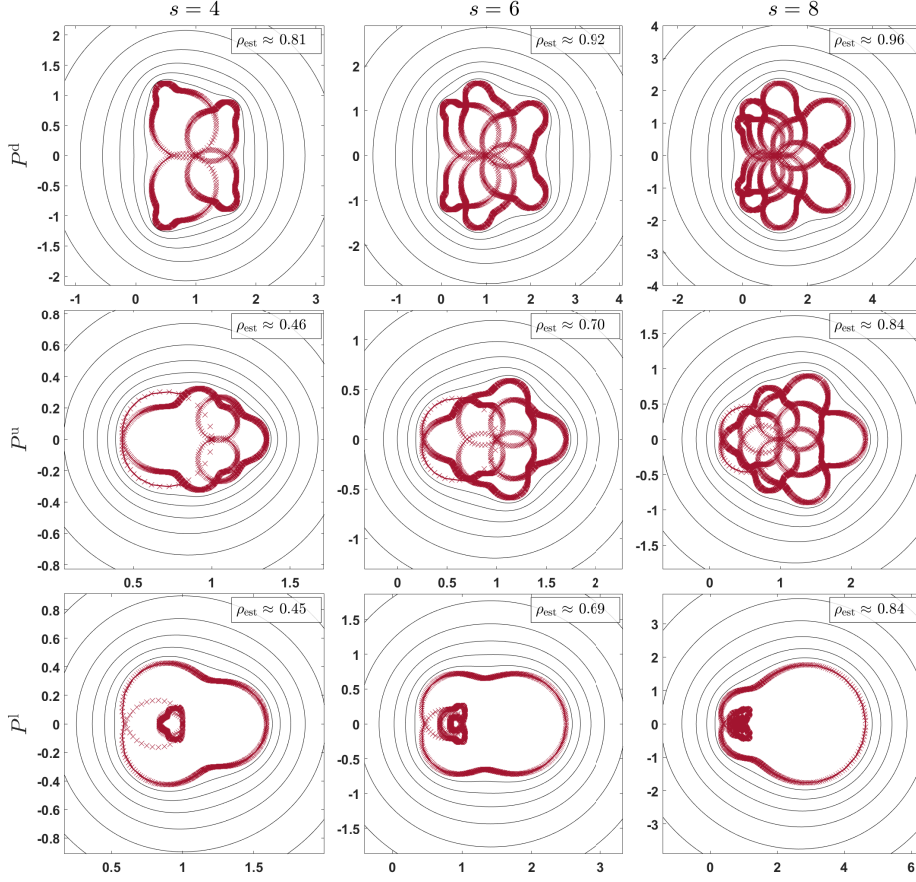


FIG. 7. In red: the eigenvalues of the preconditioned systems showing the symmetry predicted in Proposition 3.7 when taking  $\lambda_k$  as in (3.16), with  $\kappa = 0.01$  and  $h = 1/50$  and the Gauss Butcher tableau. In black: the contours of the corresponding Schwarz-Christoffel map of the exterior of these curves (or envelopes) mapped to the exterior of the unit circle, see `extermat` in [4]. We also show for each case the estimated linear convergence factor of GMRES  $\rho_{\text{est}}$ , see (3.14). To obtain these we use the techniques described above, i.e., calculating “envelopes” (which are not visible) based on suitable sparsification of the boundaries of the spectra.

*Proof.* The proof is identical for all choices of  $\star$  and we show it for  $\star = d$ . Throughout the proof we understand  $\bar{\cdot}$  as the entry-wise complex conjugation without any transposition of the vectors or matrices.

First, we notice that

$$M_{\bar{\vartheta}} = \overline{M_{\vartheta}} \quad \text{and} \quad P_{\bar{\vartheta}}^d = \overline{P_{\vartheta}^d}.$$

Next, we recall that  $(E + iF)^{-1} = (E + FE^{-1}F)^{-1} - iE^{-1}F(E + FE^{-1}F)^{-1}$  (for any  $E, F \in \mathbb{R}^{s \times s}$  and  $E$  invertible) and hence

$$(P_{\bar{\vartheta}}^d)^{-1} = \overline{(P_{\vartheta}^d)^{-1}}.$$

Recalling that for any  $X \in \mathbb{C}^{s \times s}$  and  $\mathbf{v} \in \mathbb{C}^s$ , we have  $\overline{X\mathbf{v}} = \overline{X}\overline{\mathbf{v}}$  we take the matrix

481  $X_\vartheta^d$  with an eigenpair  $(\xi_\vartheta, \mathbf{v}_\vartheta)$  and calculate

$$482 \quad X_\vartheta^d \bar{\mathbf{v}}_\vartheta = \overline{M_\vartheta} \left( \overline{P_\vartheta^d} \right)^{-1} \bar{\mathbf{v}}_\vartheta = \overline{M_\vartheta} \left( \overline{P_\vartheta^d} \right)^{-1} \mathbf{v}_\vartheta = \overline{X_\vartheta^d \mathbf{v}_\vartheta} = \overline{\xi_\vartheta} \bar{\mathbf{v}}_\vartheta,$$

483 finishing the proof.  $\square$

484 In other words, as long as  $\Psi$  is symmetrical wrt to the real axis and we sample  
 485 pairs of complex conjugate points along it, the analysis and techniques described above  
 486 can be used without any need for adjustments. We show the plots corresponding to  
 487 the discretization of the model problem (3.15) in Figure 7. We comment on some  
 488 direct generalizations next.

489 *Remark 3.8.* Importantly, some relevant, higher-dimensional problems lead to  $L$   
 490 with spectrum along *unions of 1D curves*, i.e., along  $\Psi_1, \dots, \Psi_m$ , see, e.g. [17, Sec-  
 491 tion 6]. The above techniques can be applied to each  $\Psi_i$  separately and then taking  
 492 the appropriate mix of the resulting envelopes in order to obtain GMRES estimates.

493 If  $\Psi$  is not symmetrical, then the techniques need to be adjusted when using the  
 494 Schwarz-Christoffel toolbox, as  $\Gamma_q$  is possibly non-symmetric wrt the real axis but  
 495 otherwise the results still apply.

496 We also want to comment on a similarity with the results in [16, 17]. There, the  
 497 authors addressed the question of *delay of convergence* by using similar formulations  
 498 to ours, also obtaining a GMRES problem reformulated as for a block-diagonal ma-  
 499 trix using Kronecker-product-like techniques as in Lemma 3.1. In particular, in [17,  
 500 Section 3.1] the authors use the equality

$$501 \quad \|\mathbf{r}_\ell\| = \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \left\| \varphi \left( \begin{bmatrix} X_1 & & \\ & \ddots & \\ & & X_n \end{bmatrix} \right) \mathbf{r}_0 \right\| = \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \sqrt{\sum_{j=1}^n \left\| \varphi(X_j) \mathbf{s}_0^{(j)} \right\|^2},$$

502 where  $\mathbf{s}_0^{(i)}$  is the  $i$ -th subvector of length  $s$  of  $Q^T \Pi \mathbf{r}_0$ , to obtain a lower bound

$$503 \quad (3.17) \quad \|\mathbf{r}_\ell\|^2 = \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \sum_{j=1}^n \left\| \varphi(X_j) \mathbf{s}_0^{(j)} \right\|^2 \geq \sum_{j=1}^n \min_{\substack{\varphi(0)=1 \\ \deg(\varphi) \leq \ell}} \left\| \varphi(X_j) \mathbf{s}_0^{(j)} \right\|^2$$

504 on the GMRES convergence behavior, explaining the initial stagnation phase in an  
 505 advection-diffusion problem. This way they bound the *global* minimization problem  
 506 (corresponding to solving a problem with the block-diagonal matrix  $\text{diag}(X_1, \dots, X_n)$ )  
 507 by the sum of the *local* minimization problems (each given by the small  $s$ -by- $s$  matrix  
 508  $X_j$ ). By careful analysis of the interplay of the right-hand side (or initial residual)  
 509 and the diagonal blocks in [17, Section 3.1] (there the diagonal blocks are, moreover,  
 510 tridiagonal and Toeplitz), the authors conclude

511 “... the presence of at least one system with tridiagonal Toeplitz ma-  
 trix  $T_j = \text{tridiag}(\gamma_j, \lambda_j, \mu_j)$  that is ‘close to the Jordan block’ (cf. [17,  
 Section 3.3] but see also [16]), and with  $l$  representing the index of  
 the first significant entry of the corresponding right-hand side, pre-  
 vents fast convergence of GMRES for the first  $N - l$  steps ( $N$  being  
 512 the size of the blocks  $T_j$ ) ...

... As explained in Section 3.1, the lower bound is useless for an-  
 analyzing the convergence behavior after the step  $N - l$ , possibly even  
 earlier. Hence the above approach cannot be used for quantifying any  
 possible acceleration of convergence after the initial phase. ”

We see that the approach is *fundamentally* different – both in the intended direction as well as in the results it can deliver – in spite of the fact that it works with the same technique.

We finalize this section with a remark on the *field of values* (sometimes also called the numerical range) and *pseudospectra*, which sometimes are *extremely* useful to understand and predict GMRES convergence behavior, especially if the eigenbasis of the system matrix is ill-conditioned, see, e.g., [7] and also [18, Section 5.7.3, pp. 296] and the references therein.

*Remark 3.9.* Another commonly used bound for GMRES uses the *field of values*  $\nu(C)$  or the  $\delta$ -*pseudospectrum*  $\sigma_\delta(C)$  of the system matrix  $C$ . By a direct calculation we obtain, for our model problem, the field of values as

$$\nu(MP^{-1}) = \sum_{i=1}^n \nu(X_k) \quad (\text{and analogously for } \nu(P^{-1}M)),$$

where the  $X_k$  are given as in (3.7) and the set addition is understood element-wise, i.e.,  $\nu(X_1) + \nu(X_2) = \{\alpha_1 + \alpha_2 \mid \alpha_1 \in \nu(X_1), \alpha_2 \in \nu(X_2)\}$ , or, more generally

$$\nu(MP^{-1}) \subset \kappa(Q) \sum_{i=1}^n \nu(X_k) \quad (\text{and analogously for } \nu(P^{-1}M)).$$

For the pseudospectrum we obtain an analogous formula, namely

$$\sigma_\delta(MP^{-1}) \subset \kappa(Q) \sum_{i=1}^n \sigma_\delta(X_k) \quad (\text{and analogously for } \nu(P^{-1}M)).$$

In other words, the principle of working with the small matrices  $X_k$  instead of the large matrix  $MP^{-1}$  naturally applies also to the other standard techniques for analyzing GMRES convergence behavior. However, adapting and using bounds based on field of values or the pseudospectrum of the preconditioned system for this set-up remains a topic for future research.

**4. Numerical Examples.** In this section we use the above analysis for more involved settings and also to demonstrate the convergence estimates (instead of only the convergence factor estimates). To be precise, we consider the *convergence estimates*

$$(4.1) \quad \frac{\|\mathbf{r}_\ell\|}{\|\mathbf{r}_0\|} \lesssim \min \{ \kappa_S^{\text{est}} \rho_{\text{est}}^\ell, 1 \},$$

where the estimate  $\kappa_S^{\text{est}}$  of  $\kappa_S$  is computed from the eigenbasis condition numbers of the “fake sampled” matrices  $X_{\vartheta_k}$  for  $k = 1, \dots, q$ . The convergence factor estimates reflect only the *spectral* part of the bound (2.11). Including an estimate of the term  $\kappa(S)$  in (2.11) then gives us a *convergence estimates*, which we show in this section.

We recall that the seeming independence of the preconditioner quality on the spatial mesh size  $h$  was sufficiently documented elsewhere (see [23, 20, 3, 10, 21]) and explained in Section 3 so that in our eyes, there is no need to address this direction here. Illustration of the solutions as well as further numerical experiments can be found in [21, Chapter 7]. For the sake of simplicity, we fix the number of time steps to balance the spatial and time error (see the (L2) definition in Section 2), namely we

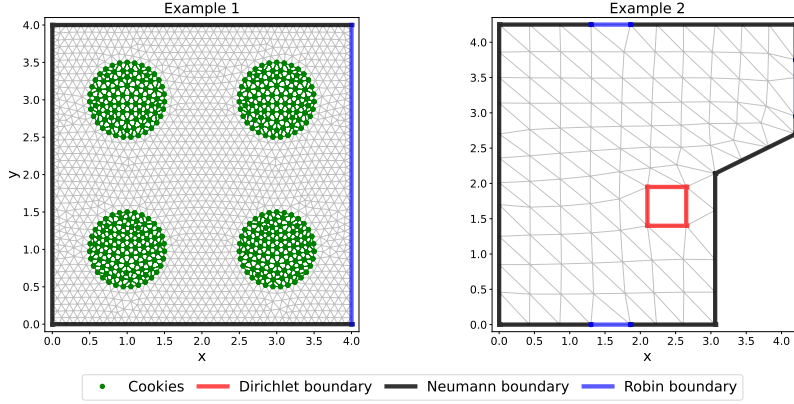


FIG. 8. The initial triangulations for Example 1 and 2 together with the boundary condition types and, for Example 1, also with highlighting the points with lower heat conductivity.

consider second order space discretization schemes,  $p$ -th order Runge-Kutta schemes and fix

$$\tau = h^{\frac{2}{p}}.$$

Last but not least, we have not set a relative residual tolerance criterion for stopping GMRES, meaning that GMRES went on until either the relative residual was on the level of machine precision or the maximum number of iterations was reached. This is not a good choice from the point of view of the solution process efficiency but since our primary focus is on studying the preconditioners, we found this reasonable.

*Diffusion problems.* We consider FEM discretizations in space<sup>9</sup> for discontinuous diffusion coefficient and for perforated domain in Example 1 and 2 with varying boundary conditions, see Figure 8.

*Example 1: Cookies in the oven.* The first problem is a simulation of baking cookies in an electrical oven projected in 2D, an idea borrowed from [15]. The cookies have a worse heat conductivity than the surrounding air (piecewise constant in space and constant in time) and the setting demands various boundary conditions, resulting in

$$\begin{aligned} \frac{\partial u}{\partial t} u &= \operatorname{div}(\sigma \nabla u) + f \quad \text{in } \Omega \times (0, T], \\ \frac{\partial u}{\partial \mathbf{n}} u &= 0 \quad \text{on } \Gamma_N \times (0, T], \quad \frac{\partial u}{\partial \mathbf{n}} u + pu = 0 \quad \text{on } \Gamma_R \times (0, T], \\ u &= 0 \quad \text{at } \Omega \times \{0\}, \end{aligned}$$

<sup>9</sup>Wherever we talk about a FEM discretization, we use linear Lagrange polynomials on conforming triangular meshes. Those are refined by the standard quadrissection of a triangle, with additional post-smoothing of the mesh.

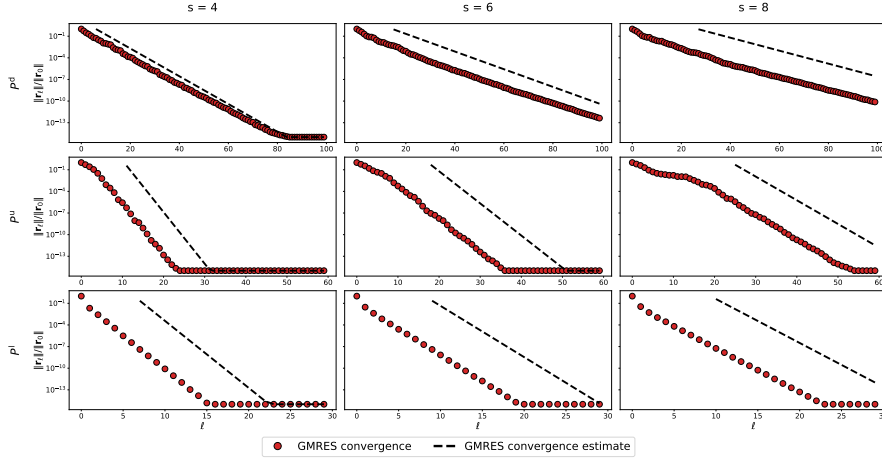


FIG. 9. The GMRES convergence behavior with the convergence estimates based on  $\rho_{\text{est}}$  for Example 1 with  $n = 26985$ .

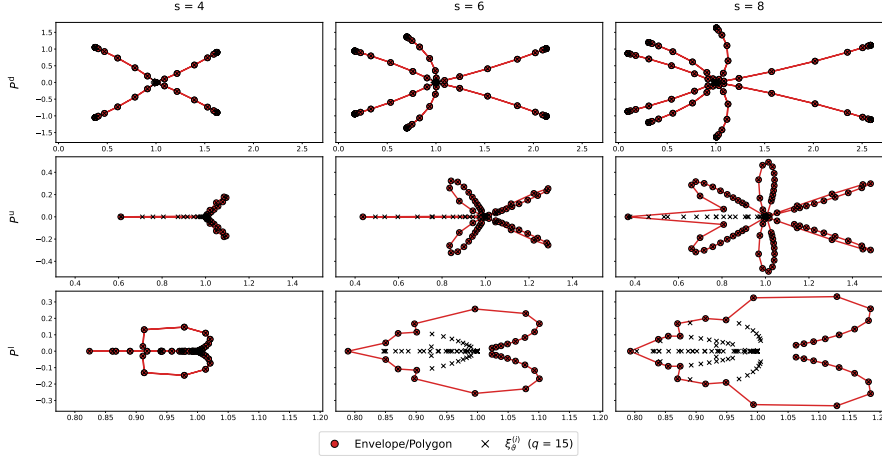


FIG. 10. The (sparsified) polygon approximations of the algebraic curves that are used in the Schwarz-Christoffel MATLAB toolbox to calculate  $\rho_{\text{est}}$  for Example 1 – for some settings these approximations correspond to the eigenvalues  $\xi_{\theta}^{(i)}$  and in some these approximations only enclose  $\xi_{\theta}^{(i)}$ .

with  $\Omega = (0, 4) \times (0, 4)$  and the boundary of  $\Omega$  is split into the Neumann and Robin parts  $\Gamma_N, \Gamma_R$ . We set the data as

$$\Gamma_N = \{x = 0\} \cup \{y = 0\} \cup \{y = 4\}, \quad \Gamma_R = \{x = 4\}, \quad p = 1, \sigma = \begin{cases} 10^3 & \text{if } (x, y) \in \text{Cookie}, \\ 1 & \text{otherwise,} \end{cases}$$

$$f(x, y, t) = \begin{cases} 3 & \text{if } \|(x, y) - (2, 2)\| \leq 1, \\ 0 & \text{otherwise,} \end{cases}$$

and show the GMRES convergence behavior with the estimates in Figure 9 as well as the sampling of the algebraic curves in Figure 10.

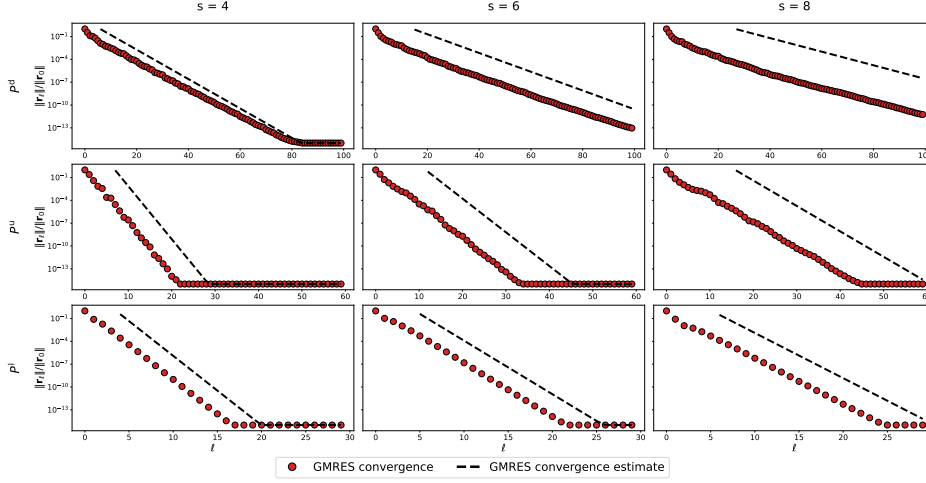


FIG. 11. The GMRES convergence behavior with the convergence estimates based on  $\rho_{\text{est}}$  for Example 2 with  $n = 26985$ .

*Example 2: The cabin heating.* The second problem uses the 2D projection of an attic room of a cabin in the western Bohemia region, whose primary heating is the chimney (bottom-right corner, modeled with a Dirichlet boundary condition changing in time), with two windows (top and bottom) and a door (right), modeled with Robin boundary conditions with Robin parameters  $p_w$  and  $p_d$ , and a good insulation otherwise, modeled with a Neumann condition. We obtain the problem

$$\begin{aligned} \frac{\partial u}{\partial t} u &= \operatorname{div}(\sigma \nabla u) \quad \text{in } \Omega \times (0, T], \\ \frac{\partial u}{\partial \mathbf{n}} u &= 0 \quad \text{on } \Gamma_N \times (0, T], \quad \frac{\partial u}{\partial \mathbf{n}} u + p u = 0 \quad \text{on } \Gamma_R \times (0, T], \\ u &= 0 \quad \text{at } \Omega \times \{0\}, \end{aligned}$$

and take the data as

$$\sigma = 1, \quad p_w = 0.1, \quad p_d = 10, \quad g_D(x, y, t) = \begin{cases} \min\{t, 0.7\} & \text{if } (x, y) \in \Gamma_D, \\ 0 & \text{otherwise,} \end{cases}$$

and show the GMRES convergence behavior with the estimates in Figure 11 as well as the sampling of the algebraic curves in Figure 12.

*Summary.* The convergence factor estimates are virtually as accurate as for the model problems in Section 3 – in Figures 9 and 11 this is clearly visible by comparing the slopes of the red and black “lines”, similarly to Figure 5. But the conditioning of the matrices  $X_{\theta_k}$  notably deteriorated as we increased  $s$ , hence worsening a bit the convergence estimates. The fact that this does not show up in the GMRES convergence behavior suggests that more delicate bounds, such as mentioned in Remark 3.9 could give a more detailed insight into the matter. However, in all cases the convergence estimates lag behind the actual convergence behavior by 10-20 iterations (which is in many if not most situations considered to be reasonably accurate).

We also showed the polygons used in the Schwarz-Christoffel toolbox. In our experience, large values of  $q$  lead to crowding problems in the SC toolbox but luckily

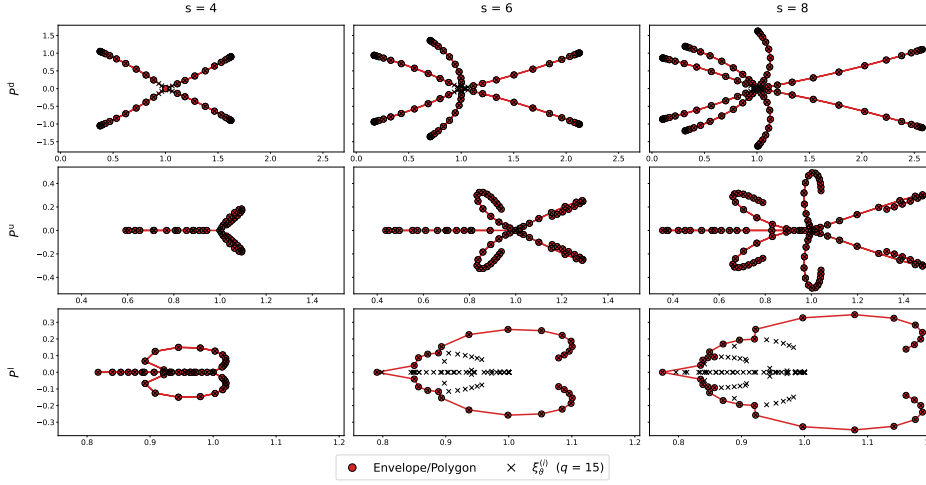


FIG. 12. The (sparsified) polygon approximations of the algebraic curves that are used in the Schwarz-Christoffel MATLAB toolbox to calculate  $\rho_{\text{est}}$  for Example 2 – for some settings these approximations correspond to the eigenvalues  $\xi_g^{(i)}$  and in some these approximations only enclose  $\xi_g^{(i)}$ .

even a very small value was usually enough. We also found that spacing the fake points  $\vartheta_k$  *logarithmically* in the corresponding interval somewhat alleviates this issue and leads to more accurate predictions of the arcs of the given algebraic curve. Nevertheless, notice that in many of the plots we excluded part of the arcs, mainly because either (a) the arcs intersected and we took the envelope of the algebraic curve (usually for the preconditioner  $P^1$ ) or (b) the points sampled along the arcs crowded sections of the arcs, which caused issues for the toolbox. In such cases we sparsified these regions by dropping some of these points. As a result, the Schwarz-Christoffel external map converged better and faster than for the problem in Section 3.1 and the contours were “ripple-free” for all of our problems, otherwise looking almost precisely as the ones in Figure 4.

*Advection problem.* We consider a centered FD discretization in space of a 2D advection problem on a unit square, i.e.,

$$(4.2) \quad \begin{aligned} \frac{\partial u}{\partial t} &= \mathbf{a} \cdot \nabla u + f \quad \text{in } \Omega \times (0, T], \\ u &= 0 \quad \text{on } \partial\Omega \times (0, T], \quad u = 0 \quad \text{at } \Omega \times \{0\}, \end{aligned}$$

with  $\Omega = (0, 1) \times (0, 1)$  and

$$\mathbf{a} = [1, 1]^T \quad \text{and} \quad f(x, y, t) = \begin{cases} 10 & \text{if } \|(x, y) - (0.5, 0.5)\| \leq 0.2, \\ 0 & \text{otherwise,} \end{cases}$$

and show the GMRES convergence behavior with the estimates in Figure 13 as well as the sampling of the algebraic curves in Figure 14.

We used a larger value of  $q = 300$  in order to capture the branches of  $\Gamma_q$  (which is not possible with  $q = 15$  but can plausibly be done with lower values than 300), and used sparsification of the envelopes to ensure smooth convergence of the SC toolbox. We see that the convergence rate estimates are again very accurate in most cases.

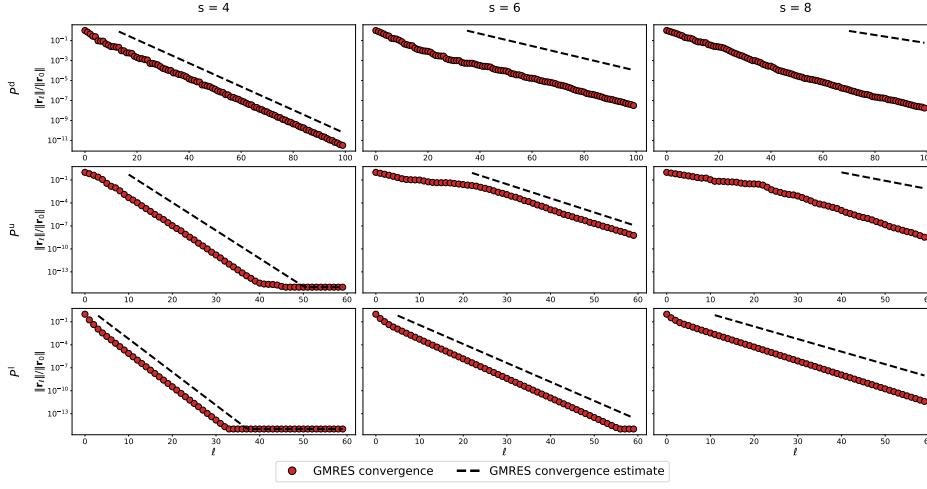


FIG. 13. The GMRES convergence behavior with the convergence estimates based on  $\rho_{\text{est}}$  for the advection problem (4.2) with  $n = 22210$ .

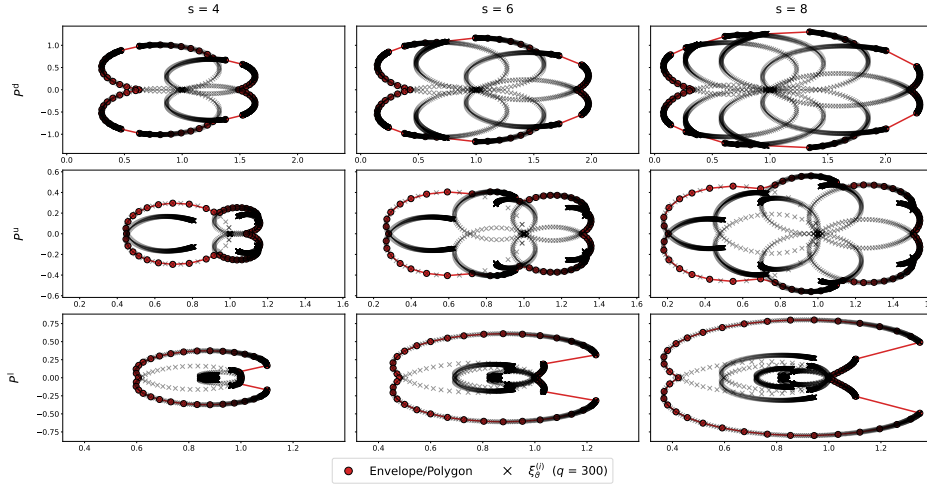


FIG. 14. The (sparsified) polygon approximations of the algebraic curves that are used in the Schwarz-Christoffel MATLAB toolbox to calculate  $\rho_{\text{est}}$  for the advection problem (4.2).

Notably, for  $s = 6$  and  $P^u$  the GMRES convergence estimate is more accurate than for the diffusive problems because GMRES convergence suffered from the non-normality of the system eigenbasis, and hence including the condition number estimate in (4.1) reflected an actual GMRES behavior. Unfortunately, for  $s = 8$  and  $P^{d,u}$  the term  $\kappa_S^{\text{est}}$  seems to fully dominate the bound.

**5. Concluding remarks.** Our main goal has been to understand the block preconditioners considered in [23, 3, 20] in more detail, and to try to explain their success and/or limitations. This goal was, in our eyes, mostly achieved but could be further improved in the sense of Remark 3.9 or by considering a more refined version of the bound (2.11), see [7, Section 2.1, equations (2.1) and (EV')] – this remains an area of interest for us for the future. Moreover, the above analysis can be directly used to



try to *optimize* Runge-Kutta methods, following the ideas in [23, 21, 10]. We also note that in practice, solving with either of the matrices  $P^{d,u,l,GSU,GS\bar{L},\dots}$  is often done with some level of *inaccuracy*, e.g., using a multigrid method. The question of interaction of this inaccuracy with the overall GMRES convergence is an important one and to the best of our knowledge has been addressed only numerically in [21, Chapter 7]. We also note that adapting the above analysis to the framework presented in [27, 26], or reformulating it from the vector equation to the matrix equation as suggested in [22], and to study in detail the comparison of these approaches for the IRK setting are attractive directions for future research.

**Acknowledgements.** Some of the ideas were stimulated by conversations with Mark Embree, Patrick Farrell, Miroslav Tůma and Petr Tichý and we would like to thank them for their inspiring comments and suggestions. We would also like to thank the anonymous reviewers for the careful reading of our manuscript and the suggestions for its improvements.

## REFERENCES

- [1] M. ARIOLI, V. PTÁK, AND Z. STRAKOŠ, *Krylov sequences of maximal length and convergence of GMRES*, BIT, 38 (1998), pp. 636–643.
- [2] P. BADDOO AND L. N. TREFETHEN, *Log-lightning computation of capacity and Green’s function*, Maple Transactions, 1 (2021).
- [3] M. R. CLINES, V. E. HOWLE, AND K. R. LONG, *Efficient order-optimal preconditioners for implicit Runge-Kutta and Runge-Kutta-Nyström methods applicable to a large class of parabolic and hyperbolic PDEs*. arXiv: <https://arxiv.org/abs/2206.08991>, 2022, <https://doi.org/10.48550/ARXIV.2206.08991>.
- [4] T. A. DRISCOLL, *A MATLAB toolbox for Schwarz-Christoffel mapping*, Tech. Report 2, 1996.
- [5] T. A. DRISCOLL, K.-C. TOH, AND L. N. TREFETHEN, *From potential theory to matrix iterations in six steps*, SIAM Rev., 40 (1998), pp. 547–578.
- [6] T. A. DRISCOLL AND L. N. TREFETHEN, *Schwarz-Christoffel mapping*, Cambridge University Press, Cambridge, First ed., 2002.
- [7] M. EMBREE, *How descriptive are GMRES convergence bounds?*, 2023, <https://arxiv.org/pdf/2209.01231.pdf>. arXiv preprint: 2209.01231.
- [8] M. EMBREE AND L. N. TREFETHEN, *Green’s functions for multiply connected domains via conformal mapping*, SIAM Rev., 41 (1999), pp. 745–761.
- [9] V. FABER, J. LIESEN, AND P. TICHÝ, *On Chebyshev polynomials of matrices*, SIAM J. on Matrix Anal. Appl., 31 (2010), pp. 2205–2221.
- [10] M. J. GANDER AND M. OUTRATA, *Spectral analysis of implicit 2-stage block Runge-Kutta preconditioners*, Linear Algebra Appl., (2023), <https://doi.org/10.1016/j.laa.2023.07.008>.
- [11] A. GREENBAUM, V. PTÁK, AND Z. STRAKOŠ, *Any nonincreasing convergence curve is possible for GMRES*, SIAM J. Matrix Anal. Appl., 17 (1996), pp. 465–469.
- [12] A. GREENBAUM, Z. STRAKOŠ, M. J. GANDER, AND M. OUTRATA, *Matrices that generate the same Krylov residual spaces*, in Recent Advances in Iterative Methods, G. H. Golub, A. Greenbaum, and M. Luskin, eds., vol. 60 of IMA Volumes in Mathematics and its Applications, Springer, 1994, pp. 95–118.
- [13] E. HILLE, *Analytic Function Theory*, vol. 2, American Mathematical Society, Chelsea Publishing, Providence, 2002.
- [14] T. KATO, *Perturbation Theory for Linear Operators*, vol. 132, Springer Berlin, Heidelberg, 2013.
- [15] D. KRESSNER AND C. TOBLER, *Low-rank tensor Krylov subspace methods for parametrized linear systems*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1288–1316.
- [16] J. LIESEN AND Z. STRAKOŠ, *Convergence of GMRES for tridiagonal Toeplitz matrices*, SIAM J. on Matrix Anal. Appl., 26 (2004), pp. 233–251.
- [17] J. LIESEN AND Z. STRAKOŠ, *GMRES convergence analysis for a convection-diffusion model problem*, SIAM J. Sci. Comput., 26 (2005), pp. 1989–2009.
- [18] J. LIESEN AND Z. STRAKOŠ, *Krylov Subspace Methods: Principles and Analysis*, Oxford University Press, Oxford, 2013.
- [19] P. MUNCH, I. DRAVINS, M. KRONBICHLER, AND M. NEYTCHIEVA, *Stage-parallel fully implicit*

- Runge–Kutta implementations with optimal multilevel preconditioners at the scaling limit, SIAM J. Sci. Comput., (2023), pp. S71–S96.
- [20] M. NEYTCHIEVA AND O. AXELSSON, *Numerical Solution Methods for Implicit Runge–Kutta Methods of Arbitrarily High Order*, in Proceedings of the Conference Algoritmy 2020, P. Frolkovič, K. Mikula, and D. Ševčovič, eds., Slovak University of Technology in Bratislava, Vydavateľstvo SPEKTRUM, 2020.
- [21] M. OUTRATA, *Schwarz methods, Schur complements, preconditioning and numerical linear algebra*, PhD thesis, University of Geneva, Math Department, 2022.
- [22] D. PALITTA AND V. SIMONCINI, *Optimality properties of Galerkin and Petrov–Galerkin methods for linear matrix equations*, Vietnam J. Math., 48 (2020), pp. 791–807.
- [23] M. M. RANA, V. E. HOWLE, K. LONG, A. MEEK, AND W. MILESTONE, *A New Block Preconditioner for Implicit Runge–Kutta Methods for Parabolic PDE Problems*, SIAM J. Sci. Comput., 43 (2021), pp. S475–S495.
- [24] T. RANSFORD, *Potential Theory in the Complex Plane*, no. 28, Cambridge University Press, Cambridge, 1995.
- [25] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, Other Titles in Applied Mathematics, SIAM, Philadelphia, Second ed., 2003.
- [26] B. S. SOUTHWORTH, O. KRZYSIK, AND W. PAZNER, *Fast solution of fully implicit Runge–Kutta and discontinuous Galerkin in time for numerical PDEs, Part II: nonlinearities and DAEs*, SIAM J. Sci. Comput., 44 (2022), pp. 636–663.
- [27] B. S. SOUTHWORTH, O. KRZYSIK, W. PAZNER, AND H. DE STERCK, *Fast solution of fully implicit Runge–Kutta and discontinuous Galerkin in time for numerical PDEs, Part I: The linear setting*, SIAM J. Sci. Comput., 44 (2022), pp. 416–443.
- [28] G. A. STAFF, K.-A. MARDAL, AND T. K. NILSSEN, *Preconditioning of fully implicit Runge–Kutta schemes for parabolic PDEs*, Modeling, Identification and Control, 27 (2006), pp. 109–123.
- [29] C. F. VAN LOAN, *The ubiquitous Kronecker product*, J. Comput. Appl. Math., 123 (2000), pp. 85–100.
- [30] G. WANNER AND E. HAIRER, *Solving Ordinary Differential Equations II : Stiff and Differential-Algebraic Problems*, Springer Berlin, Heidelberg, 1996.
- [31] G. WANNER, S. P. NØRSETT, AND E. HAIRER, *Solving Ordinary Differential Equations I : Non-Stiff Problems*, Springer Berlin, Heidelberg, 1987.