

# Preconditioning parametrized linear systems: hierarchical maps approach

Eric de Sturler and Michal Outrata

## Abstract

We propose, analyze and numerically experiment with a new type of preconditioner maps for sequences of systems of linear algebraic equations. Similarly to [5], we propose to renovate a good-quality preconditioner for the first problem  $A_1 \mathbf{x}_1 = \mathbf{b}_1$  so that we obtain a comparably good preconditioner also for the follow-up problems  $A_k \mathbf{x}_k = \mathbf{b}_k$  in the problem sequence. Complementary to [5] we use *data-sparsity* techniques, i.e., hierarchical matrix formats and propose two different preconditioner maps using these. We demonstrate the quality as well as efficiency of the proposed preconditioner maps and try to further improve on these numerically on a model problem.

## 1 Introduction

Consider a sequence of linear problems

$$A_k \mathbf{x}_k = \mathbf{b}_k, \quad \text{for } k = 1, 2, \dots, \quad (1)$$

where  $A_k \in \mathbb{R}^{N \times N}$ ,  $\mathbf{b}_k \in \mathbb{R}^N$  and denote the update matrices by  $E_{k,k+\ell}$ , i.e.,

$$A_k = A_k + E_{k-\ell,k}, \quad \text{for } \ell = 1, \dots, k-1.$$

In many applications the *efficient* solution of each of these problems requires the use of iterative method such as Krylov methods coupled with *preconditioners* (see ). Obtaining these can be computationally very demanding but overall still improves the running time significantly (see, e.g., ). Assuming we have computed a high-quality preconditioner  $P_1$  for the first system so that the GMRES method applied to

$$A_1 P_1 \mathbf{y}_1 = \mathbf{b}_1, \quad (2)$$

converges rapidly, we note that even for relatively small consecutive updates  $E_{k-1,k}$ , the GMRES convergence with fixed preconditioner  $P_1$  can deteriorate drastically, see an illustrative example in Figure 2 in Section 4. This means that the preconditioner  $P_1$  needs an update as well and in [5, Section 2], the authors proposed<sup>1</sup> a preconditioning map  $\mathcal{P}$ ,

$$\mathcal{P} : (A_1, A_k) \mapsto P_k \approx A_k^{-1} A_1, \quad (3)$$

so that

$$A_k P_k P_1 \approx A_1 P_1.$$

---

<sup>1</sup>A similar idea but for a more particular setting was already proposed in [1].

In words, we shift the focus from construction of a stand-alone preconditioner for the  $k$ -th system to construction of  $P_k$  such that  $A_k P_k$  system resembles the original one for which we already have a high-quality preconditioner. Alternatively, we can also think of this map as *renovating* the preconditioner  $P_1$  into the preconditioner  $P_k P_1$  with an important distinction that the product  $P_k P_1$  is never assembled but rather kept as two separate mat-vec routines in the preconditioned GMRES algorithm – this is preferable both because in some situations  $P_1$  is available itself only as a mat-vec routine as well as because of the computational burden of performing matrix-matrix product for large systems, which is often prohibitive. Moving forward we continue with the set-up of a right-preconditioner but the same ideas can be applied analogously to the case of a left-preconditioner. The case of split preconditioning is more interesting and will be treated separately elsewhere.

In [5], the authors study a particular type of  $\mathcal{P}$  in (3), using the sparse approximate inverse technique to obtain an efficient approximation and call the resulting  $\mathcal{P}$  the *sparse approximate map* (SAM). The key point then lies in balancing the trade-off between the *sparsity pattern* of  $P_k$ , which controls the computational complexity of its application, and the error introduced by the *approximation*, which controls the quality of the map. In [5, Section 3], the authors experimentally show that even for modest sparsity patterns (there the authors considered the patterns of powers of the system matrix) and approximation error in (3), the SAM approach can be very effective. Naturally, the area of updating preconditioners for systems of problems (possibly with only right-hand sides changing or only system matrices changing) is much richer and we refer the interested reader to the literature cited in [5] but also to [\[please recommend another reference\]](#) and the references therein. Our approach here is complementary to that in [5] – we want to use *data-sparsity* instead of the *structural sparsity* to construct a suitable preconditioner map, starting with hierarchical matrix formats such as HODLR (a non-nested format) and HSS (a nested format) but with emphasis on the generality, so that other hierarchical formats can be easily put in the place of HODLR/HSS. As a result we use a different map than the one in (3). We also analyze the GMRES convergence behavior for the preconditioned systems using our preconditioner map, doing so in a complementary way to [5], i.e., using the pseudospectra and field of values bounds, in contrast to the spectral bound in [5]. The analysis is not directly transferable but the approach is, i.e., our analysis approach can be applied to obtain analogous bounds also in [5] and, reversely, the spectral bound in [5] could be adapted to our preconditioner map as well.

The rest of the paper is structured as follows: we introduce the hierarchical approximate maps (HAMs) as well as other necessary terms in Section 2 and analyze the GMRES convergence behavior of the resulting preconditioners in Section 3, using the pseudospectral bounds. We explore this new approach numerically for a model problem in Section 4 – we start with a demonstration of the bounds in Section 4.1 and continue with the question of efficiency of the preconditioners in Section 4.2. There we explore several different directions for improving the efficiency and, finally, in summarize the contributions in Section 5.

## 2 HAM: hierarchical approximate maps

We start by looking at the *optimal* preconditioner map – obtained by assuming equality in (3) instead of approximation – for the second system and obtain  $P_2^{\text{opt}}$  as

$$P_2^{\text{opt}} = A_2^{-1} A_1 = (A_1^{-1} A_2)^{-1} = (A_1^{-1} (A_1 + E_{1,2}))^{-1} = (I + A_1^{-1} E_{1,2})^{-1}.$$

Replacing the solve with the matrix  $A_1$  with the application of the preconditioner  $P_1$  we obtain an *approximate* preconditioner map, which requires a solve with the matrix

$$\tilde{R}_2 := I + P_1 E_{1,2}.$$

To make this solve efficient we approximate  $\tilde{R}_2$  in a *hierarchical matrix format* – here we will consider two examples of such formats – HODLR and HSS, see [12] and [24] for an introduction to these – but we avoid, on purpose, using any particular properties of these so that any hierarchical format can be easily substituted for these. We consider the hierarchical formats in general in order to take advantage of the fast hierarchical solvers and thus obtain an approximation to solving with  $\tilde{R}_2$  – we denote these two operations by  $\mathcal{H}()$  for the approximation in a hierarchical format and  $\mathcal{H}_{solve}()$  for the hierarchical solve and get

$$P_2 = \mathcal{H}_{solve}(R_2), \quad \text{with} \quad R_2 := \mathcal{H}(I + P_1 E_{1,2}). \quad (4)$$

The preconditioner for  $A_2$  then becomes  $P_2 P_1$  – a product which is not assembled but rather applied piece by piece. The area of hierarchical matrices (and tensors) has yielded many very efficient approaches and algorithms in multiple different applications (see, e.g., [4, 11, 9, 3, 13, 2, 17] and also [10, 15, 25] among others) and has been a very active field of research for the last two decades. The efficiency is achieved via a multilevel scheme in which we approximate a certain off-diagonal part of the matrix in question on each level. Before generalizing to the  $k$ -th system, we recall that for hierarchical matrices we have many choices to make – the format (HODLR, HSS or some more involved formats, e.g.,  $\mathcal{H}$  and  $\mathcal{H}^2$  matrices, see [13, Sections 6–8]), the structure (usually given by the so-called cluster trees – row and column – and admissibility condition), the accuracy (denoted by  $\epsilon$ ) and the minimal block size (denoted by  $\beta$ ) to name the most common ones. Making a choice for all of these, we then obtain a hierarchical matrix of a particular *type*. As the set of all matrices of a particular type does not form a vector space with the standard algebraic operations, it is common to use the so-called *formatted* (or *hierarchical*) version of the algebraic operations, such as addition, multiplication or inversion. In the simplest form these correspond to performing the desired operation and then calculating the best hierarchical approximation of the given type of the outcome of that operation, i.e., projecting back onto the set of all matrices of the given type. Unfortunately, this would destroy the efficiency for many operations and so more sophisticated algorithms have been proposed to calculate these formatted/hierarchical operations; see [13, Section 3 and onwards]. We denote these formatted operation using the  $\circ$ -notation, e.g.,  $\oplus$  or  $\otimes$ , but shall not go into any more details, the implementation we will use is described in [18]; the effects of some of the above parameters are illustrated later in Section 4.

For the  $k$ -th system from (1), the preconditioner becomes  $P_k P_1$  with  $P_k = \mathcal{H}_{solve}(R_k)$  and we see two possible ways of defining  $R_k$ , namely

$$R_k^{(1)} := \mathcal{H}(I + P_1 E_{1,k}) \quad \text{and} \quad R_k^{(2)} := R_{k-1} \oplus \mathcal{H}(P_1 E_{k-1,k}). \quad (5)$$

Once the product of the update matrix with  $P_1$  is evaluated the rest of the operations is again done using the formatted arithmetic. Formally, we introduce two *hierarchical approximate maps* (HAMs)

$$\begin{aligned} \mathcal{P}^{(1)} : (P_1, E_{1,k}) &\mapsto P_k^{(1)} := \mathcal{H}_{solve}\left(R_k^{(1)}\right), \\ \mathcal{P}^{(2)} : (P_1, R_{k-1}, E_{k-1,k}) &\mapsto P_k^{(2)} := \mathcal{H}_{solve}\left(R_k^{(2)}\right). \end{aligned} \quad (6)$$

First, we note that  $\mathcal{P}^{(1)}$  is the natural analogue of the SAM approach in [5] but observe a key difference – evaluating the SAM map (i.e., obtaining the matrix  $R_k$ ; in [5] denoted by  $N_k$ ) does not involve  $P_1$  at all.

Second, as the construction of  $R_k^{(1,2)}$  requires<sup>2</sup> (among other steps) the matrix  $P_1 E$ , with  $E = E_{1,k}$  or  $E = E_{k-1,k}$ . If the update matrices  $E$  would each have *many* non-zero columns, then the evaluation of  $P_1 E$  requires many applications of  $P_1$  and can become the bottleneck of the entire solution process. In some applications we expect our system matrices  $A_k$  to form a “cluster”<sup>3</sup> around  $A_1$  so that both  $E_{1,k}$  and  $E_{k-1,k}$  contain roughly the same number of non-zero columns but for other applications the matrices  $A_k$  form a “string”<sup>4</sup> stretching away from  $A_1$ , where at each step the number of non-zero columns of  $E_{k-1,k}$  is manageable but already after few steps the matrix  $E_{1,k}$  accumulates large number of non-zero columns.

This highlights an important trade-off between the maps  $\mathcal{P}^{(1)}$  and  $\mathcal{P}^{(2)}$  – the map  $\mathcal{P}^{(1)}$  accumulates the information in the update matrix in contrast to the map  $\mathcal{P}^{(2)}$  that keeps some information from the previous steps in form of  $R_{k-1}$ . Moreover,  $\mathcal{P}^{(1)}$  allows us to *change* any parameters of the hierarchical type for *each* system  $A_k$ . This is not the case for  $\mathcal{P}^{(2)}$  as the formatted sum can be efficiently evaluate only if both matrices are of the same hierarchical type or there is an efficient conversion inbetween these. However the conversions are often far from efficient, depending on the hierarchical structures and formats involved.

Last but not least, we note that for  $\mathcal{P}^{(2)}$  it seems natural to keep  $P_k$  already assembled and update it *directly* in contrast to storing and updating  $R_k^{(2)}$  – let us denote this approach as  $\mathcal{P}^{(3)}$ . In context of approximate preconditioner maps this idea has been used in [1], where the update to  $R_k$  is rank one and  $P_k$  is updated directly using Sherman-Morrison-Woodbury formula. The main appeal is in reducing the computational complexities – for both HODLR and HSS formats (as well as for others) the cost of solving  $H\mathbf{x} = \mathbf{b}$  with  $H$  in a hierarchical format of rank<sup>5</sup>  $r$  scales as  $r^2$  while the cost of multiplying with the same  $H$  scales only as  $r$ . This approach is not possible with  $\mathcal{P}^{(1)}$  – there is no update to speak of – but can be derived based on  $\mathcal{P}^{(2)}$ , e.g., if we update the hierarchical inverse or the hierarchical LU factorization in some way. The fair comparison – computational-complexity-wise – is then comparing  $\mathcal{P}^{(3)}$  with using  $\mathcal{P}^{(2)}$  to obtain  $R_k^{(2)}$  and then calculating the hierarchical approximations to the *inverses of the LU factors* – the so-called *hierarchical factored approximate inverse* method/preconditioner ( $\mathcal{H}$ -FAINV), see [15]. The specifics of the hierarchical arithmetic makes it so that these operations have the same asymptotics (with respect to the system size  $N$  and the hierarchical rank  $r$ ) and the constants are comparable (see [15, Theorem 2 and Lemma 1 and below]). This means that the main appeal of the map  $\mathcal{P}^{(3)}$  – the lower complexity of application of the preconditioner – can be achieved also by combining  $\mathcal{P}^{(2)}$  with  $\mathcal{H}$ -FAINV and, moreover, the extra cost of  $\mathcal{H}$ -FAINV is comparable to one additional formatted matrix sum. All of

---

<sup>2</sup>We also note that many standard hierarchical formats can be built in a matrix-free fashion, i.e., without the need to assemble the matrix  $P_1 E$  and only by having mat-vec routine  $\mathbf{v} \mapsto P_1 E \mathbf{v}$  (e.g., the HODLR format, see [18, Section 3.3]). However, the construction sometimes requires also the action of the transpose, i.e.,  $\mathbf{v} \mapsto (P_1 E)^T \mathbf{v}$  (e.g., for the HSS format, see [18, Section 3.3]), which might not be always available, e.g., if  $P_1$  is itself only available as a mat-vec routine  $\mathbf{v} \mapsto P_1 \mathbf{v}$ . We shall not focus on this possibility here and will assume the necessity of the construction of the matrix. The matrix-free adaptation will be discussed elsewhere.

<sup>3</sup>In the sense that the number of nonzero columns in  $E_{1,k}$  does not grow substantially with  $k$ .

<sup>4</sup>In the sense that the number of nonzero columns in  $E_{1,k}$  does grow substantially faster than that of  $E_{k-1,k}$  with  $k$ .

<sup>5</sup>As noted above, the hierarchical formats keep the off-diagonal portions of the matrix in a low-rank format. Intuitively speaking, the highest rank of these is what we call the hierarchical rank of the matrix. Precise definition can, however, differ from format to format and can be found in the literature cited before.

that is, of course, under the assumption that both  $\mathcal{P}^{(3)}$  and  $\mathcal{P}^{(2)}$ -with- $\mathcal{H}$ -FAINV *achieve comparable accuracy with the same hierarchical rank*. We are not aware of any experimental results but to keep the present manuscript from growing too wide we will address this direction of updating  $P_k$  rather than  $R_k$  elsewhere. Next we turn our attention to providing some analysis for the proposed preconditioner maps.

### 3 Analysis of the HAM approach

The preconditioner maps have two points where an approximation was introduced – replacing  $A_1^{-1}$  with  $P_1$  and replacing the inverse/solve operation with its hierarchical approximation – and the analysis reflects that. The first is assumed to be under control thanks to having a very good preconditioner<sup>6</sup> while the other is under our control by a suitable choices of the hierarchical format, structure, accuracy and so on.

Looking at the  $k$ -th preconditioned system matrix  $A_k P_k^{(1)} P_1$  for the map  $\mathcal{P}^{(1)}$ , we set

$$\tilde{R}_k := I + P_1 E_{1,k} \quad (7)$$

and write

$$A_k P_k^{(1)} P_1 = A_k \mathcal{H}_{solve} \left( R_k^{(1)} \right) P_1 = B_k^{(1)} + C_k^{(1)} + \tilde{D}_k \quad (8)$$

with

$$\begin{aligned} B_k^{(1)} &:= A_k \left( \mathcal{H}_{solve} \left( R_k^{(1)} \right) - \left( R_k^{(1)} \right)^{-1} \right) P_1, \\ C_k^{(1)} &:= A_k \left( \left( R_k^{(1)} \right)^{-1} - \tilde{R}_k^{-1} \right) P_1, \\ \tilde{D}_k &:= A_k \tilde{R}_k^{-1} P_1 = A_k \left( P_1^{-1} + E_{1,k} \right)^{-1}, \end{aligned}$$

quantifying the three core approximation errors introduced in  $\mathcal{P}^{(1)}$ . The first two are controlled by the choice of the hierarchical format – the error introduced by the hierarchical LU factorization<sup>7</sup> in  $B_k^{(1)}$  (see, e.g., [11, Theorem 24]) and the error introduced by the hierarchical approximation of the matrix  $P_1 E_{1,k}$  in  $C_k^{(1)}$ . As most of the hierarchical formats keeps the diagonal entries explicitly, we in fact only approximate the product  $P_1 E_{1,k}$  and then add the identity along the diagonal, obtaining

$$C_k^{(1)} = A_k \left( R_k^{(1)} \right)^{-1} \left( \tilde{R}_k - R_k^{(1)} \right) \tilde{R}_k^{-1} P_1 = A_k \left( R_k^{(1)} \right)^{-1} \left( P_1 E_{1,k} - \mathcal{H}(P_1 E_{1,k}) \right) \tilde{R}_k^{-1} P_1. \quad (9)$$

For  $\mathcal{P}^{(2)}$  we follow the same idea and notation but adjust the approximation error due to the hierarchical matrix approximation to reflect the accumulation, i.e., we write

$$A_k P_k^{(2)} P_1 = A_k \mathcal{H}_{solve} \left( R_k^{(2)} \right) P_1 = B_k^{(2)} + C_k^{(2)} + \tilde{D}_k \quad (10)$$

---

<sup>6</sup>Similarly to above, we do not mean to say that good preconditioners always provide a good approximations of the inverse of the system matrix but rather that good preconditioners in some sense capture the essence (or a fundamental part of it) of the problem. However, the approximation quality can serve as an indicator and can be useful in absence of more refined tools.

<sup>7</sup>Notice that the error is introduced only there, the ensuing type and backward substitutions with the hierarchical factors does not introduce further approximation errors.

with

$$\begin{aligned} B_k^{(2)} &:= A_k \left( \mathcal{H}_{solve} \left( R_k^{(2)} \right) - \left( R_k^{(2)} \right)^{-1} \right) P_1, \\ C_k^{(2)} &:= A_k \left( \left( R_k^{(2)} \right)^{-1} - \tilde{R}_k^{-1} \right) P_1, \\ \tilde{D}_k &:= A_k \tilde{R}_k^{-1} P_1 = A_k \left( P_1^{-1} + E_{1,k} \right)^{-1}, \end{aligned}$$

and recalling that

$$\begin{aligned} R_k^{(2)} &= R_{k-1}^{(2)} \oplus \mathcal{H}(P_1 E_{k-1,k}) = R_{k-2}^{(2)} \oplus \mathcal{H}(P_1 E_{k-2,k-1}) \oplus \mathcal{H}(P_1 E_{k-1,k}) \\ &= I \oplus \mathcal{H}(P_1 E_{1,2}) \oplus \dots \oplus \mathcal{H}(P_1 E_{k-1,k}), \end{aligned}$$

we reformulate  $C_k^{(2)}$  analogously to (9) and obtain

$$C_k^{(2)} = A_k \left( R_k^{(2)} \right)^{-1} \left( \tilde{R}_k - R_k^{(2)} \right) \tilde{R}_k^{-1} P_1,$$

with

$$\begin{aligned} \tilde{R}_k - R_k^{(2)} &= \sum_{\ell=2}^k P_1 E_{\ell-1,\ell} - \bigoplus_{\ell=2}^k \mathcal{H}(P_1 E_{\ell-1,\ell}) \\ &= \sum_{\ell=2}^k (P_1 E_{\ell-1,\ell} - \mathcal{H}(P_1 E_{\ell-1,\ell})) + \sum_{\ell=2}^k \mathcal{H}(P_1 E_{\ell-1,\ell}) - \bigoplus_{\ell=2}^k \mathcal{H}(P_1 E_{\ell-1,\ell}). \end{aligned}$$

The matrix  $C_k^{(2)}$  captures the error accumulation due to the hierarchical approximation over the sequence of the problems. The second part, i.e., the difference of the algebraic and hierarchical sum of the hierarchical matrices, can be handled analogously to (9) but the error there is due to the *rank compression*<sup>8</sup> after the formatted sum. Altogether we see that  $\tilde{R}_k - R_k^{(2)}$  has  $2k-1$  terms that capture the hierarchical approximation error, compared to just one for  $\tilde{R}_k - R_k^{(1)}$ :  $k$  thanks to the  $k$  approximation errors of hierarchical format representation and  $k-1$  due to the recompression after the formatted additions. For long problem sequences this might make  $\mathcal{P}^{(2)}$  less appealing<sup>9</sup> or might require a HAM restart, where we compute a new preconditioner and start as if for a new system sequence.

For both  $\mathcal{P}^{(1)}$  and  $\mathcal{P}^{(2)}$  we see that the last term in (8) and (10) captures the approximation quality of  $P_1$  and is independent of the chosen map, namely we get

$$\begin{aligned} \tilde{D}_k &= A_k \left( P_1^{-1} + E_{1,k} \right)^{-1} = (A_1 + E_{1,k}) \left( P_1^{-1} + E_{1,k} \right)^{-1} = I + (A_1 - P_1^{-1}) \left( P_1^{-1} + E_{1,k} \right)^{-1} \\ &= I + (A_1 P_1 - I) (I + E_{1,k} P_1)^{-1}. \end{aligned} \tag{11}$$

Before addressing these further, we would like to emphasize a key difference between the HAM and SAM approaches in the following remark.

---

<sup>8</sup>For most hierarchical formats the efficiency stems from keeping a relatively small hierarchical rank  $r$ . However, after any formatted operation, the rank can quickly grow and hence, in practice, we usually follow the formatted operation with a hierarchical rank-compression – procedure where the hierarchical matrix is (approximately) projected onto the set of hierarchical matrices with hierarchical rank less or equal than some  $r_{\max}$ , see [18, Section 4.5] and also [13, Sections 7.2 and 8.8] for more details.

<sup>9</sup>As mentioned above, this can be sometimes balanced off by the accumulation of information in the update matrix  $E_{1,k}$  compared to the update matrices  $E_{\ell-1,\ell}$  for  $\ell = 2, \dots, k$ , which then results in a higher computational costs both for computing and applying  $P_k^{(1)}$  (assuming we keep the desired accuracy fixed).

**Remark 1.** *The SAM approach ([5]) and its analysis does not explicitly capture the interaction between the system change (here the update matrix  $E$ ) and the original preconditioned system – this interaction is “hidden” in the least square problem minimization and as a result the error terms can be fully user controlled quantities, see [5, equation (2.5) and below]. In our case, however, this interaction seems unavoidable and as a result the error contains also terms which are problem-dependent and cannot be adjusted for by the user.*

Following up on Remark 1, since we only know that  $A_k$  is non-singular for all  $k$ , it is a natural starting point for analysis to assume that the magnitude of the update is smaller than the distance of  $A_1$  to singularity, i.e., that  $\|E_{1,k}A^{-1}\| < 1$  so that indeed  $A_k$  is non-singular. Assuming that the analogue holds also for the initial preconditioner, i.e., that

$$\|E_{1,k}P_1\| =: \gamma < 1, \quad (12)$$

we can expand the inverse in  $\tilde{D}_k$  into its Neumann series, obtaining

$$\tilde{D}_k = I + (A_1P_1 - I)(I + E_{1,k}P_1)^{-1} = A_1P_1 + (A_1P_1 - I) \sum_{i=1}^{+\infty} (E_{1,k}P_1)^i =: A_1P_1 + D_k. \quad (13)$$

Alternatively, in many applications the update matrices  $E_{1,k}$  have a fixed rank  $\nu$ , i.e., we have

$$E_{1,k} = XY^T, \quad X, Y \in \mathbb{R}^{N \times \nu},$$

e.g., thanks to the localization of the update in the physical domain of the underlying PDE, see Section 4. Then, instead of using the Neumann series expansion, we use the Sherman-Morrison-Woodbury formula and, denoting  $Z := P_1^T Y$ , obtain

$$\tilde{D}_k = I + (A_1P_1 - I)(I + XZ^T)^{-1} = I + (A_1P_1 - I) \left( I - X(I + Z^T X)^{-1} Z^T \right) = A_1P_1 + D_k, \quad (14)$$

or, equivalently,

$$\tilde{D}_k = I + (A_1P_1 - I)(I + XZ^T)^{-1} = I + (A_1P_1 - I) \left( I - X(I + Z^T X)^{-1} Z^T \right) =: I + A_1P_1 - I + \hat{D}_k, \quad (15)$$

where  $\hat{D}_k$  is also of rank at most  $\nu$ . and we can write

$$A_k P_k^{(1,2)} P_1 = A_1 P_1 + B_k^{(1,2)} + C_k^{(1,2)} + D_k = I + B_k^{(1,2)} + C_k^{(1,2)} + A_1 P_1 - I + \hat{D}_k, \quad (16)$$

where both  $B_k^{(1,2)}$  and  $C_k^{(1,2)}$  are expected to be small due to the control of the hierarchical errors and we *hope* that the interaction of the update and the preconditioner captured by  $\tilde{D}_k$  can be efficiently treated with either “small in norm” approach or “low-rank” approach shown above – this is in general not guaranteed as we can have a moderate rank and  $\gamma > 1$  at the same time even for an excellent preconditioner  $P_1$ . However, if either  $\nu$  is small or  $\gamma < 1$ , then we get convergence bounds for the preconditioned GMRES for the system (1).

We continue with some direct bounds on the norm of the relevant matrices. For lack of further structure of the sequence  $A_k$ , i.e., assuming only (12), the bounds for  $B_k^{(1,2)}, C_k^{(1,2)}$  are quite crude

based on the sub-multiplicativity of the standard matrix norms, and we use the Neumann series bound to bound the norm of  $D_k$ , obtaining

$$\begin{aligned}\|B_k^{(1,2)}\| &\leq \|A_k\| \|P_1\| \left\| \mathcal{H}_{solve} \left( R_k^{(1,2)} \right) - \left( R_k^{(1,2)} \right)^{-1} \right\|, \\ \|C_k^{(1,2)}\| &\leq \|A_k\| \left( R_k^{(1,2)} \right)^{-1} \left\| \tilde{R}_k^{-1} P_1 \right\| \left\| \tilde{R}_k - \mathcal{H}(\tilde{R}_k) \right\|, \\ \|D_k\| &\leq \frac{\gamma}{1-\gamma} \|A_1 P_1 - I\|,\end{aligned}$$

where  $\|E_{1,k} P_1\| = \gamma < 1$ . We note that we expect  $A_k \left( R_k^{(1,2)} \right)^{-1} \approx A_1$  and therefore using these bounds comes with a hefty prize as the hierarchical approximation error has to balance out a possibly large terms of  $\|A_k\|$  and  $\|A_k \left( R_k^{(1,2)} \right)^{-1}\| \approx \|A_1\|$ . Assuming, in addition, that  $\|P_1 E_{1,k}\| < 1$  we can write

$$\|\tilde{R}_k^{-1} P_1\| \leq \frac{\|P_1\| \|P_1 E_{1,k}\|}{1 - \|P_1 E_{1,k}\|},$$

using, again, the Neumann series expansion bound, see (7). Altogether, we see that our understanding (and hence also the bounds) could be significantly improved if we get an insight into the interplay of the hierarchical approximation and the matrices  $A_k, P_1$ . That comes, inevitably, with a loss of generality and hence we proceed here to give general convergence bounds for the preconditioned GMRES directly. However, adapting these to more specialized settings where further information about the systems can be used seems like worthwhile direction for future research.

The GMRES convergence behavior for the type of systems as in (16) has been already studied – the perturbed system matrix in [20] and the low-rank plus small-norm perturbation to the identity matrix in [6] – using the pseudospectra bounds (PSA). Before applying these results to our situation, we recall that for any  $\delta > 0$  the  $\delta$ -pseudospectrum of a matrix  $M$ , denoted by  $\sigma_\delta(M)$  is defined as

$$\sigma_\delta(M) = \left\{ z \in \mathbb{C} \mid \|(zI - M)^{-1}\| > \frac{1}{\delta} \right\} = \{z \in \sigma(M + E) \text{ for some } E \text{ with } \|E\| < \delta\},$$

and for any  $\delta > 0$  forms a union of Jordan curves surrounding the spectrum of  $M$ , denoted by  $\sigma(M)$ , which can be recovered by taking  $\delta = 0$ . Importantly, considering the standard GMRES convergence bound we can write

$$\frac{\|\mathbf{r}_m\|}{\|\mathbf{r}_0\|} \leq \frac{L_\delta}{2\pi\delta} \min_{\substack{\deg(\varphi) \leq m \\ \varphi(0)=1}} \max_{z \in \sigma_\delta(M)} |\varphi(z)|,$$

where  $L_\delta$  denotes the arc length of the boundary of the  $\delta$ -pseudospectrum of  $M$  and  $\varphi(z)$  is a polynomial of degree up to  $m$  and normalized so that  $\varphi(0) = 1$ ; for more details on pseudospectra we refer the reader to [21] and references therein and for their use in understanding and predicting Krylov subspace methods behavior (GMRES in particular) we refer the reader to [16, Sections 4.9 and 5.7.3] but also to the recent manuscript [8, Section 2.3] and the literature cited in these. We follow the notation in [20] and denote the GMRES residuals corresponding to the initial preconditioned problem (2) by  $\mathbf{r}_m$  ( $m = 1, 2, \dots$ ) while denoting the residuals for the system  $A_k \mathbf{x}_k = \mathbf{b}_k$  with the preconditioner  $P_k^{(1,2)} P_1$  by  $\boldsymbol{\rho}_m^{(1,2)}(k) = \boldsymbol{\rho}_m^{(1,2)}$  ( $m = 1, 2, \dots$ ). Following the calculations in [20, 6], we obtain the following results.



**Proposition 1** ([20, Theorem 2.1]). *Consider a sequence of linear systems  $A_k \mathbf{x}_k = \mathbf{b}_k$  for  $k = 1, 2, \dots$  and adopting the above notation, we fix some  $k > 1$  and  $\star \in \{1, 2\}$ . Let us assume that  $\|E_{1,k}P_1\| < \gamma < 1$ , that the hierarchical format for  $P_k^\star$  was chosen so that*

$$\begin{aligned} \|\mathcal{H}_{\text{solve}}(R_k^\star) - (R_k^\star)^{-1}\| &\leq \varepsilon_B \frac{1}{\|A_k\| \|P_1\|}, \\ \|P_1 E_{1,k} - \mathcal{H}(P_1 E_{1,k})\| &\leq \varepsilon_C \frac{1}{\|A_k (R_k^\star)^{-1}\| \|\tilde{R}_k^{-1} P_1\|}, \end{aligned}$$

and, moreover, that the initial preconditioner  $P_1$  was such that

$$\|A_1 P_1 - I\| \leq \varepsilon_D \frac{1 - \gamma}{\gamma},$$

where  $\varepsilon_{B,C,D} = \varepsilon_{B,C,D}(k, \star)$  but we omit the  $k$  and  $\star$  dependency to make the exposition easier for orientation. Let  $\varepsilon(k, \star) = \varepsilon := \varepsilon_B + \varepsilon_C + \varepsilon_D < 1$ . Then for any  $\delta/2 > \varepsilon$  and all  $m = 1, 2, \dots$  we have

$$\frac{\|\rho_m^\star\|}{\|\rho_0^\star\|} \leq \frac{\|\mathbf{r}_m\|}{\|\mathbf{r}_0\|} + \varepsilon \cdot \frac{L_\delta}{\pi \delta^2} \max_{z \in \sigma_\delta(A_1 P_1)} |p_m(z)|,$$

where  $L_\delta$  denotes the arc length of the boundary of the  $\delta$ -pseudospectrum of  $A_1 P_1$  and  $p_m(z)$  is the optimal GMRES polynomial for the initial preconditioned system (2) at iteration  $m$ .

**Proposition 2** ([6, Equations 5–7]). *Consider a sequence of linear systems  $A_k \mathbf{x}_k = \mathbf{b}_k$  for  $k = 1, 2, \dots$  and adopting the above notation, we fix some  $k > 1$  and  $\star \in \{1, 2\}$ . Let us assume that  $\text{rank}(E_{1,k}) = \nu \ll N$  and that the hierarchical format for  $P_k^\star$  was chosen so that*

$$\begin{aligned} \|\mathcal{H}_{\text{solve}}(R_k^\star) - (R_k^\star)^{-1}\| &\leq \varepsilon_B \frac{1}{\|A_k\| \|P_1\|}, \\ \|P_1 E_{1,k} - \mathcal{H}(P_1 E_{1,k})\| &\leq \varepsilon_C \frac{1}{\|A_k (R_k^\star)^{-1}\| \|\tilde{R}_k^{-1} P_1\|}, \end{aligned}$$

and, moreover, that the initial preconditioner  $P_1$  was such that

$$\|A_1 P_1 - I\| \leq \varepsilon_D,$$

where  $\varepsilon_{B,C,D} = \varepsilon_{B,C,D}(k, \star)$  but we omit the  $k$  and  $\star$  dependency to make the exposition easier for orientation. Let  $\varepsilon(k, \star) = \varepsilon := \varepsilon_B + \varepsilon_C + \varepsilon_D < 1$ . Then for any  $\delta/2 > \varepsilon$  and any  $m = \nu + 1, \nu + 2, \dots$  we have

$$\frac{\|\rho_m^\star\|}{\|\rho_0^\star\|} \leq \varepsilon \cdot \frac{L_\delta}{\pi \delta^2} \max_{z \in \sigma_\delta(A_1 P_1)} |p_m(z)|,$$

where  $L_\delta$  denotes the arc length of the boundary of the  $\delta$ -pseudospectrum of  $A_1 P_1$  and  $p_m(z)$  is the optimal GMRES polynomial for the initial preconditioned system (2) at iteration  $m$ .

First, we highlight that in both Proposition 1 and 2 we describe the GMRES behavior as a function of that of the initial system. More precisely, we describe the *delay* behind the GMRES behavior for the initial preconditioned system (2). Moreover, this delay is spelled out *explicitly* as a function of  $\varepsilon$  but is, seemingly, only *implicit* as a function of the GMRES iteration  $m$  (as highlighted in [20, p. 1071–1072]). As a result, solid grasp on the convergence of the initial system

are an important piece of information to make these results useful, we comment on this further in the following section. Also, both Proposition 1 and 2 rely on some quantification of the interaction of the initial preconditioner  $P_1$  with the update matrix  $E_{1,k}$ , as highlighted in Remark 1. If these assumptions become too crude, e.g.,  $\gamma \lesssim 1$ ,  $\nu \gg 1$  or  $\|A_k\| \gg 1$ , then we expect the convergence bounds to also become quite crude.

Second, following [20], we note that both Proposition 1 and 2 give a *family* of bounds based on  $\delta > 0$ , rather than a single bound. For  $\delta$  too large, the  $\delta$ -pseudospectrum will contain the origin and due to the GMRES optimal polynomial normalization the bounds become not useful. The common knowledge is that larger values of  $\delta$  tend to be more descriptive of the initial convergence phase while the smaller values of  $\delta$  give a more accurate prediction for later stages of GMRES, see Figures 3–4 ahead or the figures in [8, Section 3.1].

Next, we want to point out that the bounds also seems appropriate for the field of values bound (FoV), similarly to [5, Section 2, eqns. (2.10–2.11)] for the SAM approach; see [16, Section 5.7.3] and [8, Section 2.2] for further references on the FoV bounds for GMRES. We address this direction in more detail in Section 4.1 below.

Last but not least, we would like to emphasize that we use the *ideal GMRES* bound, which can fail to capture the observed GMRES behavior to arbitrary level based on the interaction of the system matrix, the right-hand side and the initial guess, see [16, Section 5.7.3] for more details. Although in practice these bounds can be very useful, it is necessary to keep their limitations in mind. We continue by numerically testing and investigating the proposed preconditioner maps on a particular model example.

## 4 Numerical experiments

To illustrate and further explore the results of Section 3 we consider a sequence of 2D non-linear advection-diffusion problems on a unit square  $\Omega := [0, 1] \times [0, 1]$  with piecewise constant coefficients that slowly change, namely

$$\begin{aligned} -(p_k \cdot u_{x_1})_{x_1} - (q_k \cdot u_{x_2})_{x_2} + r \cdot u_{x_1} + s \cdot u_{x_2} &= f \quad \text{in } \Omega, \\ u &= g \quad \text{on } \partial\Omega, \end{aligned} \tag{17}$$

with constant advection coefficients  $r = s = -5$  and piecewise constant diffusion coefficients  $p_k, q_k$

$$p_k(\mathbf{x}) = q_k(\mathbf{x}) = \begin{cases} = 5 & \text{if } \mathbf{x} \in \Omega_k, \\ = 0.1 & \text{otherwise,} \end{cases}$$

where for  $k = 1, \dots, 10$  we have a small square  $\Omega_k \subset \Omega$  that is slowly moving in the  $x_2$  direction upwards as  $k$  increase, see Figure 1. Having the fixed source terms  $f$  and  $g$  as

$$\begin{aligned} f(\mathbf{x}) &= \begin{cases} = 10^4 & \text{if } \mathbf{x} \in [0.405, 0.455] \times [0.705, 0.755], \\ = 0 & \text{otherwise,} \end{cases} \\ g(\mathbf{x}) &= \begin{cases} = 2 \sin(8\pi x_1) & \text{if } x_2 = 0, \\ = 0 & \text{otherwise,} \end{cases} \end{aligned}$$

we obtain a sequence of ten linear problems for  $k = 1, \dots, 10$ .

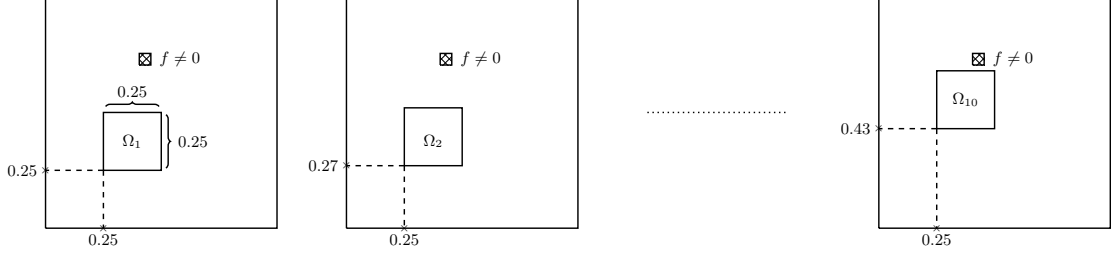


Figure 1: Geometry of the sequence of ten problems in (17).

We discretize (17) using the finite difference method with the standard 5-point stencil, obtaining a sequence of linear algebraic systems as in (1) we try to solve using the preconditioned GMRES method. We notice that all of the update matrices  $E_{\ell,k}$  have comparable number of non-zero columns, i.e., we will experiment with both  $\mathcal{P}^{(1)}$  and  $\mathcal{P}^{(2)}$ . We calculate  $P_1$  using the MATLAB's routine `ilu` with `options = 'crou'` and `droptol = 1e-6`, which we consider an expensive-to-calculate, high-quality preconditioner and we show the number of iterations using  $P_1$  for all ten problems in Figure 2 (left) based on the number of unknowns  $N$ . We see that already for moderate problem sizes there is a large increase in number of iterations for  $k \geq 2$ , making this a good test case for our purposes – we wish to make a better use of  $P_1$  for  $k \geq 2$ . Apart from this feature (which is common for many preconditioners), we do not consciously exploit any particular structure or nature of the preconditioner, i.e., it is just a common choice that is meant to be easily replaceable by a user provided preconditioner.

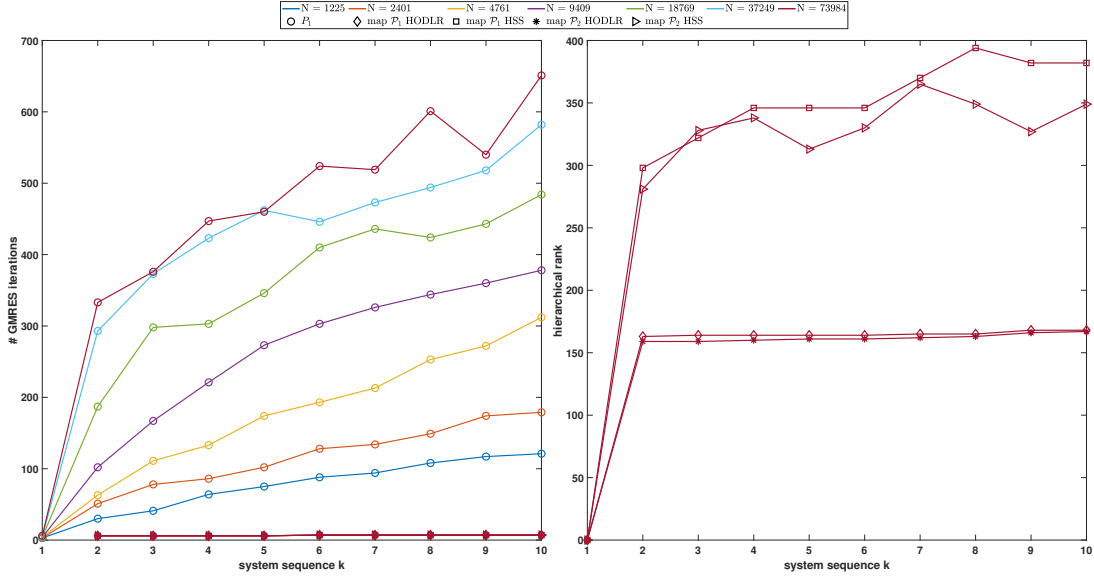


Figure 2: Left: the number of preconditioned GMRES iterations to reach the relative residual  $10^{-10}$  using the preconditioner  $P_1$  (o) for different  $N$  and also using the preconditioner maps  $\mathcal{P}^{(1)}$ ,  $\mathcal{P}^{(2)}$  for hierarchical formats HODLR ( $\diamond, *$ ) and HSS ( $\square, \triangleright$ ). Right: the hierarchical rank of the matrices  $P_k^{(1,2)}$  for the given format with  $N = 73984$ .

For all our experiments, we use the MATLAB implementation of the HODLR and HSS formats given in [18] with necessary adjustments and unless specified otherwise, the parameter settings are left on their default. Using the predefined formats HODLR and HSS directly, we see in Figure 2 (left) that both  $\mathcal{P}^{(1)}$  and  $\mathcal{P}^{(2)}$  keep the number of iterations constant<sup>10</sup>. Although this seems like good news for HAM, Figure 2 (right) shows the hierarchical rank  $r$  of the hierarchical matrices  $R_k^{(*)}$  for  $k = 1, \dots, 10$ , which determines the efficiency of the application of these preconditioners. As the  $\mathcal{H}_{solve}$  operation complexity scales as

$$\mathcal{O}(r^2 N \log^2(N)) \quad \text{and} \quad \mathcal{O}(r^2 N), \quad (18)$$

for the HODLR and HSS formats<sup>11</sup> (see [18, Table 3]) we see that applying these preconditioner is *not* efficient – in fact for both of these formats these complexities become comparable with  $\mathcal{O}(N^2)$ , which is comparable with a backward and forward substitution for some general factored preconditioner (such as  $P_1$ ).

We continue by looking first at the numerical results illustrating the above analysis in Section 4.1 and then continue in Section 4.2 by numerically investigating possible adjustments to the hierarchical formats so that the considered preconditioners become efficient.

## 4.1 HAM bounds

**Pseudospectra** We start with a closer look at the bounds in Proposition 1 and 2 – both of them contain the quantity

$$\frac{\varepsilon}{\delta/2} \tilde{S}_m(\delta), \quad \text{with} \quad \tilde{S}_m(\delta) := \frac{L_\delta}{2\pi\delta} \max_{z \in \sigma_\delta(A_1 P_1)} |p_m(z)|, \quad (19)$$

where  $L_\delta$  denotes the arc length of the boundary of the  $\delta$ -pseudospectrum of  $A_1 P_1$ ,  $p_m(z)$  is the optimal GMRES polynomial for the initial preconditioned system (2) at iteration  $m$  and  $\tilde{S}_m(\delta)$  constitutes an upper bound on the preconditioned GMRES convergence for that problem. Clearly, this quantity needs to be evaluated or further estimated in order to obtain the desired convergence bounds as was already highlighted as a possible bottleneck, see [20, p. 1071–1072]. However, our set-up gives us a powerful tool – we have already run the preconditioned GMRES method for the initial system. In the context of the original work, this is not the case as one is interested in solving *only* the perturbed system, likely not even having access to the unperturbed one. As a result we have at our disposal the  $(m+1)$ -by- $m$  matrix  $H_{m+1,m}$  (with entries  $h_{ij}$ ) coming out of the GMRES applied to (2).

If GMRES terminated at iteration  $\tilde{\mu}$  with  $h_{\tilde{\mu}+1,\tilde{\mu}} = 0$ , then we found an invariant Krylov subspace, which by definition contains all relevant information for understanding the GMRES convergence and the system matrix in (19) (in our case  $A_1 P_1$ ) can be equivalently replaced by  $H_{\tilde{\mu},\tilde{\mu}}$ , i.e., the  $\tilde{\mu}$ -th leading principal submatrix of  $H_{\tilde{\mu}+1,\tilde{\mu}}$ . In practice, this is often ill-advised as we are satisfied with the relative residual being small enough at some earlier iteration  $\mu < \tilde{\mu}$  and as a result we have  $h_{\mu+1,\mu} \neq 0$ . Nevertheless, we use the pseudospectra of  $H_{\mu,\mu}$  (or  $H_{\mu+1,\mu}$ ) as an *approximation* to those of the system matrix *before*  $h_{\mu+1,\mu} = 0$ , obtaining only a GMRES convergence

<sup>10</sup>We show this only for the problem with the largest size – highlighted by the color – but this was true for any problem and size we have experimented with.

<sup>11</sup>Notice that the favorable complexity for the HSS format is balanced out by the notably higher hierarchical rank. This is to be expected has been our experience in general and, as a result, makes both formats competitive.

*estimates* rather than bounds, see [8, Section 4 and Theorem 4.1] and the references therein. In practice, this becomes an efficient tool for evaluating (19) – from the beginning we aim to make good use of a very efficient preconditioner  $P_1$ , meaning that  $\mu \ll N$  and hence  $H_{\mu,\mu}$  (or  $H_{\mu+1,\mu}$ ) are small matrices. As a result, calculating the harmonic Ritz values (which characterize the optimal GMRES polynomial  $p_m(z)$  for  $m = 1, \dots, \mu$ , see [16, Section 5.7.1]) as well as the  $\delta$ -pseudospectra<sup>12</sup>  $\sigma_\delta(H_{\mu,\mu})$  (or  $\sigma_\delta(H_{\mu+1,\mu})$ ) becomes *computationally insignificant* compared to the GMRES method to be run for  $k = 2, 3, \dots$  and so does the evaluation of (19). Hence, we can evaluate these *estimates* for a given  $\varepsilon$  *a-priori*, in theory allowing us to choose appropriate values of  $\varepsilon_{B,C,D}(k, \star)$ . We show the  $\delta$ -pseudospectra of  $H_{\mu,\mu}$  together with the estimates  $S_m(\delta, \mu)$  of  $\tilde{S}_m(\delta)$  with

$$S_m(\delta, \mu) := \frac{L_\delta}{2\pi\delta} \max_{z \in \sigma_\delta(H_{\mu,\mu})} |p_m(z)|, \quad m = 1, 2, \dots,$$

where  $\mu$  depends on the given GMRES relative residual tolerance  $\tau_{\text{GMRES}}$ . We take  $\tau_{\text{GMRES}} = 10^{-8}$  in Figure 3 and as machine precision in Figure 4. We see that the estimates as well as the pseudospectra do not suffer much by stopping GMRES earlier, although they are different and taking  $\tau_{\text{GMRES}} = 10^{-8}$  in Figure 3 resulted in some of the estimates to underestimate the relative residual at iteration  $m = 3$ .

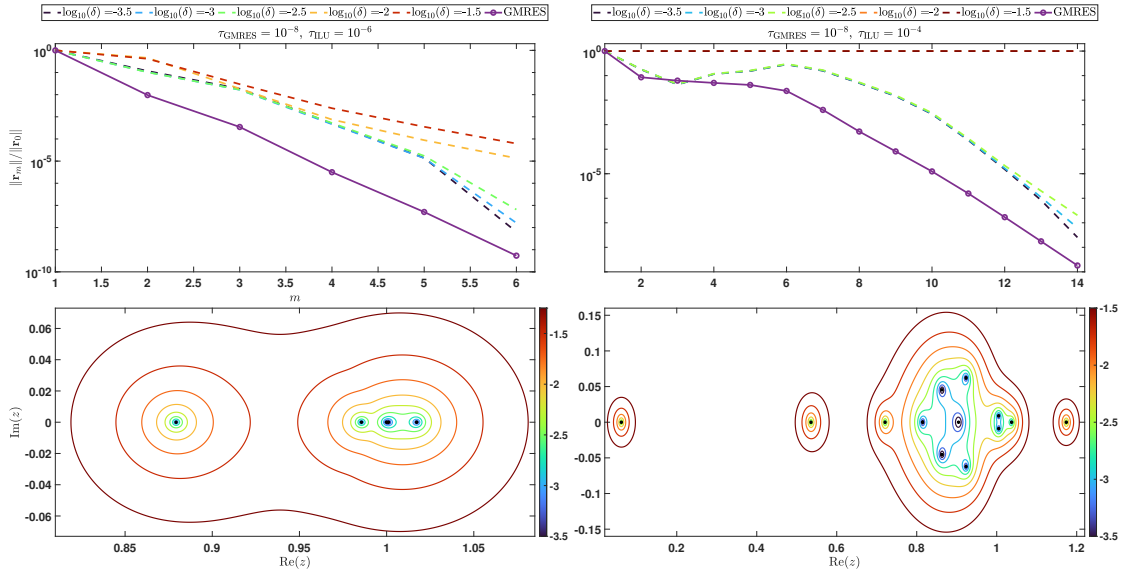


Figure 3: Top: the GMRES convergence together with the PSA bounds  $C_m(\delta)$  for various  $\delta$ ; Bottom: the boundaries of the  $\delta$ -pseudospectra inducing the bounds together with the Ritz values (\*). We fixed  $N = 73984$  and obtained  $P_1$  via the MATLAB routine `ilu(A1)` with `options.type = 'crout'` and `options.droptol =  $\tau_{\text{ILU}}$` . The colorbar for the bottom row gives the decadic logarithm of  $\delta$ , as suggested by the legend at the top.

<sup>12</sup>As of now, the EigTool package ([22]) used in the community as the default tool to calculate pseudospectra does not support calculation with large sparse matrices due to compatibility issues with newer releases of MATLAB and hence we cannot investigate how well the pseudospectra  $\sigma_\delta(H_{\mu,\mu})$  (or  $\sigma_\delta(H_{\mu+1,\mu})$ ) approximate  $\sigma_\delta(A_1 P_1)$ . However, precisely this process, i.e., approximating the large sparse matrix with the small dense matrix  $H_{m+1,m}$  coming out of  $m$  steps of the Arnoldi process, was the fundamental process on which the EigTool was built and it also has become a “default workaround” in the community to obtain the pseudospectra of large sparse matrices until the functionality of EigTool gets restored, see [8, 23, 22].

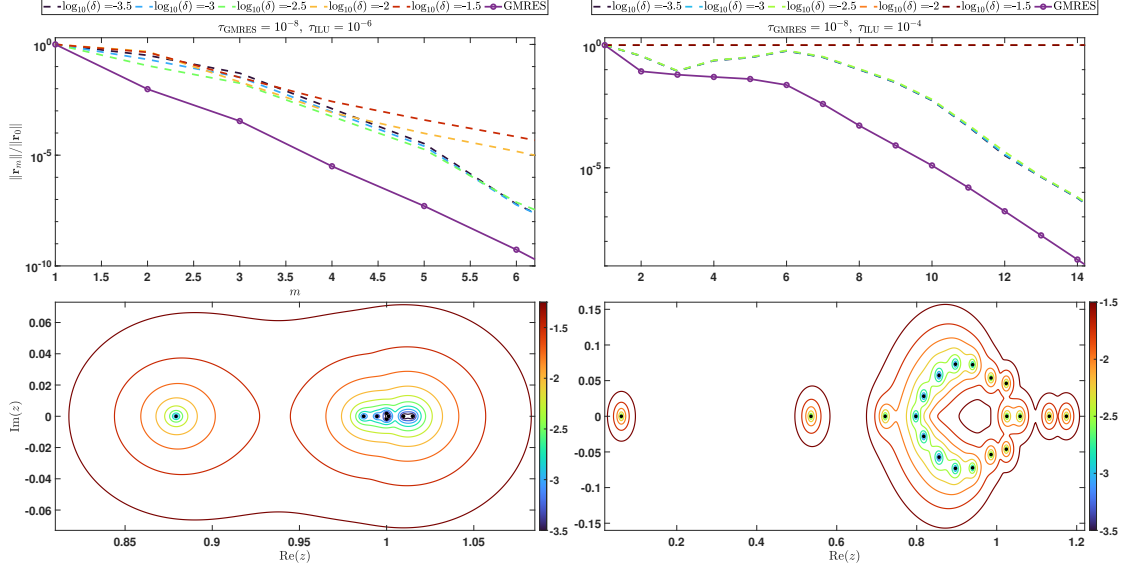


Figure 4: Top: the GMRES convergence together with the PSA bounds  $C_m(\delta)$  for various  $\delta$ ; Bottom: the boundaries of the  $\delta$ -pseudospectra inducing the bounds together with the Ritz values (\*). We fixed  $N = 73984$  and obtained  $P_1$  via the MATLAB routine `ilu(A1)` with `options.type = 'crout'` and `options.droptol =  $\gamma_{\text{ILU}}$` . The colorbar for the bottom row gives the decadic logarithm of  $\delta$ , as suggested by the legend at the top.

Returning to the quantities of interest in Propositions 1 and 2, these correspond to the same curves in Figures 3 and 4, only multiplied by  $2\varepsilon/\delta < 1$ , corresponding to a downward shift by  $|\log_{10}(2\varepsilon/\delta)|$  when looking at the top rows in Figures 3 and 4, obtaining the delay after the GMRES convergence applied to (2) (if Proposition 1 applies) or after  $\nu$  steps of GMRES (if Proposition 2 applies).

This highlights that the hurdle to overcome is in fact the  $\varepsilon$  (and/or  $\nu$ ) parameters we have to settle for when using Propositions 1 and 2 – as highlighted in Section 3, the general scope of the analysis of the preconditioner is the primary reason for the bounds being quite crude, leading to quite an overestimation in this example – we obtain, by a direct calculation,  $\gamma \gg 1$  and  $\nu \approx 50 - 70$ , meaning that the convergence bounds are no good in describing GMRES convergence behavior before iteration 50. This is a meaningful drawback to improve on and we aim to do that elsewhere, in full length. Here, instead, we touch upon the field of values bound below and then turn our attention to the issue of *efficiency* raised at the beginning of Section 4, thus finalizing the introduction of the HAM approach.

**Field of values** Building upon Section 3, we look at the contributions of the terms in (8) and (10) to enlarging the field of values of the preconditioned system. We use the MATLAB toolbox function `fv` to bound and visualize the field of values (FoVs) of matrices, see [14]. We show the FoVs of the three parts of (16) in Figure 5 for  $\mathcal{P}^{(1)}$  and in Figure 6 for  $\mathcal{P}^{(2)}$ .

We see that the FoVs of  $B_k^{(1,2)}, C_k^{(1,2)}$  are centered around the origin, as these capture the approximation error that converges to zero matrix as we decrease  $\epsilon$  while the FoVs of  $D_k$  are centered around the point  $1 + 0i$  as expected based on (11). However, we see that for both  $\mathcal{P}^{(1,2)}$ ,

the overall FoVs of the preconditioned system for  $k \geq 2$  *includes the origin* for  $\epsilon = 10^{-4}, 10^{-2}$  – making the bound based on them useless. We note that these could still be, in theory, salvaged by the “cutting-out” technique of Crouziex and Greenbaum, see [7]. However, this approach is, in our opinion, *highly* unlikely to be successful here as the FoVs are dominated by the second line terms, i.e., by FoVs of  $C_k^{(1,2)}$ , which are *centered around the origin*. We also note that the dominating term has always been, in our experience, either the FoVs corresponding to  $C_k^{(1,2)}$  – the term due to the error of the hierarchical approximation<sup>13</sup> – or to  $D_k$  – the term due to the quality of  $P_1$ . Moreover, as  $N$  grows, the diameter of the FoVs of  $C_k^{(1,2)}$  grows significantly as well, i.e., the term corresponding to  $D_k$  dominated only for problems of small size. Last, we note that changing the format to HSS either did not change the FoV meaningfully or even enlarged the diameters of the FoVs in all of our numerical experiments.

Although these results are negative, they are not completely surprising – the FoV bounds are generally expected to be at best as descriptive as the bounds based on pseudospectra but *often much less so*, see [8, Section 3.1 – No Example: Only FOV descriptive]. We also note that based on (8) and (10) we can give a proper bound, as suggested already in [5, Section 2, eqn. (2.8)] for SAM but based on the above insight, we see that such bound is not of much interest in general.

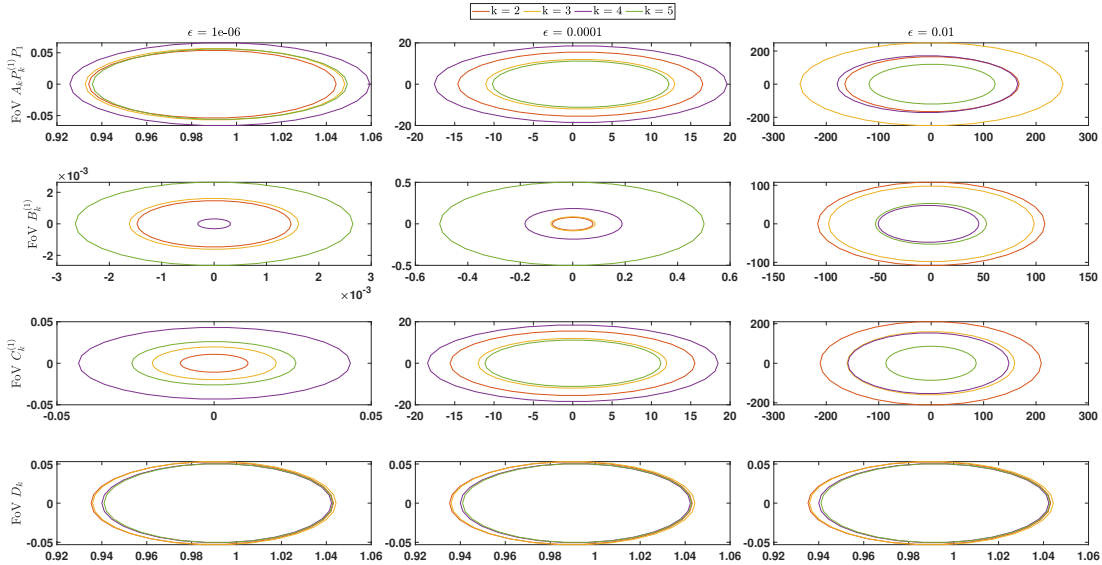


Figure 5: The bounds of FoVs for the three parts of the preconditioned systems for  $k = 2, \dots, 5$  for the map  $\mathcal{P}^{(1)}$  different values of  $\epsilon$ . We show the results only for the HODLR format and fixed  $\beta = 128$  and  $N = 18769$ .

<sup>13</sup>Notice that this term differs meaningfully between the maps  $\mathcal{P}^{(1,2)}$  – we predicted this in the analysis in Section 3 and now confirmed it in Figures 5 and 6.

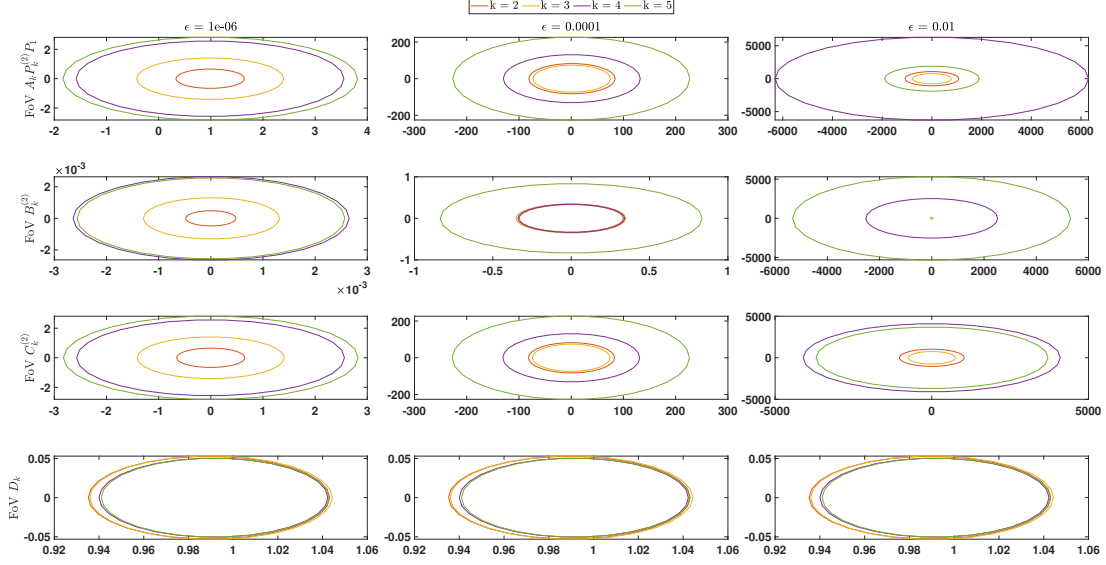


Figure 6: The bounds of FoVs for the three parts of the preconditioned systems for  $k = 2, \dots, 5$  for the map  $\mathcal{P}^{(2)}$  different values of  $\epsilon$ . We show the results only for the HODLR format and fixed  $\beta = 128$  and  $N = 18769$ .

## 4.2 HAM efficiency

We return to the issue highlighted at the beginning of this section, namely that applying the HAM preconditioners  $P^{(1,2)}$  is, complexity-wise, comparable to a Gaussian elimination for a pre-factored preconditioner. First, let us emphasize, that part of the problem we are addressing is the fact that *we do not have* that pre-factored preconditioner but we acknowledge that such a high complexity is not feasible if we want to have a competitive preconditioner scheme. To understand why the hierarchical rank becomes so high, we take a closer look at the predefined settings of the MATLAB `hm-toolbox` by Massei and Kressner we used, see [18].

In their implementation, both HODLR and HSS use the basic balanced binary tree as the cluster tree with the standard admissibility condition and therefore either of these formats is further characterized by only two parameters – the minimal block-size  $\beta \in \mathbb{N}$  and the accuracy threshold  $\epsilon$ . Increasing  $\beta$  corresponds to stopping the hierarchical blocking of the matrix at an earlier stage and therefore preserving larger diagonal blocks of the original matrix exact. Decreasing  $\epsilon$  corresponds to requiring a more accurate overall approximation within the chosen format by virtue of increasing *the hierarchical rank* of the approximation (for more detailed description, we refer to the literature cited in Section 3 and to [18] and the references therein). The default values are set to  $\beta = 256$  and  $\epsilon = 10^{-12}$ . While these might be perfectly reasonable as a standalone, the first natural step in our context is investigate the effect of relaxing these on the hierarchical rank of the resulting matrices. We note that although there is some relevant theory linking  $\epsilon$  with the hierarchical rank  $r$ , see, e.g., , to the best of our knowledge all analysis in this direction is considered *in the limit as  $\epsilon \rightarrow 0$*  – exactly the opposite of what we want to look at. Hence we resort to numerical investigation.

**Effect of  $\beta$  and  $\epsilon$**  First, we note that changing the block size  $\beta$  had little to no effect on the quality of the preconditioners. We ran extensive amount of experiments, varying both  $\beta \in [32, 2048]$



as well as the system size  $N \in (10^3, 10^5)$  for both maps  $\mathcal{P}^{(1)}, \mathcal{P}^{(2)}$  and for both formats HODLR and HSS and found that both the rank as well as the number of GMRES iterations are quite stable with respect to changing  $\beta$ , even though these do change a little more for  $\beta = 1024, 2048$  – lowering the number of iterations and varying the hierarchical rank (notably but not significantly). Looking at the runtimes, our experience was that  $\beta$  and  $N$  should be kept *proportional* so that the cluster tree has more or less constant number of levels. We are, however, aware of implementation limitations of the `hm-toolbox` in MATLAB and more experiments and analysis is needed in this direction.

In contrast to that, changing the accuracy  $\epsilon$  had a significant impact on both the number of iteration of preconditioned GMRES as well as on the hierarchical rank and we illustrate this in Figures 7–8. We see a notable difference for most of the values of  $\epsilon$  – smaller values result in more accurate approximations and hence lower number of iterations but require a higher hierarchical rank. Altogether, we see that the number of preconditioned GMRES iterations remains more or less constant for  $k = 1, \dots, 10$  with much lower hierarchical rank compared to Figure 2 and with a *significant* improvement compared to using  $P_1$ , even for a quite poor accuracy  $\epsilon = 0.1$ , see Figure 2.

Comparing the maps  $\mathcal{P}^{(1)}$  and  $\mathcal{P}^{(2)}$ , we see that  $\mathcal{P}^{(1)}$  results in lower number of GMRES iterations but requires a higher hierarchical rank to do it. Recalling (18), a fairer comparison is to compare the number of iterations times the cost per iteration. Denoting the number of GMRES iterations with the map  $\mathcal{P}^{(i)}$  by  $n^{(i)}$  and the resulting hierarchical rank by  $r^{(i)}$ , the map  $\mathcal{P}^{(1)}$  is more efficient provided that

$$(n^{(2)}/n^{(1)})^{-1/2} < r^{(2)}/r^{(1)}. \quad (20)$$

However, even in this better comparison, the map  $\mathcal{P}^{(1)}$  is favorable to the counterpart, in spite of the higher hierarchical rank. This is, perhaps, to be expected due to the lack of error accumulation compared  $\mathcal{P}^{(2)}$ . Since the number of nonzero columns for  $E_{1,k}$  and  $E_{k-1,k}$  is roughly the same, there is no direct counterweight to this benefit – the somewhat higher hierarchical ranks of  $\mathcal{P}^{(1)}$  are well compensated for.

Comparing the HODLR and HSS formats, we see that the hierarchical rank *as well as* the number of GMRES iterations for the HSS formats tends to be higher for the HSS format for both  $\mathcal{P}^{(1)}, \mathcal{P}^{(2)}$ . However, recalling the complexities in (18), we see that this does not necessarily make the HSS format less efficient as the extra term of  $\log^2(N)$  in the HODLR complexity can offset these differences. In fact, for  $N \gtrsim 5000$ , the HSS format runtimes<sup>14</sup> of the results plotted on Figures 7–8 are comparable or even significantly lower than if using the HODLR format, as  $N$  increases.

---

<sup>14</sup>These need to be taken with the same caveats as above with respect to performance of the `hm-toolbox` and MATLAB. Also, the  $\mathcal{H}$ -FAINV approach would likely change this as the scaling of the application of the preconditioner with respect to the hierarchical rank would become linear compared the currently quadratic one.

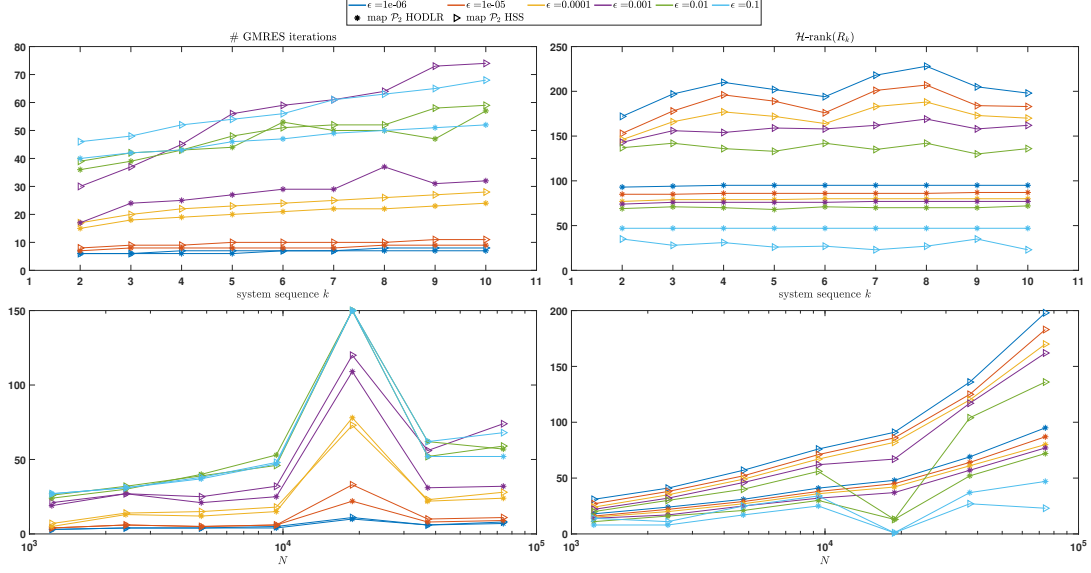


Figure 8: Number of iterations of the preconditioned GMRES (left) and the hierarchical rank of the matrix  $R_k$  (right) plotted against the system sequence (top) and system size (bottom) for the map  $\mathcal{P}^{(2)}$  for both HODLR and HSS formats and various values of  $\epsilon$ . Here we fixed  $\beta = 256$  and  $N = 73984$  for the top row and  $k = 10$  for the bottom row.

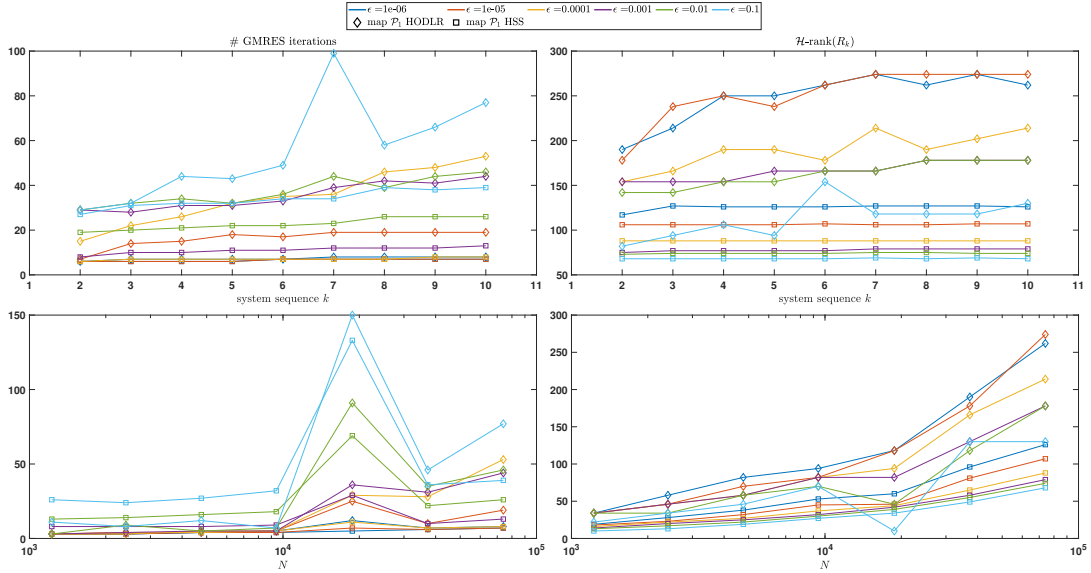


Figure 7: Number of iterations of the preconditioned GMRES (left) and the hierarchical rank of the matrix  $R_k$  (right) plotted against the system sequence (top) and system size (bottom) for the map  $\mathcal{P}^{(1)}$  for both HODLR and HSS formats and various values of  $\epsilon$ . Here we fixed  $\beta = 256$  and  $N = 73984$  for the top row and  $k = 10$  for the bottom row.

Last but not least, we see that with growing  $N$  both the number of iterations as well as the

hierarchical ranks grow. This is not surprising as incomplete factorization preconditioners, such as  $P_1$ , are known to exhibit this behavior, i.e., we cannot hope for anything better in terms of GMRES iterations in general. Importantly, looking at the hierarchical rank growth, we are, unfortunately, still on the same  $N^2$ -like trajectory for almost all values of  $\epsilon$ , regardless of the format and map used<sup>15</sup>. In order to address this issue, we look for further approximations with the goal of decreasing the hierarchical rank as well as the cost of construction of the preconditioner. For the SAM approach, the authors explored a similar direction by prescribing a *fixed sparsity pattern* of the preconditioner (and hence putting the complexity of the application of the preconditioner under the user control) and it turned out that the GMRES iteration count did not suffer significantly even from restrictive sparsity patterns (compared to the exact solution). We consider a similar approach in the following two sections – first *structural* sparsification, where some non-zero entries are dropped, and then *data* sparsification, where we truncate the hierarchical ranks.

**Structural sparsification** In our model problem the update matrix  $E$  (standing for either  $E_{1,k}$  or  $E_{k-1,k}$  depending on the chosen map) has only relatively few non-zero columns, each of which has only very few non-zero entries, as is the case in *many* applications. Having  $E$  structurally sparse, the idea of further sparsification seems at first natural in order to reduce the amount of information that needs to be captured by the hierarchical format approximation. In fact, there are also two options for this – we can either sparsify the update matrix  $E$  itself or the matrix  $P_1 E$  before we apply calculate its hierarchical approximation<sup>16</sup> – and those are the only options as once we calculate  $\mathcal{H}(P_1 E)$ , there is no possibility for dropping non-zero elements of that matrix. We investigate the efficiency of both of the options numerically by dropping all entries that are smaller in magnitude than a fixed tolerance  $\tau$  and show representative examples of our experience in Figure 9 (sparsifying the matrix  $E$ ) and in Figure 10 (sparsifying the matrix  $P_1 E$ ).

First, we see that there is an important distinction between sparsifying  $E$  and  $P_1 E$  – the former in fact does decrease the hierarchical rank – sometimes even significantly – while the latter does the opposite – it increases the hierarchical rank, and quite noticeably. This is a key difference that can be heuristically understood on the following example: consider the matrices

$$M = \begin{bmatrix} 10 & 5 & 2.5 \\ 1 & 0.5 & 0.25 \\ 0.1 & 0.05 & 0.025 \end{bmatrix} = \begin{bmatrix} 10 \\ 1 \\ 0.1 \end{bmatrix} \begin{bmatrix} 1 & 0.5 & 0.25 \end{bmatrix} \quad \text{and} \quad M_{\text{sp}} = \begin{bmatrix} 10 & 5 & 2.5 \\ 1 & 0.5 & 0 \\ 0.1 & 0 & 0 \end{bmatrix},$$

where  $M_{\text{sp}}$  is obtained by sparsifying  $M$  with the drop tolerance  $\tau = 1/3$ . When we sparsify  $P_1 E$ , the analogue of this occurs blockwise and *destroys* the low-rank structure that is naturally present due to the ellipticity of the underlying differential operator: the approximability in general stems from large blocks of the matrix being almost constant with only few outliers and hence dropping entries can disrupt these almost constant blocks and even a relatively small drop tolerance can make the hierarchical rank explode, see [4] but also [19, Chapters 1–4] and the references therein. Notice that this problem does not appear when we sparsify *before* we apply  $P_1$ , i.e., before we let the discretized differential operator act on the update matrix – quite on the contrary. In the case of

<sup>15</sup>We note that the problem with  $N = 18769$  was in some sense exceptional as the system matrices were consistently easier to approximate (resulting in a lower hierarchical rank) but resulted in a worse preconditioner (hence higher number of iterations).

<sup>16</sup>Notably, the later option would be harder to implement in the case of matrix-free assembly of the hierarchical format mentioned in Section 3. However, as we will see, there doesn't seem to be a reason to be concerned about this issue.

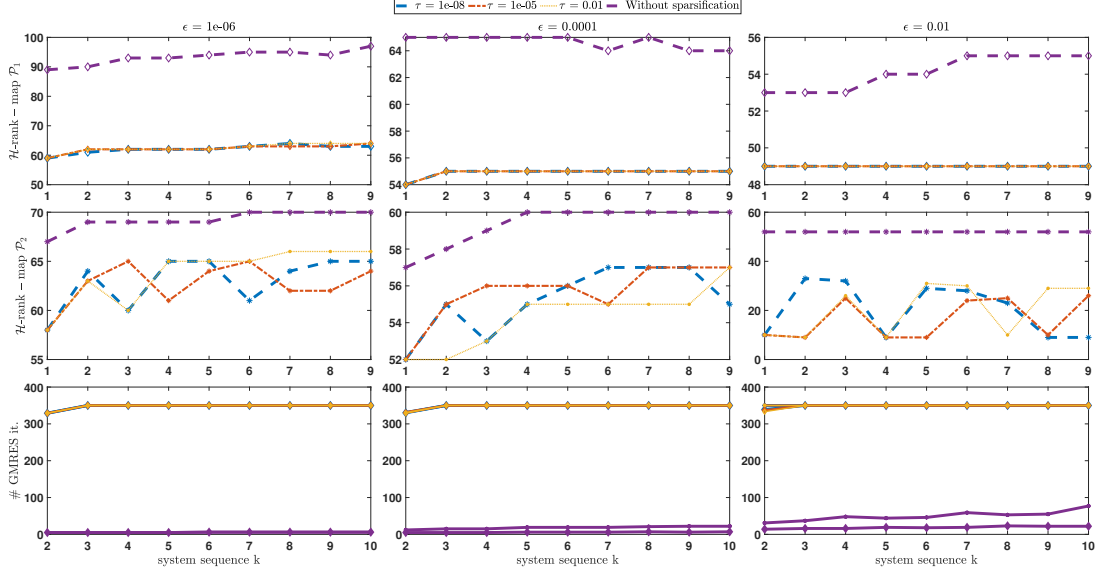


Figure 9: We show the hierarchical ranks of the resulting matrices  $R_k$  for the HODLR format when sparsifying the matrix  $E_{1,k}$  ( $E_{k-1,k}$ ) for the map  $\mathcal{P}^{(1)}$  ( $\mathcal{P}^{(2)}$ ) in the first row (second row) and the number of GMRES iterations (capped at 350) with the sparsified preconditioner compared to no sparsification. We fixed  $\beta = 256$  and  $N = 37249$ .

sparsifying  $E$ , our initial intuition of decreasing the amount of information to approximate follows through. To be more precise, the non-zero columns of  $P_1 E$  are linear combinations of columns of  $P_1$  with coefficients stored in the columns of  $E$  and hence the sparsification of  $E$  moderates the number of the terms in these linear combinations. Assuming  $P_1$  mimics the well-approximability in a hierarchical format of  $A_1^{-1}$ , having fewer of these columns combined naturally leads to lower hierarchical rank.

However, we see that the information decrease has a *drastic* effect on the quality of the preconditioner in terms of GMRES iterations – we capped the maximal number of GMRES iterations at 350 and only the second system converged, taking about 330 iterations to reach the  $10^{-10}$  relative residual. Compared to the non-sparsified preconditioners and keeping (20) in mind, we see that the structural sparsification – either of  $E$  or  $P_1 E$  – is unlikely to improve the overall efficiency of our preconditioners, based on both, the shown data as well as our general experience over many more experiments. We note that these observations remained true when varying other relevant parameters, such as the format (for brevity we do not show the analogous plots for the HSS format), accuracy  $\epsilon$ , the drop tolerance  $\tau$  and the size of the system  $N$  (with the caveat that for small systems, these effects are significantly damped).

**Data sparsification** As the structural sparsification did not yield the desired effect, we turn our attention to *data-sparsification* – truncation of the hierarchical rank of the considered matrices. This direction is the proper analogue of the fixed (structural) sparsity pattern in the SAM approach mentioned above – we replace the structural sparsity in both the map and in its further restriction by data-sparsity. The truncation for the HODLR format then consists of blockwise truncation, following the blocking of the format. For the HSS format, things are a bit more involved as the

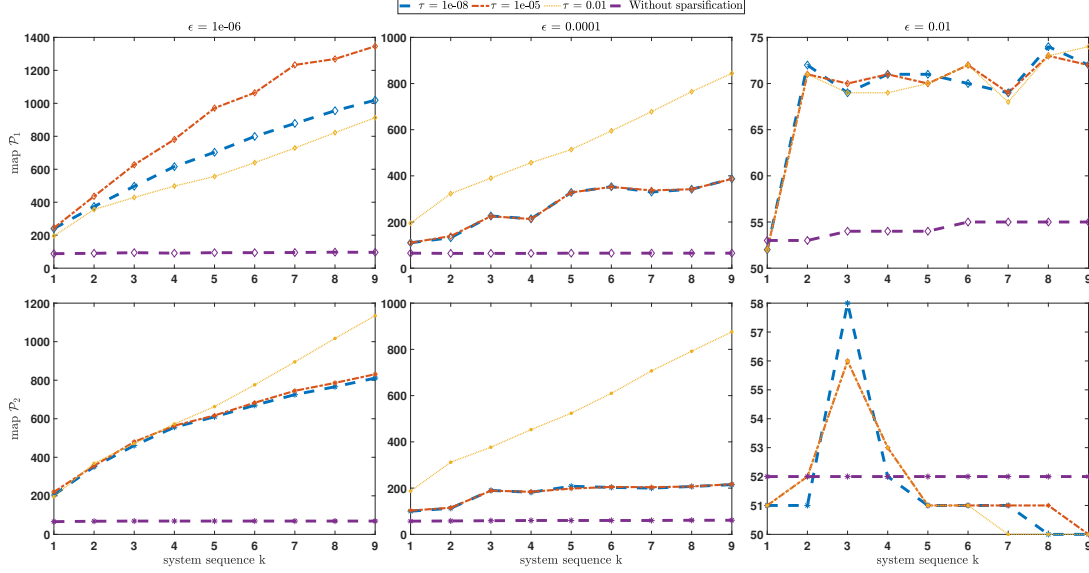


Figure 10: The hierarchical ranks of the resulting matrices  $R_k$  for the HODLR format when sparsifying the matrix  $P_1 E_{1,k}$  or  $P_1 E_{k-1,k}$ . We fixed  $\beta = 256$  and  $N = 37249$ .

format itself is more delicate – following the notation of the `hm-toolbox` in [18] we calculated the SVD of the *core matrices*  $S_{ij}^{(\ell)}$  stored in the `B12`, `B21` attributes of the HSS class and truncated it (as well as the other corresponding matrices in the format in order to reap the computational efficiency benefits of the truncation).

Similarly to the analysis in Section 3, we have only one option for hierarchical rank truncation to  $r_{\max}$  for the map  $\mathcal{P}_1$  – the matrix  $R_k$  – but we have two options for the map  $\mathcal{P}_2$  – the matrix  $\mathcal{H}(P_1 E_{k-1,k})$  or the matrix  $R_k$ . Clearly, only the second option guarantees the rank to be lower than  $r_{\max}$  as the formatted sum  $R_{k-1} \oplus \text{ranktrunc}(\mathcal{H}(P_1 E_{k-1,k}))$  can increase the hierarchical rank – in fact for this option we can bound the hierarchical rank is by  $2^k r_{\max}$  (and since the complexities scale with the square<sup>17</sup> of the rank, it can be up to  $4^k$  times as expensive). In practice we have never observed anything close to this but we have observed that the hierarchical rank does somewhat increase as  $k$  grows and sometimes can even outgrow the hierarchical rank without any truncation – for the second option, that is. We illustrate this as well as the performance of the resulting preconditioners in Figure 11.

<sup>17</sup>As mentioned above, this can be relaxed to linear instead of square by adding the  $\mathcal{H}$ -FAINV step.

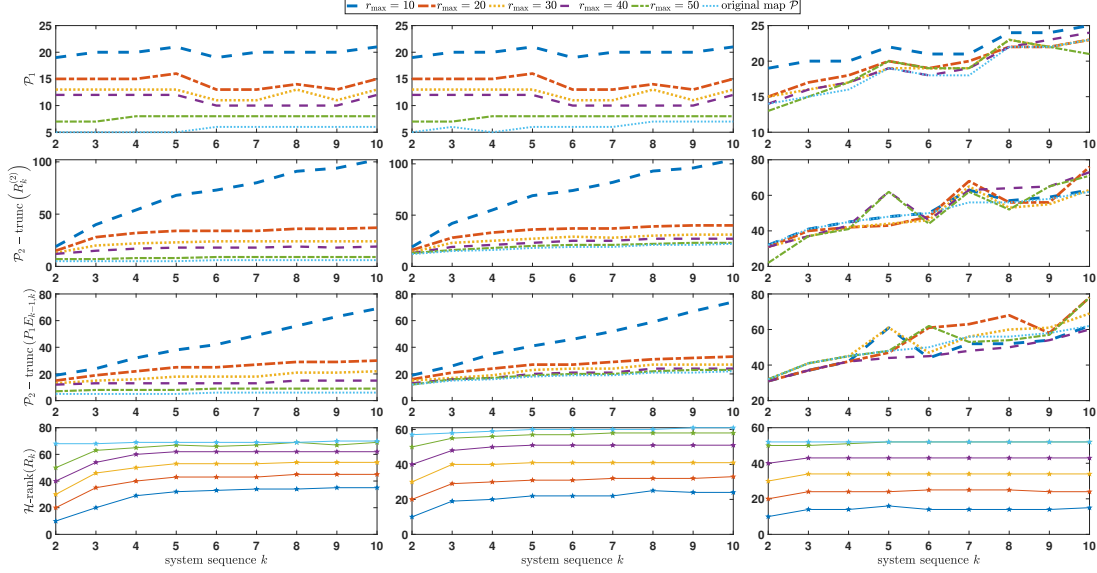


Figure 11: The number of GMRES iterations with fixed hierarchical ranks for  $\mathcal{P}^{(1)}$  (first row) and  $\mathcal{P}^{(2)}$  (second row) and for  $\mathcal{P}^{(2)}$  when truncating the hierarchical rank of only  $P_1 E_{k-1,k}$  (third row). The resulting hierarchical rank of  $R_k^{(2)}$  when truncating only  $P_1 E_{k-1,k}$  as well as that of the hierarchical rank of  $R_k^{(2)}$  without truncation (original map) is shown below (fourth row). We show the results only for the HODLR format and fixed  $\beta = 256$  and  $N = 37249$ .

We see that the number of GMRES iterations compared to the non-truncated preconditioner increases and this increase can be quite mild (plot in the (1,1) position) or very pronounced (plot in the (2,2) position) depending on the way the truncation is carried out and on the map. On one hand, we have observed that the map  $\mathcal{P}_1$  is quite stable with respect to the hierarchical rank truncation, and decreasing the hierarchical rank 2-5 times still gives convergence in a comparable, say double, number of iterations. Recalling (20), this corresponds to a significant improvement in efficiency of the truncated preconditioner. On the other hand, we observed that the map  $\mathcal{P}_2$  is more sensitive to the rank truncation and especially so if the truncation is applied to  $\mathcal{H}(P_1 E)$  rather than to  $R_k$ . Also, we observe that the higher the accuracy (i.e., the smaller the  $\epsilon$ ), the more diverse can these differences be – partly also because higher accuracy translates naturally to higher hierarchical ranks, which in turn leaves more space for highlighting the effects.

An explanation or at least a deeper insight into the interaction of the data-sparsification and the preconditioner maps are left as an open direction for future research.

## 5 Conclusion and future work

We have proposed and studied – both analytically and experimentally – a new type of preconditioner maps for sequences of systems of linear algebraic equations, using data-sparsity or, to be more precise, using the hierarchical matrix formats HODLR and HSS. The theoretical results give a complementary approach to analysis of preconditioner maps compared to [5] and give a direct insight into which quantities can ensure a solid GMRES convergence. We illustrated the need for preconditioner maps on a model example and pursued further avenues to make our maps more

efficient so that application of the preconditioner is manageable. Further analysis as well as analysis for more particular settings is of clear interest and remains open for future research.

## References

- [1] K. Ahuja, B. K. Clark, E. de Sturler, D. M. Ceperley, and J. Kim. Improved scaling for quantum Monte Carlo on insulators. *SIAM Journal on Scientific Computing*, 33(4):1837–1859, 2011.
- [2] J. Ballani and D. Kressner. *Matrices with hierarchical low-rank structures*, chapter Exploiting Hidden Structure in Matrix Computations: Algorithms and Applications, pages 161–209. Springer, 2016.
- [3] M. Bebendorf and T. Fischer. On the purely algebraic data-sparse approximation of the inverse and the triangular factors of sparse matrices. *Numerical Linear Algebra with Applications*, 18(1):105–122, 2011.
- [4] M. Bebendorf and W. Hackbusch. Existence of  $\mathcal{H}$ -matrix approximants to the inverse FE-matrix of elliptic operators with  $L^\infty$ -coefficients. *Numerische Mathematik*, 95(1):1–28, 2003.
- [5] A. Carr, E. de Sturler, and S. Gugercin. Preconditioning parametrized linear systems. *SIAM Journal on Scientific Computing*, 43(3):A2242–A2267, 2021.
- [6] A. K. Carr, E. de Sturler, and M. Embree. Analysis of GMRES for low-rank and small-norm perturbations of the identity matrix. volume 22, page e202200267. Wiley Online Library, 2023.
- [7] M. Crouzeix and A. Greenbaum. Spectral sets: numerical range and beyond. *SIAM Journal on Matrix Analysis and Applications*, 40(3):1087–1101, 2019.
- [8] M. Embree. How descriptive are GMRES convergence bounds? 2022.
- [9] B. Engquist and L. Ying. Sweeping preconditioner for the Helmholtz equation: Hierarchical matrix representation. *Communications on Pure and Applied Mathematics*, 64(5):697–735, 2011.
- [10] L. Grasedyck, R. Kriemann, and S. Le Borne. Parallel black box  $\mathcal{H}$ -LU preconditioning for elliptic boundary value problems. *Computing and Visualization in Science*, 11(4-6):273–291, 2008.
- [11] L. Grasedyck, R. Kriemann, and S. Le Borne. Domain decomposition based-LU preconditioning. *Numerische Mathematik*, 112(4):565–600, 2009.
- [12] W. Hackbusch. A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices. Part I: Introduction to  $\mathcal{H}$ -matrices. *Computing*, 62(2):89–108, 1999.
- [13] W. Hackbusch. *Hierarchical Matrices: Algorithms and Analysis*, volume 49. Springer Berlin, Heidelberg, 2015.
- [14] N. J. Higham. The Matrix Computation Toolbox for MATLAB (version 1.0). Technical report, 2002.

- [15] R. Kriemann and S. Le Borne.  $\mathcal{H}$ -FAINV: Hierarchically factored approximate inverse preconditioners. *Computing and Visualization in Science*, 17(3):135–150, 2015.
- [16] J. Liesen and Z. Strakoš. *Krylov Subspace Methods: Principles and Analysis*. Oxford University Press, Oxford, 2013.
- [17] P. G. Martinsson. A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1251–1274, 2011.
- [18] S. Massei, L. Robol, and D. Kressner. hm-toolbox: MATLAB software for HODLR and HSS matrices. *SIAM Journal on Scientific Computing*, 42(2):43–68, 2020.
- [19] M. Outrata. *Schwarz methods, Schur complements, preconditioning and numerical linear algebra*. PhD thesis, University of Geneva, Math Department, 2022.
- [20] J. A. Sifuentes, M. Embree, and R. B. Morgan. GMRES convergence for perturbed coefficient matrices, with application to approximate deflation preconditioning. *SIAM Journal on Matrix Analysis and Applications*, 34(3):1066–1088, 2013.
- [21] L. N. Trefethen and M. Embree. *Spectra and Pseudospectra: The Behaviour of Non-Normal Matrices and Operators*. Princeton University Press, Princeton, New Jersey, 2005.
- [22] T. G. Wright. Eigtool. <http://www.comlab.ox.ac.uk/pseudospectra/eigtool/>, 2002.
- [23] T. G. Wright and L. N. Trefethen. Large-scale computation of pseudospectra using ARPACK and `eigs`. *SIAM Journal on Scientific Computing*, 23(2):591–605, 2001.
- [24] J. Xia, S. Chandrasekaran, M. Gu, and X. S. Li. Fast algorithms for hierarchically semiseparable matrices. *Numerical Linear Algebra with Applications*, 17(6):953–976, 2010.
- [25] X. Ye, J. Xia, and L. Ying. Analytical low-rank compression via proxy point selection. *SIAM Journal on Matrix Analysis and Applications*, 41(3):1059–1085, 2020.