

# REDUCING PARALLEL COMMUNICATION IN ALGEBRAIC MULTIGRID THROUGH SPARSIFICATION\*

AMANDA BIENZ<sup>†</sup>, ROBERT D. FALGOUT<sup>‡</sup>, WILLIAM GROPP<sup>§</sup>, LUKE N. OLSON<sup>¶</sup>,  
AND JACOB B. SCHRODER<sup>||</sup>

**Abstract.** Algebraic multigrid (AMG) is an  $\mathcal{O}(n)$  solution process for many large sparse linear systems. A hierarchy of progressively coarser grids is constructed that utilize complementary relaxation and interpolation operators. High-energy error is reduced by relaxation, while low-energy error is mapped to coarse-grid matrices and reduced there. However, large parallel communication costs often limit parallel scalability. As the multigrid hierarchy is formed, each coarse matrix is formed through a triple matrix product. The resulting coarse grids often have significantly more nonzeros per row than the original fine-grid operator, thereby generating high parallel communication costs associated with sparse matrix-vector multiplication (SpMV) on coarse levels. In this paper, we introduce a method that systematically removes entries in coarse-grid matrices after the hierarchy is formed, leading to improved communication costs. We sparsify by removing weakly connected or unimportant entries in the matrix, leading to improved solve time. The main trade-off is that if the heuristic identifying unimportant entries is used too aggressively, then AMG convergence can suffer. To counteract this, the original hierarchy is retained, allowing entries to be reintroduced into the solver hierarchy if convergence is too slow. This enables a balance between communication cost and convergence, as necessary. In this paper we present new algorithms for reducing communication and present a number of computational experiments in support.

**Key words.** multigrid, algebraic multigrid, non-Galerkin multigrid, high performance computing

23 AMS subject classifications. 65F50, 65N55, 65F08, 65F10, 65Y05, 68Q10

**1. Introduction.** Algebraic multigrid (AMG) [17, 7, 18] is an  $\mathcal{O}(n)$  linear solver for standard discretizations of elliptic differential equations [2, 23, 20]. We consider AMG as a solver for the symmetric, positive definite matrix problem

$$27 \quad (1) \qquad \qquad \qquad Ax = b,$$

with  $A \in \mathbb{R}^{n \times n}$  and  $x, b \in \mathbb{R}^n$ . AMG consists of two phases, a setup and a solve phase. The setup phase defines a sequence or *hierarchy* of  $L$  coarse-grid and interpolation operators,  $A_1, \dots, A_L$  and  $P_0, \dots, P_{L-1}$  respectively. The solve phase iteratively improves the solution through relaxation and coarse-grid correction.

---

\*This work is partially supported by the National Science Foundation Graduate Research Fellowship award DGE-1144245 and the Air Force Office of Scientific Research under grant FA9550-12-1-0478. This research is also part of the Blue Waters sustained-petascale computing project, which is supported by the National Science Foundation (awards OCI-0725070 and ACI-1238993) and the state of Illinois. Blue Waters is a joint effort of the University of Illinois at Urbana-Champaign and its National Center for Supercomputing Applications. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 (LLNL-JRNL-673388).

<sup>†</sup>Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801, bienz2@illinois.edu, <http://web.engr.illinois.edu/~bienz2>

<sup>†</sup>Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, CA 94550, rfalgout@llnl.gov, <http://people.llnl.gov/falgout2>

<sup>8</sup>Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801,  
wgropp@illinois.edu, <http://wgropp.cs.illinois.edu>

<sup>¶</sup>Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801, lukeo@illinois.edu, <http://lukeo.cs.illinois.edu>

The focus of this paper is on the communication complexity of AMG in a distributed memory, parallel setting. To be clear, we refer to the *communication complexity* as the time cost of interprocessor communication, while referring to the *computational complexity* as the time cost of the floating point operations. The *complexity* or *total complexity* is then the cost of the algorithm, combining the communication and computational complexities.

There is a trade-off between per-iteration complexity and the resulting convergence — this is controlled by the setup phase in AMG. Indeed, more accurate interpolation leads to more entries in  $P_\ell$  and this often improves the overall convergence. However, this also leads to more entries in  $A_{\ell+1}$  through the Galerkin product,  $A_{\ell+1} = P_\ell^T A P_\ell$ , which is the most common way to form  $A_{\ell+1}$  in AMG. As we will see, the number of entries in  $A_{\ell+1}$  is a good proxy for the complexity of level  $\ell + 1$ . In contrast, sparser interpolation and fast coarsening reduce the complexity of a single iteration of an AMG cycle through fewer entries in  $A_{\ell+1}$  and fewer overall AMG levels, but can lead to a deterioration in convergence [23, 20].

The sparse matrices,  $A_1, \dots, A_L$ , in the multigrid hierarchy are, by design, smaller in dimension, yet are often more dense as the level increases. As an example of this, Table 1 shows the properties of a hierarchy for a 3D Poisson problem with a 27-point finite element stencil on a  $100 \times 100 \times 100$  grid. Classical AMG (Ruge-Stüben) is used for this example<sup>1</sup>. As the problem size decreases on coarse levels, the average number of nonzero entries per row increases. Here we denote by `nnz` the number of nonzero entries in the respective matrix. Figure 1 highlights this example, where we see that both the density and pattern of nonzeros increase on lower levels in the hierarchy.

level $\ell$	matrix size $n$	nonzeros <code>nnz</code>	nonzeros per row
			$\text{nnz}/n$
0	1 000 000	26 463 592	26
1	124 984	3 645 644	29
2	23 042	1 466 006	64
3	2991	198 043	66
4	570	32 680	57
5	117	4705	40

Table 1: Matrix properties using classical AMG for a 3D Poisson problem.

54

In parallel, coarse levels that are more dense correlate with an increase in parallel communication costs [4]. Figure 2 shows this by plotting the time spent on each level in an AMG hierarchy during the solve phase. The time grows substantially on coarse levels, which is attributed to increased communication costs from a decrease in sparsity. This effect is common in AMG methods; Figure 2 shows two examples.

In this paper we introduce a method for controlling the communication complexity in AMG. The method increases the sparsity of coarse-grid operators ( $A_\ell$ ,  $\ell = 1, \dots, L$ ) by eliminating entries in  $A_\ell$ . This results in an improved balance between convergence and per-iteration complexity in comparison to the standard algorithm. In addition, we develop an adaptive method which allows nonzero entries to be reintroduced into the AMG hierarchy, thus recovering convergence if entry elimination is too aggressive.

<sup>1</sup>See PyAMG [?] using `ruge_stuben_solver(poisson((100,100,100), type='FE'))` for more details. Similar results are seen using a 7-point finite difference discretization.



Fig. 1: Matrix sparsity pattern using classical AMG for three levels in the hierarchy:  $\ell = 3, 4, 5$ . The full matrix properties are given in Table 1.

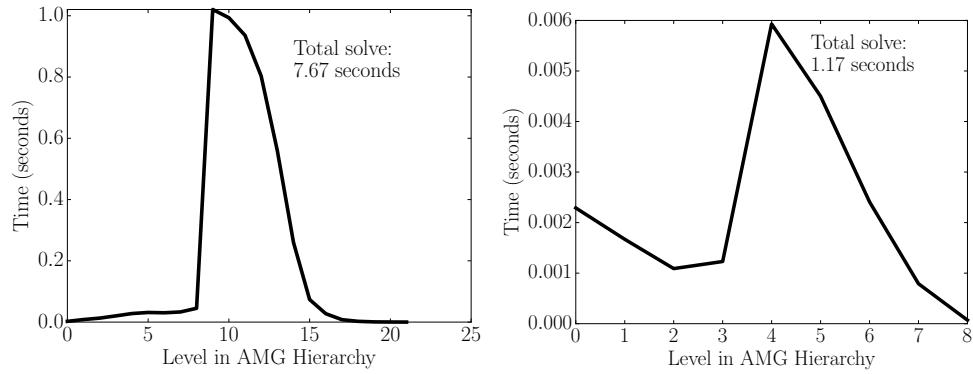


Fig. 2: Left: Time spent on each level of the hierarchy during a single iteration of classical parallel AMG for a 3D Poisson problem with *hypre* using Falgout coarsening [8] and hybrid symmetric Gauss-Seidel relaxation. Right: Repeat experiment, but using aggressive HMIS coarsening. The total time is much lower; however, the qualitative feature of expensive coarse levels remains.

66 In the context of this paper, we define *sparsity* and *density* in terms of the average  
 67 number of nonzeros per row (or equivalently, the average degree of a node in the  
 68 graph of the matrix). In particular, density of a matrix  $A_\ell$  of size  $n_\ell$  is defined  
 69 to be  $\text{nnz}(A_\ell)/n_\ell$ . The performance of AMG is closely correlated with this metric,  
 70 especially communication costs. In addition, note that if a matrix  $A_\ell$  is “sparser” or  
 71 “denser” under this definition, it is also the case under the more traditional density  
 72 metric,  $\text{nnz}(A_\ell)/n_\ell^2$ . Another advantage is that this measure yields a meaningful  
 73 comparison between matrices of different sizes. For example, a goal of our algorithm  
 74 is to generate coarse matrices as close as possible in terms of sparsity structure to the  
 75 fine grid matrix.

76 There are a number of existing approaches to reduce per-iteration communication  
 77 complexity at the cost of convergence. Aggressive coarsening, such as HMIS [20] and  
 78 PMIS [20], rapidly coarsens each level of the hierarchy, leading to a reduction in both

79 the number and the density of coarse operators. While these coarsening strategies  
 80 reduce the cost of each iteration or cycle in the AMG solve phase, they do so at  
 81 the cost of accuracy, often resulting in reduced convergence. Likewise, interpolation  
 82 schemes such as *distance-two interpolation* [19], improve the convergence for aggressive  
 83 coarsening, but also result in an increase in complexity.

84 Figure 2 shows that the time per level during the solve phase is reduced in comparison  
 85 to standard coarsening, even though the same number of processes and problem  
 86 size per core are used. The use of HMIS coarsening is the only difference in problem  
 87 settings between these two runs. Regardless, while aggressive coarsening may reduce  
 88 the total work required during an iteration of AMG, the problem of expensive coarse  
 89 levels still persists.

90 Another strategy for reducing communication complexity in AMG consists of sys-  
 91 tematically improving sparsity in the interpolation operators [19]. Removing nonzeros  
 92 from the interpolation operators reduces the complexity of the coarse-grid operators,  
 93 however this process can also have an unpredictable impact on coarse-level perfor-  
 94 mance if used too aggressively. Sparsity can alternatively be improved by considering  
 95 a sparse approximation to the transformation that relates the fine-grid matrix  $A$   
 96 purely injected to the coarse grid with the less sparse but generally more desirable  
 97 Galerkin coarse-grid matrix. This transformation along with the injected matrix  $A$  go  
 98 on to form a coarse grid that is sparser than Galerkin but still yields good multigrid  
 99 convergence for the several cases [6].

100 The typical approach to building coarse-grid operators,  $A_\ell$ , is to form the Galerkin  
 101 product with the interpolation operator:  $A_{\ell+1} = P_\ell^T A_\ell P_\ell$ . This ensures a *projection*  
 102 in the coarse-grid correction process and a guarantee on the reduction in error in each  
 103 iteration of the AMG solve phase (although the factor by which the error is reduced  
 104 may be small for a poorly converging method). Yet it is the triple matrix product in  
 105 the Galerkin construction that leads to a growth in the number of nonzeros in coarse-  
 106 grid matrices. As such, there are several approaches to constructing coarse operators  
 107 that do not use a Galerkin product and are termed *non-Galerkin* methods. These  
 108 methods have been formed in a classical AMG setting [12] and also in a smoothed  
 109 aggregation [21] context. In general, these methods selectively remove entries from  
 110 coarse-grid operators, reducing the complexity of the multigrid cycle. Assuming the  
 111 appropriate entries are removed from coarse-grid operators, the result is a reduction  
 112 in complexity with little impact on convergence.

113 An alternative to limiting communication complexity is to directly determine the  
 114 coarse-grid stencil, an approach used in geometric multigrid. For instance, simply  
 115 rediscritizing the PDE on a coarse-level results in the same stencil pattern as for the  
 116 original finest-grid operator, thus avoiding any increase in the number of nonzeros in  
 117 coarse-grid matrices. More sophisticated approaches combine geometric and algebraic  
 118 information and include BoxMG [9, 10] and PFMG [3], where a stencil-based coarse-  
 119 grid operator is built. Additionally, collocation coarse grids (CCA) [22] have been used  
 120 on coarse levels to effectively limit the number of nonzeros. Yet, all these methods rely  
 121 on geometric properties of the problem being solved. One exception is the extension of  
 122 collocation coarse grids to algebraic multigrid (ACCA) [11], which has shown similar  
 123 performance to smoothed aggregation AMG.

124 Another approach is to consider sparsification of the graph of the sparse matrix.  
 125 There are several major approaches, grouped by how the difference between the orig-  
 126 inal graph and the sparsified graph is measured. One major approach finds graphs  
 127 where the weight of cuts in the original graph and the sparsified graph are close [?].  
 128 In another, called spectral sparsification [?, ?], edges of the graph are removed if the

129 resulting graph Laplacian remains spectrally close to the original. Sparsification of  
 130 graphs has been an active area of research recently [?], concentrating on developing  
 131 *near* linear-time algorithms for the sparsification. While the immediate impact on the  
 132 coarse-level matrices in a multigrid hierarchy has not been studied, this could point  
 133 to an additional improvement to the methods presented in this paper.

134 The approach developed in this paper is to form a coarse-level operator that does  
 135 not satisfy a Galerkin product by modifying *existing* hierarchies. The novel benefit of  
 136 the proposed approach is that it is applicable to most AMG methods, requires no ge-  
 137 ometric information, and provides a mechanism for recovery if the dropping heuristic  
 138 is chosen too aggressively (see Section 6). This paper is outlined as follows. Section 2  
 139 describes standard algebraic multigrid as well as the method introduced in [12]. Sec-  
 140 tion 3 introduces two new methods for reducing the communication complexity of  
 141 AMG: Sparse Galerkin and Hybrid Galerkin. Parallel performance models for these  
 142 methods are described in Section 4, and the parallel results are displayed in Section 5.  
 143 An adaptive method for controlling the trade-off between communication complex-  
 144 ity and convergence is described in Section 6. Finally, Section 7 makes concluding  
 145 remarks.

146 **2. Algebraic Multigrid.** In this section we detail the AMG setup and solve  
 147 phases. We let the *fine-grid* operator  $A$  be denoted with a subscript as  $A_0$ .

148 Algorithm 1 describes the setup phase and begins with **strength**, which iden-  
 149 tifies the strongly connected edges<sup>2</sup> in the graph of  $A_\ell$  to construct the strength-of-  
 150 connection matrix  $S_\ell$ . From this,  $P_\ell$  is built in **interpolation** to interpolate vectors  
 151 from level  $\ell + 1$  to level  $\ell$ , with the goal to accurately interpolate (smooth) error not  
 152 sufficiently reduced by relaxation. For classical AMG, **interpolation** first forms a  
 153 disjoint splitting of the fine-level index set  $\{1, \dots, n\} = C \cup F$ , where  $C$  is the set of  
 154 indices on the coarse level and where  $F$  is the set of indices for variables that reside  
 155 only on the fine level. The size of the coarse grid is then given by  $n_{\ell+1} = |C|$ , and an  
 156 interpolation operator,  $P_\ell : \mathbb{R}^{n_{\ell+1}} \rightarrow \mathbb{R}^{n_\ell}$ , is constructed using  $S_\ell$  and  $A_\ell$  to compute  
 157 sparse interpolation formulas that are accurate for algebraically smooth functions.  
 158 Finally, the coarse-grid operator is created through a Galerkin triple-matrix product,  
 159  $A_{\ell+1} = P_\ell^T A_\ell P_\ell$ . In a two-level setting, this ensures the desirable property that the  
 160 coarse-grid correction process  $I - P_\ell A_{\ell+1}^{-1} P_\ell^T A_\ell$  is an  $A_\ell$ -orthogonal projection that  
 161 ensures the best coarse-grid correction of the error in the range of interpolation (in the  
 162  $A_\ell$  norm). When an approximation to  $A_{\ell+1}$  is introduced, this projection property is  
 163 lost, leading to a possible deterioration in convergence.

164 The density of each coarse-grid operator  $A_{\ell+1}$  depends on that of the interpolation  
 165 operator  $P_\ell$ . Even interpolation operators with modest numbers of nonzeros typically  
 166 lead to increasingly dense coarse-grid operators [12, 13]. Algorithm 1 addresses this  
 167 with the optional step **sparsify**, which triggers the sparsification steps developed  
 168 in this paper. The approach [12] also fits within this framework, which we detail in  
 169 Section 2.1.

170 The solve phase of AMG, described in Algorithm 2 as a V-cycle, iteratively im-  
 171 proves an initial guess  $x_0$  through the use of the residual equation  $A_0 e_0 = r_0$ , where  
 172  $e_0$  and  $r_0$  are the fine-grid error and residual, respectively. High energy error in the  
 173 approximate solution is reduced through relaxation in **relax** — e.g. Jacobi or Gauss-

---

<sup>2</sup>A degree-of-freedom  $i$  is strongly connected to  $j$  if algebraically smooth error (error not effec-  
 tively reduced by relaxation) varies slowly between them. Strength information critically informs  
 AMG how to coarsen [?] and how to interpolate [18]; while a variety of strength measures abound [?],  
 the standard strength measure is sufficient for the problems tested.

**Algorithm 1:** amg\_setup

---

**Input:**  $A_0$ : fine-grid operator  
**max\_size:** threshold for max size of coarsest problem  
**nongalerkin:** (optional) non-Galerkin method  
 $\gamma_1, \gamma_2, \dots$  (optional) drop tolerances for each level

**Output:**  $A_1, \dots, A_L$ ,  
 $P_0, \dots, P_{L-1}$

---

```

while size( $A_\ell$ ) > max_size
   $S_\ell = \text{strength}(A_\ell)$            {Strength-of-connection of edges}
   $P_\ell = \text{interpolation}(A_\ell, S_\ell)$  {Construct interpolation and injection}
   $A_{\ell+1} = P_\ell^T A_\ell P_\ell$         {Galerkin product}

  if nongalerkin
     $\hat{A}_{\ell+1} = \text{sparsify}(A_{\ell+1}, A_\ell, P_\ell, S_\ell, \gamma_\ell)$  {Remove nonzeros in  $A_{\ell+1}$ }
     $A_{\ell+1} = \hat{A}_{\ell+1}$ 

```

---

174 Seidel. The remaining error is reduced through coarse-grid correction: a combination  
175 of restricting the residual equation to a coarser level, followed by interpolating and  
176 correcting with the resulting approximate error. The coarsest-grid equation is com-  
puted with **solve**, using a direct solution method.

**Algorithm 2:** amg\_solve

---

**Input:**  $x_0$ : fine-level initial guess  
 $b_0$ : fine-level right-hand side  
 $A_0, \dots, A_L$   
 $P_0, \dots, P_{L-1}$

**Output:**  $x_0$ , fine-level approximation

---

```

for  $\ell = 0, \dots, L-1$  do
   $\text{relax}(A_\ell, x_\ell, b_\ell, \nu_1)$            {Pre-smooth  $\nu_1$  times}
   $b_{\ell+1} = P_\ell^T(b_\ell - A_\ell x_\ell)$        {Restrict residual}
   $x_L = \text{solve}(A_L, b_L)$                  {Coarsest-level direct solve}
  for  $\ell = L-1, \dots, 0$  do
     $x_\ell = x_\ell + P_\ell x_{\ell+1}$              {Interpolate and correct}
     $\text{relax}(A_\ell, x_\ell, b_\ell, \nu_2)$            {Post-smooth  $\nu_2$  times}

```

---

177  
178 The dominant computational kernel in Algorithm 2 is the sparse matrix-vector  
179 (SpMV) product, found in **relax** and interpolation/restriction. Typically relaxation  
180 dominates since  $A_\ell$  is larger and denser than  $P_\ell$ . Thus, the performance on level  $\ell$  of  
181 the solve phase depends strongly on the performance of a single SpMV with  $A_\ell$ .

182 When performing parallel sparse matrix operations, a matrix  $A$  is distributed  
183 evenly across processes using a contiguous, row-wise partition, as shown in Figure 3.  
184 The local portion of the matrix is split into two groups: the diagonal block, containing  
185 all columns of  $A$  that correspond to local elements of the vector; and the off-diagonal

block, corresponding to elements of the vector that are stored on other processes. For a SpMV, all off-process elements in the vector that correspond to matrix nonzeros must be communicated. Therefore, the density of a matrix contributes to the cost of communication complexity in the SpMV operation. This implies that the less sparse AMG coarse levels yield large communication costs and, often, an inefficient solve phase [12, 13].

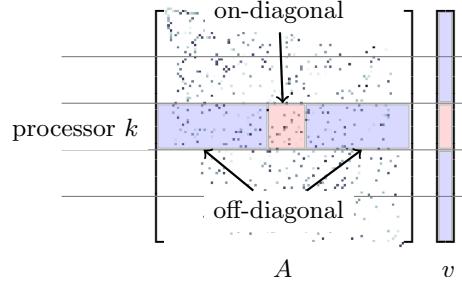


Fig. 3: Matrix  $A$  and vector  $v$  distributed across processes in a row-wise partition.

191

192     **2.1. Method of non-Galerkin coarse grids.** In this section we introduce  
 193 terminology related to [12], which we term the method of *non-Galerkin coarse grids*  
 194 or non-Galerkin for the remainder of the paper. In this case, coarse-level operators  
 195 do not satisfy the Galerkin relationship where  $A_{\ell+1} = P_\ell^T A_\ell P_\ell$  for each level  $\ell$ . The  
 196 coarse-grid operators are first formed through the Galerkin product, followed by a  
 197 sparsification step that generates  $\hat{A}_{\ell+1}$  — see the call to `sparsify` in Algorithm 1. As  
 198 motivated in the previous section, fewer nonzeros in the coarse-grid operator reduce  
 199 the communication requirements. The sparser matrix  $\hat{A}_{\ell+1}$  replaces  $A_{\ell+1}$  and is  
 200 then used when forming the remaining levels of the hierarchy, creating a dependency  
 201 between  $\hat{A}_{\ell+1}$  and all successive levels as shown in Figures 6a and 6b. Thus, this  
 202 approach does not preserve a coarse-grid correction corresponding to an  $A$ -orthogonal  
 203 projection, as described in Section 2.

204     In the following we use `edges`( $A$ ), for a sparse matrix  $A$ , to represent the set  
 205 of edges in the graph of  $A$ . That is, `edges`( $A$ ) =  $\{(i, j) \text{ such that } A_{i,j} \neq 0\}$ , where  
 206  $A_{i,j} = (A)_{i,j}$  is the  $(i, j)$ <sup>th</sup> entry of  $A$ . In addition, we denote  $\hat{P}_\ell$  as the *injection*  
 207 interpolation operator that injects from level  $\ell + 1$  to the  $C$  points on level  $\ell$  so that  
 208  $\hat{P}_\ell$  is defined as the identity over the coarse points, leaving  $\hat{P}_\ell$  zero over the  $F$ -points.

209     The `sparsify` method for reducing the nonzeros in a matrix is described in Al-  
 210 gorithm 3, where the level subscripts are dropped for readability. The algorithm  
 211 selectively removes small entries outside a minimal sparsity pattern<sup>3</sup> given by  $\mathcal{M}_\ell$   
 212 where  $\text{edges}(\mathcal{M}_\ell) = \text{edges}(\hat{P}_\ell^T A_\ell P_\ell + P_\ell^T A_\ell \hat{P}_\ell)$ . For a given tolerance  $\gamma$ , any entry  
 213  $A_{i,j}$  with  $(i, j) \notin \mathcal{M}$  and  $|A_{i,j}| < \gamma \max_{k \neq i} |A_{i,k}|$  is considered insignificant and is re-  
 214 moved. When entry  $A_{i,j}$  is removed, the value of  $A_{i,j}$  is lumped to other entries that  
 215 are strongly connected to  $A_{i,j}$ , and  $A_{i,j}$  is set to zero. This reduces the per-iteration  
 216 communication complexity and heuristically targets spectral equivalence between the  
 217 sparsified operator and the Galerkin operator [12, 21].

<sup>3</sup>The goal of the minimal sparsity pattern is to maintain, at the minimum, a stencil as wide for the coarse grid as exists for the fine grid. This is a critical heuristic for achieving spectral equivalence between the sparsified operator and the Galerkin operator. The current  $\mathcal{M}$  achieves this in many cases. It is possible in some cases to reduce  $\mathcal{M}$  further. See [12] for more details.

**Algorithm 3:** sparsify from [12]**Input:**  $A_c$ : coarse-grid operator $A$ : fine-grid operator $P$ : interpolation $\hat{P}$ : injection $S$ : classical strength matrix $\gamma$ : sparse dropping threshold parameter**Output:**  $\hat{A}_c$ , a sparsified  $A_c$ 

$$\begin{array}{ll} \mathcal{M} = \text{edges}(\hat{P}^T A P + P^T A \hat{P}) & \{\text{Edges in the minimal sparsity pattern}\} \\ \mathcal{N} = \emptyset & \{\text{Edges to keep in } A_c\} \\ \hat{A}_c = \mathbf{0} & \{\text{Initialize sparsified } A_c\} \end{array}$$

```
for  $(A_c)_{i,j} \neq 0$  do
  if  $(i,j) \in \mathcal{M}$  or  $|(A_c)_{i,j}| \geq \gamma \max_{k \neq i} |(A_c)_{i,k}|$ 
     $\mathcal{N} \leftarrow \mathcal{N} \cup \{(i,j), (j,i)\}$  {Add strong edges or the required pattern}
```

```
for  $(A_c)_{i,j} \neq 0$  do
  if  $(i,j) \in \mathcal{N}$ 
     $(\hat{A}_c)_{i,j} = (A_c)_{i,j}$ 
  else
     $\mathcal{W} = \{k \mid S_{j,k} \neq 0, (i,k) \in \mathcal{N}\}$  {Find strong neighbors in the keep list}
    for  $k \in \mathcal{W}$  do
       $\alpha = \frac{|S_{j,k}|}{\sum_{m \in \mathcal{W}} |S_{j,m}|}$  {Relative strength to } k
       $(\hat{A}_c)_{i,k} \leftarrow (\hat{A}_c)_{i,k} + \alpha (A_c)_{i,j}$ 
       $(\hat{A}_c)_{k,i} \leftarrow (\hat{A}_c)_{k,i} + \alpha (A_c)_{i,j}$ 
       $(\hat{A}_c)_{k,k} \leftarrow (\hat{A}_c)_{k,k} - \alpha (A_c)_{i,j}$ 
```

**Algorithm 3b:** Diagonal Lumping – Alternative **for** loop (§ 3.1)

```
for  $(A_c)_{i,j} \neq 0$  do
  1  $\text{ismax} \leftarrow |(A_c)_{i,j}| = \max_{k \neq i} |(A_c)_{ik}|$  and  $(i,k) \notin \mathcal{N} \forall k \neq i$  and  $\sum_j A_{i,j} = 0$ 
  if  $(i,j) \in \mathcal{N}$  or  $\text{ismax}$  {Keep if entry is the single, maximum nonzero}
  2  $(\hat{A}_c)_{i,j} = (A_c)_{i,j}$ 
  else  $(\hat{A}_c)_{i,i} \leftarrow (\hat{A}_c)_{i,i} + (A_c)_{i,j}$  {Otherwise add to the diagonal}
```

218 There is a trade-off between the communication requirements and the conver-  
 219 gence rate. Each entry in the matrix has a communication cost that is dependent  
 220 on the number of network links that the corresponding message travels in addition  
 221 to network contention. In addition, each entry in the matrix also influences conver-

gence of AMG, with large entries generally having larger impact (although this is not uniformly the case). Any entry that has an associated communication cost outweighing the impact on convergence should be removed. However, while it is possible to predict this communication cost based on network topology and message size, the entry's contribution to convergence cannot be easily predetermined. When dropping via non-Galerkin coarse grids, if the chosen drop tolerance is too large, too many entries are removed and convergence deteriorates. Because the ideal drop tolerance is problem dependent and cannot be predetermined, it is likely that the chosen drop tolerance is suboptimal.

Figures 4a and 4b show the convergence and communication complexity, respectively, of various AMG hierarchies for solving a 3D Poisson problem with the method of non-Galerkin coarse grids. For both figures, the 27-point Laplacian was solved on 8192 processes with 10,000 degrees-of-freedom per process. The original Galerkin hierarchy converges in the fewest number of iterations, but has the highest communication complexity. Non-Galerkin removes an ideal number of nonzeros from coarse-grid operators (labeled *ideal*) when no entries are removed from the first coarse level, and all successive levels have a drop tolerance of 1.0. In this case, the communication complexity of the solver is greatly reduced with little effect on convergence. However, if the first coarse level is also created with a drop tolerance of 1.0, essential entries are removed (labeled *too many*). While the complexity of the hierarchy is further reduced, the method fails to converge.

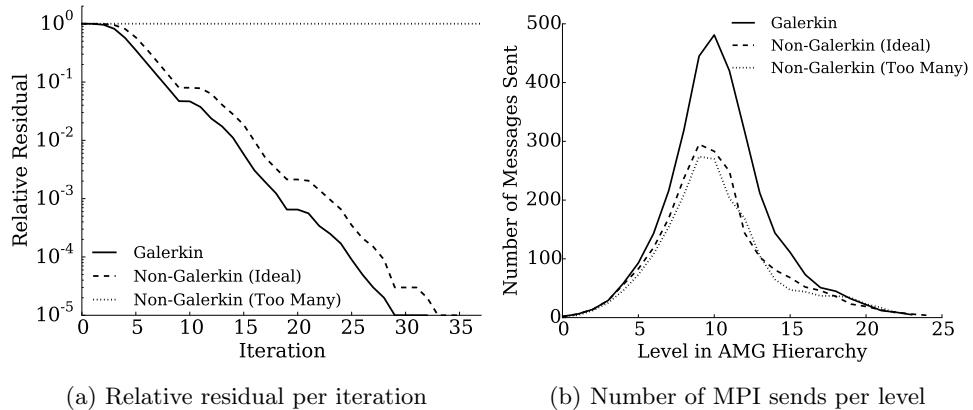


Fig. 4: Sparsity is improved in the AMG hierarchy for the 27-Point Laplacian on 8192 processes with 10,000 degrees-of-freedom per core. The time spent on middle levels of the AMG hierarchy is decreased (right) with little change to the residual after each iteration (left).

If a large drop tolerance is chosen in the non-Galerkin method, the effect on convergence can be determined after one or two iterations of the solve phase. At this point, if convergence is poor, eliminated entries can be re-introduced into the matrix. However, with this method, convergence improvements cannot be guaranteed. As shown in Algorithm 1, sparsifying on a level affects all coarser-grid operators. Hence, adding entries back into the original operator does not influence the impact of their removal on all coarser levels. Figure 5 shows how re-adding entries is inef-

250 effective by plotting the required communication costs versus the achieved convergence  
 251 for both Galerkin and non-Galerkin AMG solve phases for the same 3D Poisson problem.  
 252 The data set *Non-Galerkin (added back)* is generated by removing entries with a  
 253 drop tolerance of 1.0 (everything outside of  $\mathcal{M}$ ) on the first coarse-grid operator and  
 254 0.0 (retaining everything) on all successive levels. This results in a non-convergent  
 255 method. We then add these removed entries back into the first coarse-grid opera-  
 256 tor, but this does not reintroduce the entries which were removed from coarser grid  
 257 operators as a result of the non-Galerkin triple-matrix product  $P_\ell^T \hat{A}_\ell P_\ell$ . Figure 5  
 258 shows that this hierarchy requires little coarse-level communication after all entries  
 259 have been reinstated to the first coarse-grid operator. However as the required entries  
 260 are not added back into all coarser grid operators, the method still fails to converge.

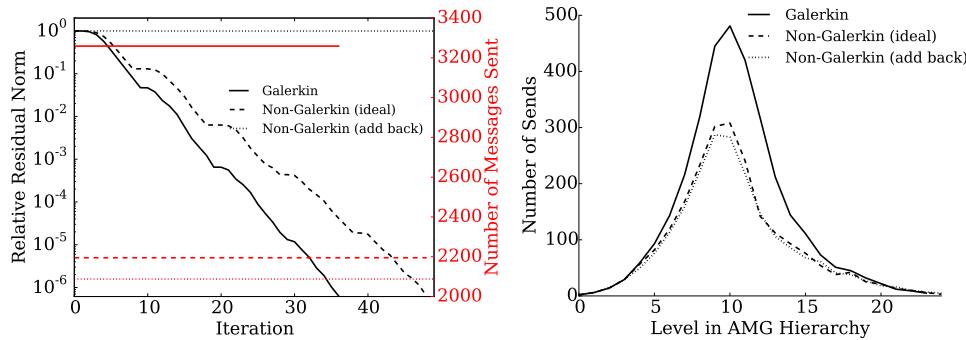


Fig. 5: Convergence vs. communication of Galerkin and non-Galerkin hierarchies for the 27-point Poisson problem on 8192 processes, with 10,000 degrees-of-freedom per process. Relative residual per AMG iteration (black) vs the number of MPI sends per iteration (red) (left), and (maximum) number of sends per level in AMG hierarchy (right)

261

262 **3. Sparse and Hybrid Galerkin approaches.** In this section we present two  
 263 methods as alternatives to the method of non-Galerkin coarse grids. The methods con-  
 264 sist of forming the entire Galerkin hierarchy before sparsifying each operator, yielding  
 265 a lossless approach for increasing sparsity in the AMG hierarchy. The first method,  
 266 which is called the Sparse Galerkin method is described in Algorithm 4 (see Line 1).  
 267 Sparse Galerkin creates the entire Galerkin hierarchy as usual. The hierarchy is then  
 268 thinned as a post-processing step to remove relatively small entries outside of the  
 269 minimal sparsity pattern  $\mathcal{M} = \hat{P}^T A P + P^T A \hat{P}$  using `sparsify`.

270 The second method that we introduce is called Hybrid Galerkin since it combines  
 271 elements of Galerkin and Sparse Galerkin to create the final hierarchy. The method  
 272 is again lossless, and is outlined in Algorithm 4 (see Line 2). After the Galerkin  
 273 hierarchy is formed, small entries outside are removed, this time using a modified,  
 274 minimal sparsity pattern of  $\mathcal{M} = \hat{P}^T \hat{A} P + P^T \hat{A} \hat{P}$ .

275 The Sparse and Hybrid Galerkin methods retain the structure of the original  
 276 Galerkin hierarchy. Consequently, these methods introduce error only into relax-  
 277 ation and residual calculations. The remaining components of each V-cycle in the  
 278 solve phase (see `amg_solve`), such as restriction and interpolation are left unmodified.  
 279 Therefore, the grid transfer operators do not depend on any sparsification, as shown

**Algorithm 4: sparse\_hybrid\_setup**


---

**Input:**  $A_0$ : fine-grid operator  
**max\_size:** threshold for max size of coarsest problem  
 $\gamma_1, \gamma_2, \dots$  drop tolerances for each level  
**sparse\_galerkin:** Sparse Galerkin method  
**hybrid\_galerkin:** Hybrid Galerkin method

**Output:**  $\hat{A}_1, \dots, \hat{A}_L$

```

 $A_1, \dots, A_L, P_0, \dots, P_{L-1} = \text{amg\_setup}(A_0, \text{max\_size}, \text{False})$ 
 $\hat{A}_0 = A_0$ 
for  $\ell \leftarrow 1$  to  $L$  do
  if sparse_galerkin
     $\hat{A}_{\ell+1} = \text{sparsify}(A_{\ell+1}, A_\ell, P_\ell, S_\ell, \gamma_\ell)$  {Increase using the Sparse Method}
  else if hybrid_galerkin
     $\hat{A}_{\ell+1} = \text{sparsify}(A_{\ell+1}, \hat{A}_\ell, P_\ell, S_\ell, \gamma_\ell)$  {Increase using the Hybrid Method}

```

---

280 in Figure 6. Here, we see that the Sparse Galerkin method does not use the modified  
 281 (or sparsified) operators to create the next coarse-grid operator in the hierarchy. Con-  
 282 versely, Hybrid Galerkin uses the newly modified operator to compute the sparsity  
 pattern  $\mathcal{M}$  for the next coarse-grid operator.

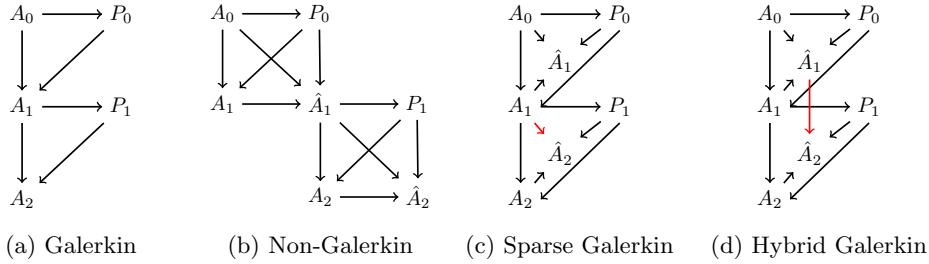


Fig. 6: Dependencies for forming each operator in the various AMG hierarchies. The difference between Sparse and Hybrid Galerkin dependencies is highlighted in red.

283  
 284 The new Sparse Galerkin and Hybrid Galerkin methods reduce the per-iteration  
 285 cost in the AMG solve cycle as less communication is required by each sparse, coarse-  
 286 grid operator. However, high-energy error may also be relaxed at a slower rate,  
 287 yielding a reduction in the convergence factor. As a result, the solve phase is more  
 288 efficient when the reduction in communication outweighs the change in convergence  
 289 factor.

290 Similar to the method of non-Galerkin, it is difficult to predict the impact of  
 291 removing entries from  $A_c$  in Algorithm 3 on the relaxation process. However, as the  
 292 structure of the Galerkin hierarchy is retained, the convergence factor of the solve  
 293 phase can be controlled on-the-fly. In our approach, differences between  $A_c$  and  $\hat{A}_c$   
 294 are stored while forming the sparse approximations. Subsequently, if the conver-  
 295 gence factor falls below a tolerance, entries can be reintroduced into the hierarchy,

296 allowing improvement of the convergence factor up to that of the original Galerkin  
 297 hierarchy (see Section 6).

298 **3.1. Diagonal Lumping.** A significant amount of work is required in Algo-  
 299 rithm 3 to improve the sparsity of each coarse operator. When forming non-Galerkin  
 300 coarse grids, this additional setup cost is hidden by the reduced cost of the triple  
 301 matrix product with the sparsified matrix  $P^T \hat{A} P$ . However, as the entire Galerkin  
 302 hierarchy is initially formed as usual in our new methods (Algorithm 4) the additional  
 303 work greatly reduces the scalability of the setup phase, as shown in Section 5.2. This  
 304 significant cost suggests using an alternative method for sparsification of coarse-grid  
 305 operators. When reducing the number of nonzeros from coarse-grid operators with  
 306 Sparse Galerkin or with Hybrid Galerkin, the structure of the Galerkin hierarchy  
 307 remains intact, allowing a more flexible treatment of increasing sparsity in the ma-  
 308 trix. For instance, one option is to remove entries by lumping to the diagonal rather  
 309 than strong neighbors, as described in Algorithm 3b. This variation of **sparsify** is  
 310 beneficial for several reasons, including: a much cheaper setup phase when compared  
 311 to Algorithm 3; potential to reduce the cost of the solve phase; reduced storage con-  
 312 straints for adaptive solve phases (see Algorithm 5); and retaining positive-definiteness  
 313 of coarse operators.

314 Algorithm 3b replaces the **for** loop in Algorithm 3. For each nonzero entry in the  
 315 matrix, the algorithm first checks if the entry is the maximum element in the row and  
 316 if all other entries in the row are selected for removal (see Line 1). In this case, the  
 317 nonzero entry is not removed if there is a zero row sum.

318 The method of diagonal lumping (Algorithm 3b) results in a cheaper setup phase  
 319 than Algorithm 3. The original non-Galerkin **sparsify** requires each removed en-  
 320 try to be symmetrically lumped to significant neighbors. As a result, the process  
 321 of calculating the associated strong connections requires a large amount of compu-  
 322 tation through a costly loop over neighbors of neighbors. Furthermore, to maintain  
 323 symmetry, all matrix entries that are not stored locally must be updated, requiring  
 324 a significant amount of interprocessor communication. Lumping these entries to the  
 325 diagonal eliminates both the computational and communication complexities.

326 Eliminating the requirement of lumping to strong neighbors yields potential for  
 327 removing a larger number of entries from the hierarchy, further reducing the commu-  
 328 nication costs of the solve phase. The original version of Algorithm 3 requires that  
 329 an entry must have strong neighbors to be removed, as its value is lumped to these  
 330 neighbors.

331 While relaxing the restrictions of the original non-Galerkin **sparsify** provides  
 332 more opportunity to remove entries, the diagonal lumping also negatively influences  
 333 convergence in some cases. Indeed, the energy of the operator can be increased as a  
 334 result of the diagonal lumping, leading to a decrease in (spectral) equivalence with  
 335 the original operator. To mitigate this, if convergence suffers, entries can be easily  
 336 reintroduced into the hierarchy, improving convergence during the solve phase.<sup>4</sup> As  
 337 removed entries are only added to the diagonal, the storage of both the sparse matrix  
 338 along with removed entries is minimal. In addition, these entries can be restored by  
 339 inserting their values to the original positions, and subtracting these values from the

---

<sup>4</sup>Another mitigating factor occurs when diagonal lumping is done after the construction of the hierarchy, when  $\hat{A}_\ell$  will be used only for relaxation. At this point,  $\hat{A}_\ell$  must only effectively damp high energy modes and leave low energy modes largely untouched. Since diagonal lumping of relatively small entries (as controlled by  $\gamma$ ) largely impacts spectral equivalence for low energy modes, this mitigates any effect from reduced spectral equivalence.

340 associated diagonal entries as shown in Algorithm 5. The process of reintroducing  
 341 these entries requires no interprocessor communication as well as a low amount of  
 342 local computational work.

343 Diagonal lumping also preserves matrix properties such as symmetric positive-  
 344 definiteness (SPD). As described in the following theorem, if the sparsity of a diag-  
 345 onally dominant SPD matrix is increased using diagonal lumping, the resulting matrix  
 346 remains SPD. Consequently, Sparse and Hybrid Galerkin with diagonal lumping can  
 347 be used in preconditioning many methods such as conjugate gradient. It is important  
 348 to note, that while SPD matrices are an attractive property for AMG, AMG meth-  
 349 ods do not guarantee diagonal dominance of the coarse-grid operators. Yet, in many  
 350 instances this property is preserved, for example for more standard elliptic operators.  
 351

352 **THEOREM 1.** *Let  $A$  be SPD and diagonally dominant. If  $\hat{A}$  is produced by Algo-*  
 353 *rithm 3b, then it is symmetric positive semi-definite and diagonally dominant.*

354 *Proof.* Let  $A$  be SPD with diagonal dominance,

355 (2) 
$$|A_{i,i}| \geq \sum_{k \neq i} |A_{i,k}|, \forall i.$$

356 Symmetry of  $\hat{A}$  is guaranteed from the symmetry of both  $A$  and the  $\mathcal{N}$  from Algo-  
 357 rithm 3. For all off-diagonal entries  $(i,j), (j,i) \in \mathcal{N}$ ,

358 (3) 
$$\hat{A}_{i,j} = A_{i,j} = A_{j,i} = \hat{A}_{j,i},$$

359 by Line 2 in Algorithm 3b and the symmetry of  $A$ .

360 The positive-definiteness is guaranteed by the diagonal dominance and a Gersh-  
 361 gorin disc argument. The proof proceeds by starting with the matrix  $A$  and then  
 362 considering the change made to  $A$  by the elimination of each entry. Initially, all the  
 363 Gershgorin discs of  $A$  are strictly on the right-side of the origin, thus implying that all  
 364 eigenvalues are non-negative. Then, assume that we eliminate some arbitrary entry  
 365  $A_{i,j}$ ,  $(i,j) \in \mathcal{N}$ . This results in row  $i$  being updated

366 (4) 
$$A_{i,i} \leftarrow A_{i,i} + A_{i,j} \quad \text{and} \quad A_{i,j} \leftarrow 0$$

367 If  $A_{i,j} > 0$ , then the center of the Gershgorin disc is shifted to the right, and the radius  
 368 shrinks, thus keeping the disc to the right of the origin and preserving definiteness.  
 369 If  $A_{i,j} < 0$ , then the center of the disc is shifted to the left by  $|A_{i,j}|$ , but the radius  
 370 of the disc also shrinks by  $|A_{i,j}|$ . This also keeps the disc to the right of the origin  
 371 and preserves semi-definiteness. Furthermore, since each disc is never shifted to the  
 372 left half plane, diagonal dominance is also preserved. The proof then proceeds by  
 373 considering all of the entries to be eliminated.  $\square$

374 **REMARK 3.1.** *If any row of  $A$  is strictly diagonally dominant, as often happens*  
 375 *with Dirichlet boundary conditions, then  $\hat{A}$  will be positive definite. Essentially, Al-*  
 376 *gorithm 3b never shifts a Gershgorin disc to the left, so  $\hat{A}$  can have no 0 eigenvalue.*

377 **4. Parallel Performance.** In this section we use a parallel performance model  
 378 to illustrate the per-level costs associated with each of the six methods:

379 **Galerkin** - Classic coarsening in AMG, as outlined in Algorithm 1;

380 **Non-Galerkin** - The base algorithm presented in [12];

381 **Sparse Galerkin** - The new Algorithm 4 with `sparse_galerkin` and full lumping  
 382 from Algorithm 3;

- 383 **Sparse Galerkin (Diag)** - The new Algorithm 4 with `sparse_galerkin` and diagonal lumping from Algorithm 3b;  
 384 **Hybrid Galerkin** - The new Algorithm 4 with `hybrid_galerkin` and full lumping from Algorithm 3; and  
 385 **Hybrid Galerkin (Diag)** - The new Algorithm 4 with `hybrid_galerkin` and diagonal lumping from Algorithm 3b.

389 The solve phase of AMG (see Algorithm 2) is largely comprised of sparse matrix-vector multiplication, thus we model each method by assessing the cost of performing 390 a SpMV on each level of the hierarchy. We focus on the operators  $A_\ell$ , as the work 391 required for this matrix is more costly than the restriction and interpolation operations. Specifically, we employ an  $\alpha$ - $\beta$  model to capture the cost of the parallel SpMV 392 based on the number of nonzeros in  $A$ . We denote  $p$  as the number of processors,  $\alpha$  393 as the latency or startup cost of a message, and  $\beta$  as the reciprocal of the network 394 bandwidth [13, 14]. In addition,  $\text{nnz}_p$  represents the average number of nonzeros local 395 to a process, while  $s_p$  and  $n_p$  are the maximum number of MPI sends and message size 396 across all processors. Finally, we use  $c$  to represent the cost of a single floating-point 397 operation. With this we model the total time as  
 398

399 (5) 
$$T = 2c\text{nnz}_p + \max_p s_p(\alpha + \beta n_p).$$

400 For the model parameters above we use the Blue Waters supercomputer at the University 401 of Illinois at Urbana-Champaign [1, 5]. The latency and bandwidth were measured 402 through the HPCC benchmark [16], yielding  $\alpha = 1.8 \times 10^{-6}$  and  $\beta = 1.8 \times 10^{-9}$ . 403 Since the achieved floprate depends on matrix size, we determine the value of  $c$  by 404 timing the local SpMV. Specifically, letting  $\text{nnz}_{\text{local}}$  be the number of nonzeros local 405 to the processor and  $T_{\text{local}}$  the time to perform the local portion of the SpMV, we 406 compute  $c = T_{\text{local}} / 2\text{nnz}_{\text{local}}$  for each matrix in the hierarchy.

407 The minimal per-level cost associated with the non-Galerkin and Sparse/Hybrid 408 Galerkin methods occurs when entries are removed with a drop tolerance of  $\gamma =$   
 409 1.0. Using the model, (5), this is highlighted in Figure 7 for both the Laplace and 410 rotated anisotropic diffusion problems (a full description of these problems is given 411 in Section 5). We see that both non-Galerkin and Hybrid Galerkin have potential 412 to minimize the per-level cost. However, when the per-level cost is minimized, the 413 convergence of AMG often suffers. Therefore, less-aggressive drop tolerances such as 414  $\gamma < 1.0$  may remove fewer entries, increasing the per-level cost, but due to better 415 convergence will improve the overall cost of the solve phase. Indeed for the rotated 416 anisotropic diffusion problem, this is the case, where we reduce  $\gamma$  at fine levels in the 417 hierarchy in order to retain convergence (not shown for brevity).  
 418

419 **5. Parallel results for Sparse and Hybrid Galerkin.** In this section we 420 highlight the parallel performance of the Sparse and Hybrid Galerkin methods. We 421 consider scaling tests on the familiar 3D Laplacian since this is a common multigrid 422 problem used to establish a baseline. In order to test problems where AMG convergence 423 is suboptimal, we consider a 2D rotated anisotropic diffusion problem. Finally, 424 we test our methods on a suite of matrices from the Florida Sparse Matrix Collection. 425 All computations were performed on the Blue Waters system at the University of 426 Illinois at Urbana-Champaign [1]. Each method was implemented and solved with 427 `hypre` [2, 15], using default parameters unless otherwise specified. In summary, we 428 compare the solve and setup times for the six methods considered in Section 4 while 429 preconditioning a Krylov method such as CG or GMRES in each test.

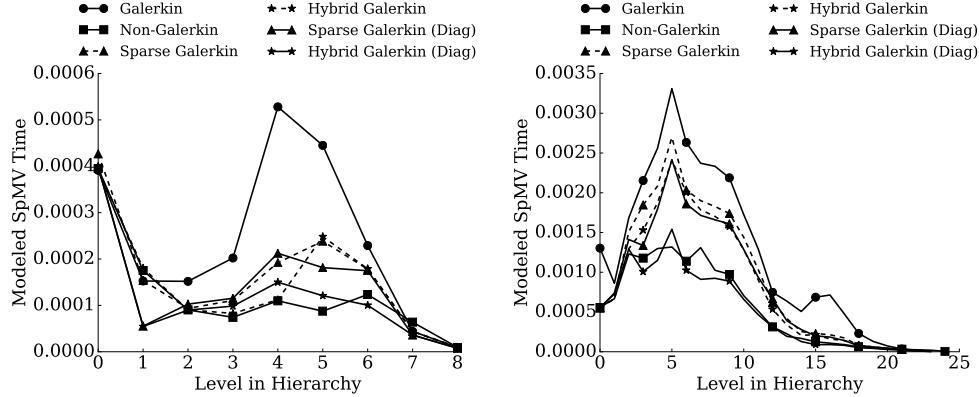


Fig. 7: Modeled minimal cost of a single SpMV on each level of the AMG hierarchy for Laplace (left) and rotated anisotropic diffusion (right), for an aggressive drop tolerance of 1.0 on each level.

430        The drop tolerances for each method vary by level, using a combination of 0.0,  
 431        0.01, 0.1, and 1.0 across the coarse levels. Six combinations of these drop tolerances  
 432        are tested for the various test cases, and the series yielding the minimum solve time  
 433        for each is selected. *Note:* At 100,000 cores, the best drop tolerances from the second  
 434        largest run size are used due to large costs associated with running 6 drop tolerances  
 435        at this core count. Details of the drop tolerances used in all the below tests are found  
 436        in the Supplemental Materials due to length.<sup>5</sup>

437        Generally, the drop tolerances are aggressive for the simple isotropic 3D Laplacian  
 438        example, where for instance in Figures 8 and 9, the Sparse Galerkin (Diag) method  
 439        used the values of [0.0, 0.1], i.e., no dropping on the first coarse level and then 0.1  
 440        on all subsequent levels. For the more complex examples such as rotated anisotropic  
 441        diffusion, the drop tolerances are less aggressive and usually begin on coarser levels.  
 442        For instance in Figures 8 and 9, the Sparse Galerkin (Diag) method used values of  
 443        [0.0, 0.0, 0.0, 0.1], i.e., no dropping occurs until the third coarse level, where 0.1 is used  
 444        from that point onward on all coarser levels. Overall, the drop tolerance is similar to  
 445        other multigrid parameters, such as the strength-of-connection drop tolerance. Some  
 446        experimentation with a problem type is required, but thereafter a general, conservative  
 447        parameter choice can be made for subsequent use.

448        We consider the diffusion problem

$$449 \quad (6) \quad -\nabla \cdot K \nabla u = 0,$$

450        with two particular test cases for our simulations. Also, as problems with less structure  
 451        result in increased density on coarse levels, we consider a subset from the Florida  
 452        sparse matrix collection. The test problems are as follows:

453        **Laplace** - Here, we use  $K = I$  on the unit cube with homogeneous Dirichlet bound-  
 454        ary conditions. Q1 finite elements are used to discretize the problem using a uni-  
 455        form mesh, leading to a familiar 27-point stencil. The preconditioner formed for

<sup>5</sup>In addition, the modifications to *hypre v2.11.0* are stored at <https://github.com/lukeolson/hybre/releases/tag/SISC-sparse-hybrid-galerkin> and <https://github.com/lukeolson/hybre/releases/tag/SISC-sparse-hybrid-galerkin-adaptive>.

456 the 3D Laplacian uses aggressive coarsening (HMIS) and distance-two (extended  
 457 classical modified) interpolation. The interpolation operators were formed with a  
 458 maximum of five elements per row, and hybrid symmetric Gauss-Seidel was the  
 459 relaxation method.

460 **Rotated Anisotropic Diffusion** - In this case, we consider a diffusion tensor with  
 461 homogeneous Dirichlet boundary conditions of the form  $K = Q^T D Q$ , where  $Q$  is  
 462 a rotation matrix and  $D$  is a diagonal scaling defined as

$$463 \quad (7) \quad Q = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix} \quad D = \begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix}.$$

464 Q1 finite elements are used to discretize a uniform, square mesh. In the following  
 465 tests we use  $\theta = \frac{\pi}{8}$  and  $\epsilon = 0.001$ . In each case, the preconditioner uses Falgout  
 466 coarsening [8], extended classical modified interpolation and hybrid symmetric  
 467 Gauss-Seidel.

468 **Florida Sparse Matrix Collection** - We consider a subset of all real, symmetric,  
 469 positive definite matrices from the Florida sparse matrix collection with size over  
 470 1,000,000 degrees-of-freedom. In addition we consider only the cases where GM-  
 471 RES preconditioned with Galerkin AMG converges in fewer than 100 iterations.  
 472 Each problem uses HMIS coarsening and so-called *extended+i* interpolation if  
 473 possible. In some cases, however, Galerkin AMG does not converge with these  
 474 options; in these cases Falgout coarsening and modified classical interpolation  
 475 are used. Relaxation for all systems is hybrid symmetric Gauss-Seidel. *Note:*  
 476 When necessary for convergence, some *hypre* parameters, such as the minimum  
 477 coarse-grid size and strength tolerance, vary from the default.

478 The following results demonstrate that the diagonally lumped Sparse and Hybrid  
 479 Galerkin methods are able to perform comparably to non-Galerkin. Non-Galerkin  
 480 and Sparse/Hybrid Galerkin all significantly reduce the per-iteration cost by reducing  
 481 communication on coarse levels. Since the method of non-Galerkin is multiplicative  
 482 in construction, the setup times are often much lower in comparison to standard  
 483 Galerkin. However, Sparse and Hybrid do not observe this benefit since they are  
 484 constructed in a post-processing step. While the per-iteration work is decreased for all  
 485 methods, the convergence suffers for the case of rotated anisotropic diffusion problems  
 486 with non-Galerkin at large scales. However, Sparse and Hybrid Galerkin converge at  
 487 rates similar to the original Galerkin hierarchy, yielding speedup in total solve times.  
 488 Moreover, in a strong scaling study, we observe that Hybrid Galerkin is competitive,  
 489 particularly at large core counts.

490 **5.1. Improving sparsity in AMG Hierarchies.** The significant number of  
 491 nonzeros on coarse levels creates large, relatively dense matrices near the middle of  
 492 the AMG hierarchy, yielding large communication costs for each SpMV performed on  
 493 these levels. As the solve phase of AMG consists of many SpMVs on each level of the  
 494 hierarchy, the time spent on coarse levels can increase dramatically. Sparse, Hybrid,  
 495 and non-Galerkin can all reduce both the cost associated with communication as well  
 496 as the time spent on each level during a solve phase.

497 Figure 8 shows the time spent on each level of the hierarchy during a single  
 498 iteration of AMG, for both test cases with 10,000 degrees-of-freedom per core using  
 499 8192 cores. Both the method of non-Galerkin coarse grids, as well as the Sparse and  
 500 Hybrid Galerkin methods, reduce the time required on levels near the middle of the  
 501 hierarchy. Non-Galerkin has a larger impact on the time spent on middle levels of the  
 502 hierarchy for the Laplace problem than Sparse and Hybrid Galerkin. However, for the

503 anisotropic problem, diagonally-lumped Hybrid Galerkin reduces level-time similarly  
 504 to non-Galerkin. This is due to a large reduction in the number of messages required  
 505 in each SpMV as shown in Figure 9. The reduction in total size of all messages  
 communicated is relatively small.

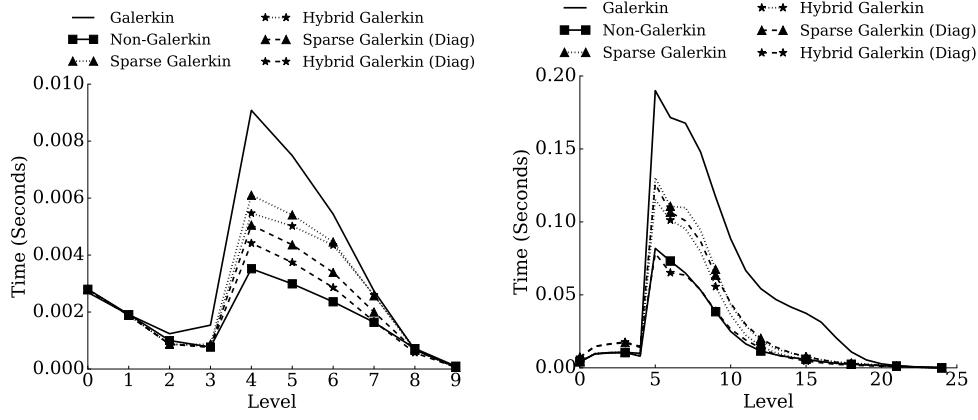


Fig. 8: Time spent on each level of the AMG hierarchy during a single iteration of the solve phase for **Laplace** (left) and **Rotated Anisotropic Diffusion** (right), each with 10,000 degrees-of-freedom per core.

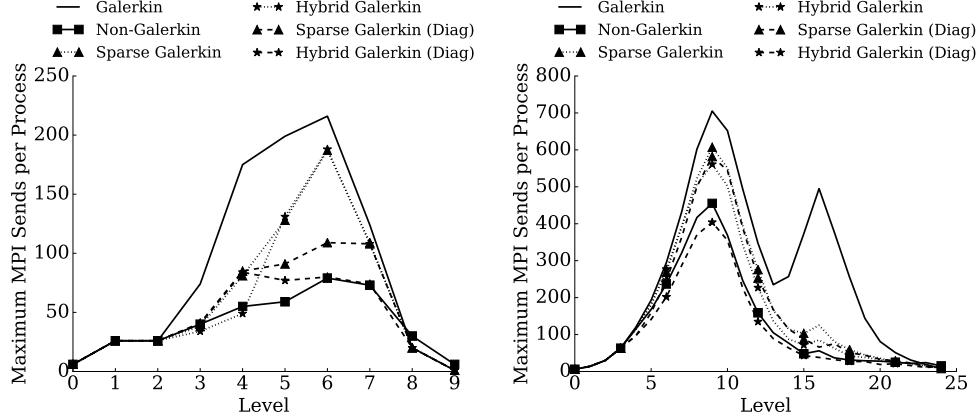


Fig. 9: Number of sends required to perform a single SpMV on each level of the AMG hierarchy for: **Laplace** (left) and **Rotated Anisotropic Diffusion** (right), each with 10,000 degrees-of-freedom per core.

506  
 507     The increase in time spent on each level, as well as the associated communication  
 508 costs of these levels, becomes more pronounced at higher processor counts in a strong  
 509 scaling study. Figure 10 illustrates this by plotting the per-level times required during  
 510 a single iteration of AMG, as well as the number of messages communicated during  
 511 a SpMV for the rotated anisotropic diffusion problem with 1,250 degrees-of-freedom  
 512 per core using 8192 cores. Compared with the 10,000 degrees-of-freedom per core in

513 Figures 8 and 9 there is a sharper increase in time required for levels near the middle  
of the hierarchy due to the increasing dominance of communication complexity.

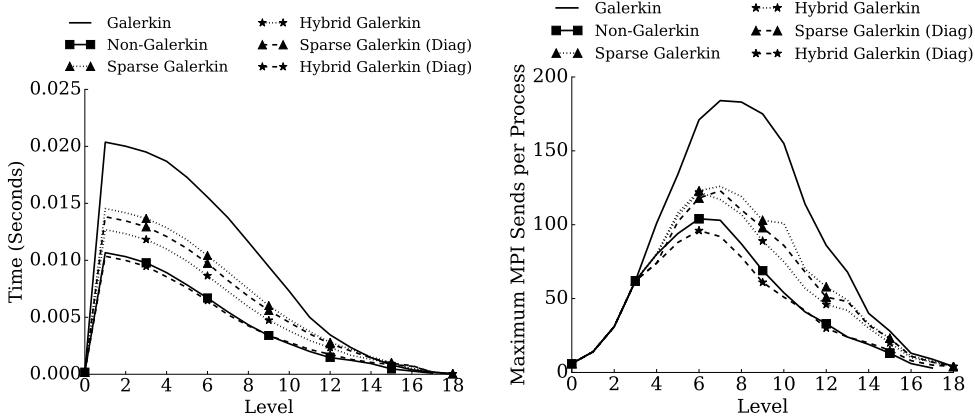


Fig. 10: For each level of the AMG hierarchy, time per iteration of AMG (left) and number of messages sent during a single SpMV (right) for the **Rotated Anisotropic Diffusion** problem with 1,250 degrees-of-freedom per core.

514

515     **5.2. Costs of Weakly Scaled Setup Phases.** Each sparsification method  
516 can lead to reduced communication costs in the middle of the hierarchy. However,  
517 removing insignificant entries from coarse-grid operators requires additional work in  
518 the setup phase. In the non-Galerkin method, setup times are reduced since the  
519 increased sparsity is used directly in the triple-matrix product required to form each  
520 successive coarse-grid operator. However, for the new methods, Sparse and Hybrid  
521 Galerkin, the entire Galerkin hierarchy is first constructed so that the sparsify process  
522 on each level requires additional work. Figure 11 shows the times required to setup  
523 an AMG hierarchy for rotated anisotropic diffusion, with Laplace setup times scaling  
524 in a similar manner. While there is a slight increase in setup cost associated with the  
525 Sparse and Hybrid Galerkin hierarchies, this extra work is nominal. Therefore, while  
526 the majority of this additional work is removed when using diagonal lumping, the  
527 differences in work required in the setup phase between these two lumping strategies  
528 is insignificant for the problems being tested.

529     **5.3. Weak Scaling of GMRES Preconditioned by AMG.** In this section  
530 we investigate the *weak* scaling properties of the methods. Figure 12 shows both the  
531 average convergence factor and total time spent in the solve phase for a weak scaling  
532 study with rotated anisotropic diffusion problems at 10,000 degrees-of-freedom per  
533 core using GMRES preconditioned by AMG. GMRES is used over CG because Algo-  
534 rithm 3 guarantees symmetry but not positive-definiteness of the preconditioner. In  
535 many cases, positive-definiteness is preserved, but when using more aggressive drop  
536 tolerances, we have observed this property being lost. While the convergence of both  
537 diagonally-lumped Sparse and Hybrid Galerkin remain similar to that of Galerkin,  
538 the non-Galerkin method converges more slowly. Therefore, while non-Galerkin and  
539 diagonally-lumped Hybrid Galerkin yield similar communication requirements, GM-  
540 RES preconditioned by Hybrid Galerkin performs significantly better as fewer itera-  
541 tions are required.

REMARK 5.1. With the chosen drop tolerances, non-Galerkin does not converge for this anisotropic problem at 100,000 cores. In this case, nothing was dropped from the first three coarse levels of the hierarchy. On the fourth coarse level a drop tolerance of 0.01 was used, and the fifth was sparsified with a tolerance of 0.1. The remaining levels were sparsified with a drop tolerance of 1.0. This was determined to be the best tested drop tolerance sequence for smaller run sizes, and multiple drop tolerance sequences were not tuned at this large problem size due to the significant costs. However, a better drop tolerance could yield a convergent non-Galerkin method at this scale.

The efficiency of weakly scaling to  $p$  processes is defined as  $E_p = \frac{T_1}{T_p}$ , where  $T_1$  is the time required to solve the problem on a single process and  $T_p$  is the time to solve on  $p$  processes. The efficiency of solving weakly scaled rotated anisotropic diffusion problems with non-Galerkin, Sparse Galerkin, and Hybrid Galerkin, relative to the efficiency of Galerkin AMG, are shown in Figure 13. While both the original and diagonally-lumped Sparse and Hybrid Galerkin methods scale more efficiently than Galerkin, the poor convergence of non-Galerkin on large run sizes yields a reduction in relative efficiency.

**5.4. Strong Scaling of GMRES Preconditioned by AMG.** We next consider the rotated anisotropic diffusion system with approximately 10,240,000 unknowns using cores ranging from 128 to 100,000. Therefore, the simulation is reduced from 80,000 degrees-of-freedom per core when run on 128 cores, to just over 100 degrees-of-freedom per core on 100,000 cores. Computation dominates the total cost of solving a problem partitioned over relatively few processes, as each process has a large amount of local work. However, as the problem is distributed across an increasing number of processes, the local work decreases while communication requirements increase. Therefore, the time required to solve a problem is reduced with strong scaling, but only to the point where communication complexity begins to dominate. The efficiency of strongly scaling to  $p$  processes is defined as  $E_p = \frac{T_1}{pT_p}$ . Figure 14 shows the efficiency of solving a strongly scaled rotated anisotropic diffusion problem with GMRES preconditioned by the various sparse methods. In each case we observe

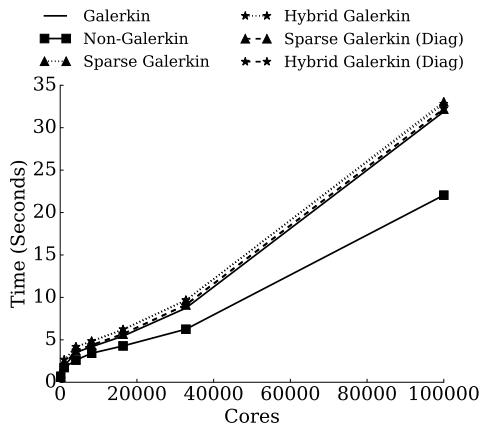


Fig. 11: Time required to setup AMG hierarchy for **Rotated Anisotropic Diffusion** with 10,000 degrees-of-freedom per core.

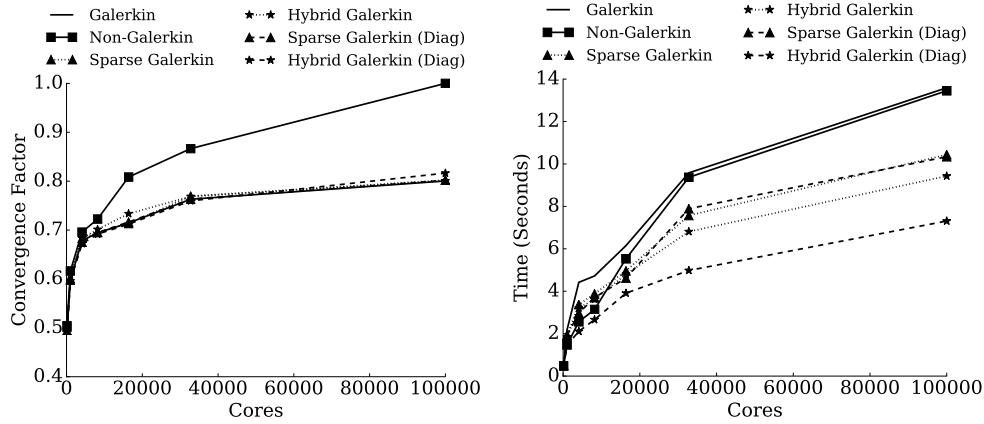


Fig. 12: Convergence factors (left) and times (right) for weak scaling of **Rotated Anisotropic Diffusion** (10,000 degrees-of-freedom per core), solved by preconditioned GMRES. For large problem sizes, non-Galerkin AMG does not converge, and timings indicate when the maximum iteration count was reached.

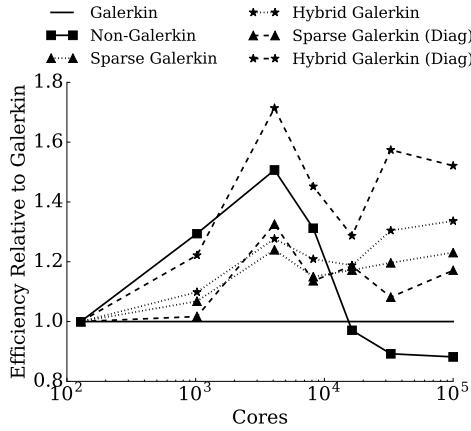


Fig. 13: Efficiency of solving weakly scaled **Rotated Anisotropic Diffusion** at 10,000 degrees-of-freedom per core with various methods, relative to that of the Galerkin hierarchy.

571 improvements over standard Galerkin.

572 A strong scaling study is also performed on the subset of matrices from the  
 573 Florida sparse matrix collection. These problems were tested on 64, 128, 256, and  
 574 512 processes. Figure 15 shows the time required to perform a single V-cycle for each  
 575 of the matrices in the subset, relative to the time required by Galerkin AMG. All  
 576 methods reduce the per-iteration times for each matrix in the subset. Furthermore,  
 577 the total time required to solve each of these matrices is also reduced, as shown in  
 578 Figure 16. While Sparse Galerkin provides some improvement, the Hybrid and non-  
 579 Galerkin methods are comparable, particularly at high core counts.

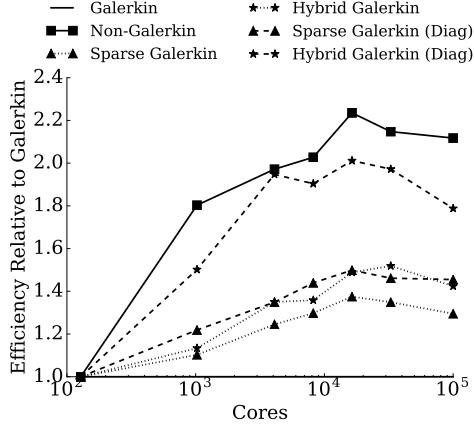


Fig. 14: Efficiency of non-Galerkin and Sparse/Hybrid Galerkin methods in a strong scaling study, **relative to Galerkin AMG** for **Rotated Anisotropic Diffusion**.

580        **5.5. Diagonal Lumping Alternative and PCG.** Diagonal lumping retains  
 581 positive definiteness of diagonally-dominant coarse-grid operators as described in The-  
 582 oreom 1. Therefore, as the preconditioned conjugate gradient (PCG) method requires  
 583 both the matrix and preconditioner to be symmetric and positive-definite, the Laplace  
 584 and anisotropic diffusion problems are solved by conjugate gradient preconditioned by  
 585 the diagonally-lumped Sparse and Hybrid Galerkin hierarchies. Figure 17 shows the  
 586 solve phase times for solving the weakly scaled rotated anisotropic diffusion problem  
 587 with PCG. As with GMRES, both the Sparse and Hybrid Galerkin preconditioners  
 588 decrease the time required in the AMG solve phase during a weak scaling study.

589        **6. Adaptive Solve Phase.** The previous results describe the case where good  
 590 drop tolerances were known a priori for **sparsify**. However, as the appropriate drop  
 591 tolerance changes with problem type, problem size, and even level of the AMG hier-  
 592 archy, a good drop tolerance is often not easily realized. When the drop tolerance is  
 593 too small, few entries are removed from the hierarchy and the communication com-  
 594 plexity remains the same. However, if the drop tolerance is too large, the solver is  
 595 non-convergent, as described in Section 2.1.

596        In this section we consider an *adaptive* method that attempts to add entries back  
 597 into the hierarchy as a deterioration in convergence is observed. This is detailed in  
 598 Algorithm 5. The algorithm initializes a Sparse or Hybrid Galerkin hierarchy and  
 599 proceeds by executing  $k$  iterations of a preconditioned Krylov method — e.g. PCG.  
 600 If the convergence is below a tolerance, the coarse levels are traversed until a coarse  
 601 grid operator is found on which entries were removed with a drop tolerance greater  
 602 than 0.0. Entries are then added back to this coarse-grid operator, reducing the drop  
 603 tolerance by a factor of 10. Any new drop tolerance below  $\gamma_{\min} = 0.01$  is rounded  
 604 down to 0.0. This continues until entries have been reintroduced into  $s$  coarse-grid  
 605 operators. At this point, the Krylov method continues, using the most recent value  
 606 for  $x$  unless the previous iterations diverged from the true solution.

607        This entire process is then repeated until convergence. The adaptive solve phase  
 608 requires additional iterations over Galerkin AMG, as initial iterations of this method  
 609 may not converge. However, the goal of this solver is to guarantee convergence sim-

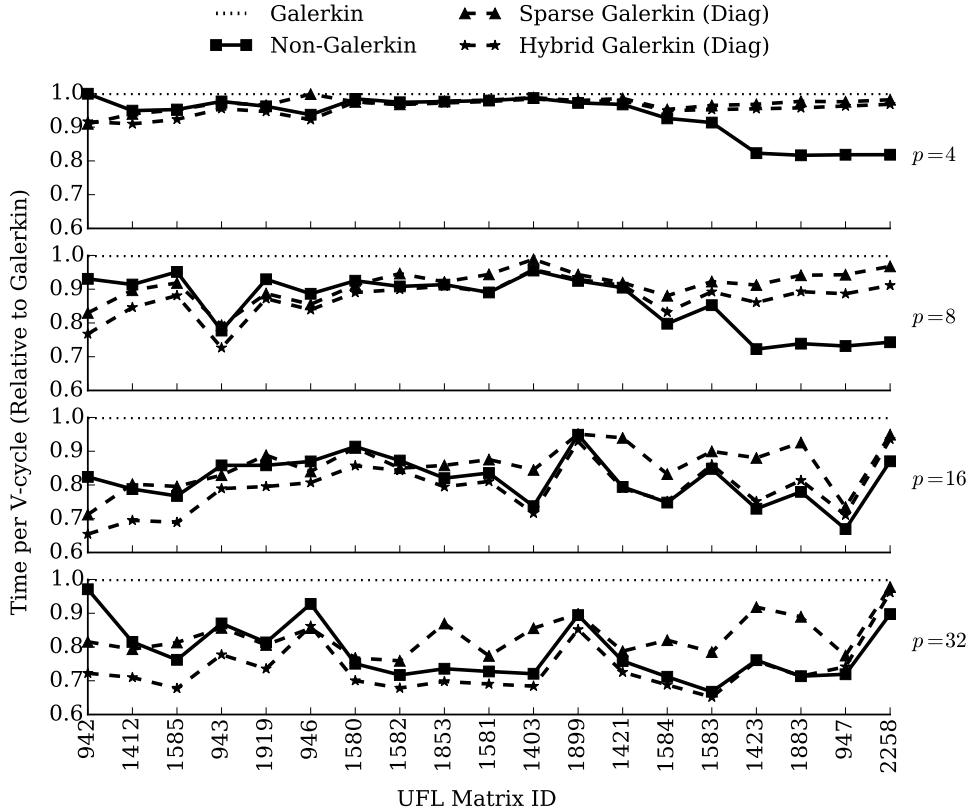


Fig. 15: Time (relative to Galerkin) per iteration for each matrix in the Florida Sparse Matrix Collection, using  $p = 64, 128, 256$ , and  $512$ .

610 ilar to Galerkin AMG. Speed-up over Galerkin AMG is still dependent on choosing  
 611 reasonable initial drop tolerances.

612 If the Krylov method is not flexible, such as PCG or GMRES, then it must be  
 613 restarted after the preconditioner has been edited. On the other hand, a flexible  
 614 method, such as FGMRES, would not have to be restarted, but requires greater  
 615 memory storage (and possibly also more computational work) than PCG. While the  
 616 new algorithms are agnostic to the Krylov scheme, we use restarted PCG in order to  
 617 directly compare schemes.

618 EXAMPLE 6.1. As an example, consider the case of a hierarchy with 6 levels using  
 619 drop tolerances of  $[0, 0.01, 0.1, 1.0, 1.0, 1.0]$  — i.e.,  $\hat{A}_1$  retains all entries from  $A_1$ ,  
 620  $\hat{A}_2$  and  $\hat{A}_3$  result from `sparsify` with  $\gamma = 0.01$  and  $\gamma = 0.1$ , etc. Suppose that  
 621 `adaptive_solve` with  $k = 3$  and  $s = 2$  results in 3 iterations of PCG and a large  
 622 residual. The adaptive solve finds the first level containing a sparsified coarse grid  
 623 matrix, namely  $\hat{A}_2$ . The drop tolerance on this level is changed from 0.01 to 0.0, and  
 624 the original coarse matrix  $A_2$  is sparsified with the new drop tolerance. Furthermore,  
 625 since  $s = 2$  the drop tolerance on level 3 is reduced from 0.1 to 0.01, and  $A_3$  is also  
 626 sparsified. PCG then restarts with the new hierarchy. If convergence continues to

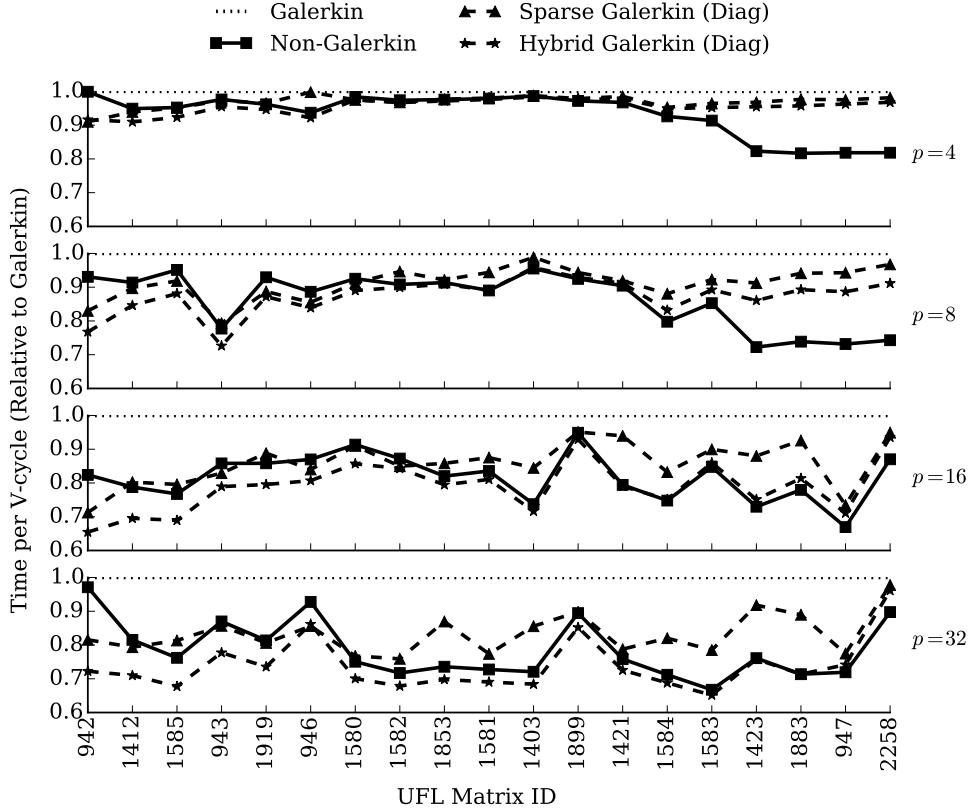


Fig. 16: Time (relative to Galerkin) per AMG solve for each matrix in the Florida Sparse Matrix Collection, using  $p = 64, 128, 256$ , and  $512$ .

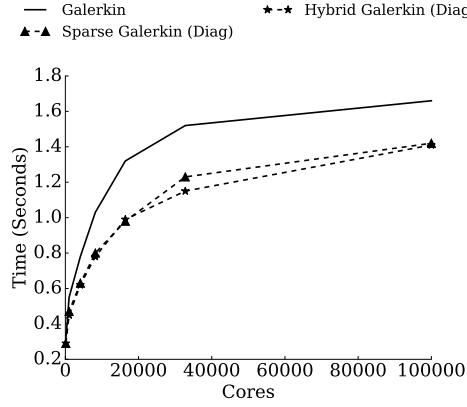


Fig. 17: Weak scaling solve time for **Rotated Anisotropic Diffusion**, solved by PCG preconditioned by various AMG hierarchies.

**Algorithm 5:** adaptive\_solve

---

<b>Input:</b>	$A, b, x_0$	
	$\hat{A}_1, \dots, \hat{A}_L$	Sparse/Hybrid Galerkin coarse grid matrices
	$A_1, \dots, A_L$	original Galerkin coarse grid matrices
	$P_0, \dots, P_{L-1}$	
	$k$	PCG iterations before convergence test
	$s$	AMG levels per update
	$\gamma_0, \dots, \gamma_L$	sparsification drop tolerance used at each level
	$\text{tol}$	convergence tolerance
	$\text{sparse\_galerkin}$	Sparse Galerkin method
	$\text{hybrid\_galerkin}$	Hybrid Galerkin method
<b>Output:</b>	$x$	

---

$x = x_0$   
 $r_0 = b - Ax_0$

```

while  $\|r\|/\|r_0\| \leq \text{tol}$ 
   $M = \text{preconditioner}(\text{amg\_solve}, \hat{A}_1, \dots, \hat{A}_L, P_0, \dots, P_{L-1})$ 
   $x = \text{PCG}(A, b, x, k, M)$  {Call  $k$  steps of preconditioned CG}
   $r = b - Ax$ 
  if  $\frac{\|r\|}{\|r_0\|} \leq \text{tol}$ 
     $\text{continue}$ 
  else
    for  $\ell = 0, \dots, L$  do
      if  $\gamma_\ell > 0$ 
         $\ell_{\text{start}} \leftarrow \ell$  {Find finest level that uses dropping}
    for  $\ell = \ell_{\text{start}} \dots \ell_{\text{start}} + s$  do
       $\gamma_\ell = \begin{cases} \frac{\gamma_\ell}{10}, & \text{if } \frac{\gamma_\ell}{10} > \gamma_{\min} \\ 0, & \text{otherwise} \end{cases}$  {Determine new dropping parameter}
       $\gamma_{\min}$  is the  $\min(\gamma_0 \dots \gamma_L)$ 
      if  $\text{sparse\_galerkin}$  {Re-add entries at the new dropping tolerance}
         $\hat{A}_\ell = \text{sparsify}(A_\ell, A_{\ell-1}, P_{\ell-1}, S_{\ell-1}, \gamma_\ell)$ 
      else if  $\text{hybrid\_galerkin}$  {Re-add entries at the new dropping tolerance}
         $\hat{A}_\ell = \text{sparsify}(A_\ell, \hat{A}_{\ell-1}, P_{\ell-1}, S_{\ell-1}, \gamma_\ell)$ 

```

---

627 suffer after 3 iterations, the hierarchy is updated again, but since  $\hat{A}_2$  has  $\gamma_2 = 0.0$ ,  
628 entries are reintroduced into coarse matrices  $\hat{A}_3$  and  $\hat{A}_4$  instead.

629 Using Algorithm 5, Figure 18 shows both the relative residual of the system after  
630 each iteration as well as the communication costs of PCG using three different AMG  
631 hierarchies: standard Galerkin, Sparse Galerkin with diagonal lumping and aggressive  
632 dropping, and Sparse Galerkin with diagonal lumping modified with adaptivity. For  
633 the adaptive case, we purposefully choose an overly aggressive initial drop tolerance  
634 so that entries can be added back multiple times and one coarse level at a time to  
635 show the effect on convergence and communication. Initially, when the drop tolerance  
636 is aggressive, the associated communication costs are low, but the resulting PCG it-

637 erations do not converge; this provides a baseline. As sparse entries are reintroduced  
 638 into the hierarchy, convergence improves, while only slightly increasing the associated  
 639 communication cost. When entries are reintroduced into the hierarchy, the precondi-  
 640 tioner for PCG changes, and hence, the method must be restarted. After restarting  
 the method, convergence improves.

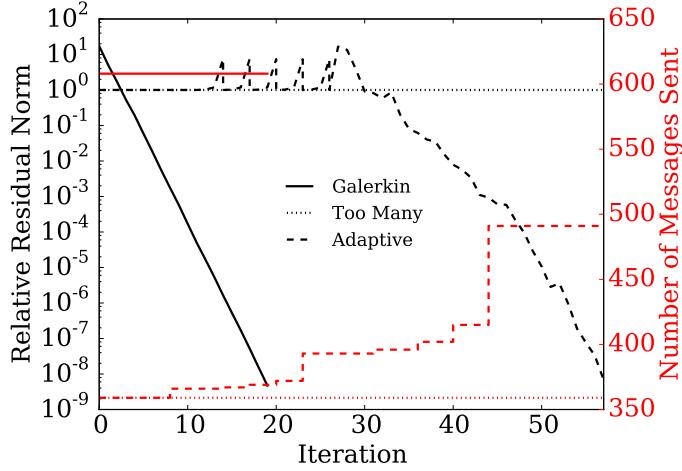


Fig. 18: Relative residual (black) and number of MPI sends (red) per iteration when solving the **Laplace** problem with: (1) PCG using Galerkin AMG; (2) Hybrid Galerkin with aggressive dropping (labeled Too Many); (3) Hybrid Galerkin solved with Algorithm 5, using  $k = 3$ ,  $s = 1$ , and  $\gamma_0 \dots \gamma_L$  set to the same drop tolerances as the aggressive case.

641

642 **7. Conclusion.** We have introduced a lossless method to reduce the work re-  
 643 quired in parallel algebraic multigrid by removing weak or unimportant entries from  
 644 coarse-grid operators after the multigrid hierarchy is formed. This alternative to  
 645 the original method of non-Galerkin coarse grids is similarly capable of reducing the  
 646 communication costs on coarse levels, yielding an overall reduction in solve times.  
 647 Furthermore, this method retains the original Galerkin hierarchy, allowing many of  
 648 the restrictions of non-Galerkin to be relaxed. As a result, removed entries are easily  
 649 lumped directly to the diagonals, greatly reducing setup costs, while also reducing  
 650 communication complexity during the solve phase. Furthermore, as entries are added  
 651 to the diagonal, entries removed from the matrix are stored and adaptively reintro-  
 652 duced into the hierarchy if necessary for convergence. Hence, the trade-off between  
 653 convergence and the communication costs is controlled at solve-time with little addi-  
 654 tional work.

655

## REFERENCES

- 656 [1] *Blue Waters*. <https://bluewaters.ncsa.illinois.edu/>.  
 657 [2] *HYPRE: High performance preconditioners*. <http://www.llnl.gov/CASC/hypre/>.  
 658 [3] S. F. ASHBY AND R. D. FALGOUT, *A parallel multigrid preconditioned conjugate gradient al-*  
 659 *gorithm for groundwater flow simulations*, Nuclear Science and Engineering, 124 (1996),  
 660 pp. 145–159. UCRL-JC-122359.

- [4] ALLISON H. BAKER, MARTIN SCHULZ, AND ULRIKE M. YANG, *On the performance of an algebraic multigrid solver on multicore clusters*, in Proceedings of the 9th International Conference on High Performance Computing for Computational Science, VECPAR'10, Berlin, Heidelberg, 2011, Springer-Verlag, pp. 102–115.
- [5] BRETT BODE, MICHELLE BUTLER, THOM DUNNING, TORSTEN HOEFLER, WILLIAM KRAMER, WILLIAM GROPP, AND WEN-MEI HWU, *The Blue Waters super-system for super-science*, in Contemporary High Performance Computing: From Petascale Toward Exascale, Jeffrey S. Vetter, ed., vol. 1 of CRC Computational Science Series, Taylor and Francis, Boca Raton, 1 ed., 2013, pp. 339–366.
- [6] MATTHIAS BOLTEN, THOMAS K. HUCKLE, AND CHRISTOS D. KRAVVARITIS, *Sparse matrix approximations for multigrid methods*, Linear Algebra and its Applications, (2015), pp. –.
- [7] A. BRANDT, S. F. MCCORMICK, AND J. W. RUGE, *Algebraic multigrid (AMG) for sparse matrix equations*, in Sparsity and Its Applications, D. J. Evans, ed., Cambridge Univ. Press, Cambridge, 1984, pp. 257–284.
- [8] E. CHOW, R. D. FALGOUT, J. J. HU, R. S. TUMINARO, AND U. M. YANG, *A survey of parallelization techniques for multigrid solvers*, in Frontiers of Parallel Processing for Scientific Computing, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2005.
- [9] JE DENDY, *Black box multigrid*, Journal of Computational Physics, 48 (1982), pp. 366–386.
- [10] J. DENDY, *Black box multigrid for nonsymmetric problems*, Appl. Math. Comput., 13 (1983), pp. 261–283.
- [11] IRAD YAVNEH ERAN TREISTER, RAN ZEMACH, *Algebraic collocation coarse approximation (acca multigrid)*, in 12th Copper Mountain Conference on Iterative Methods, 2012.
- [12] ROBERT D. FALGOUT AND JACOB B. SCHRODER, *Non-Galerkin coarse grids for algebraic multigrid*, SIAM Journal on Scientific Computing, 36 (2014), pp. C309–C334.
- [13] HORMOZD GAHVARI, ALLISON H. BAKER, MARTIN SCHULZ, ULRIKE MEIER YANG, KIRK E. JORDAN, AND WILLIAM GROPP, *Modeling the performance of an algebraic multigrid cycle on hpc platforms*, in Proceedings of the International Conference on Supercomputing, ICS '11, New York, NY, USA, 2011, ACM, pp. 172–181.
- [14] HORMOZD GAHVARI, MARK HOEMMEN, JAMES DEMMEL, AND KATHERINE YELICK, *Benchmarking sparse matrix-vector multiply in five minutes*, in SPEC Benchmark Workshop 2007, Austin, TX, January 2007.
- [15] VAN EMDEN HENSON AND ULRIKE MEIER YANG, *Boomeramg: A parallel algebraic multigrid solver and preconditioner*, Appl. Numer. Math., 41 (2002), pp. 155–177.
- [16] PIOTR LUSCZEK, JACK J. DONGARRA, DAVID KOESTER, ROLF RABENSEIFNER, BOB LUCAS, JEREMY KEPNER, JOHN MCCALPIN, DAVID BAILEY, AND DAISUKE TAKAHASHI, *Introduction to the HPC challenge benchmark suite*, (2005).
- [17] S. F. MCCORMICK AND J. W. RUGE, *Multigrid methods for variational problems*, SIAM J. Numer. Anal., 19 (1982), pp. 924–929.
- [18] J. W. RUGE AND K. STÜBEN, *Algebraic multigrid (AMG)*, in Multigrid Methods, S. F. McCormick, ed., Frontiers Appl. Math., SIAM, Philadelphia, 1987, pp. 73–130.
- [19] HANS DE STERCK, ROBERT D. FALGOUT, JOSHUA W. NOLTING, AND ULRIKE MEIER YANG, *Distance-two interpolation for parallel algebraic multigrid*, Numerical Linear Algebra with Applications, 15 (2008), pp. 115–139.
- [20] HANS DE STERCK, ULRIKE MEIER YANG, AND JEFFREY J. HEYS, *Reducing complexity in parallel algebraic multigrid preconditioners*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 1019–1039.
- [21] ERAN TREISTER AND IRAD YAVNEH, *Non-Galerkin multigrid based on sparsified smoothed aggregation*, SIAM Journal on Scientific Computing, 37 (2015), pp. A30–A54.
- [22] ROMAN WIENANDS AND IRAD YAVNEH, *Collocation coarse approximation in multigrid*, SIAM Journal on Scientific Computing, 31 (2009), pp. 3643–3660.
- [23] ULRIKE MEIER YANG, *On long-range interpolation operators for aggressive coarsening*, Numerical Linear Algebra with Applications, 17 (2010), pp. 453–472.