

1
2

SCALABLE ASYNCHRONOUS DOMAIN DECOMPOSITION SOLVERS*

3 CHRISTIAN GLUSA[†], ERIK G. BOMAN[‡], EDMOND CHOW[§], SIVASANKARAN
4 RAJAMANICKAM[¶], AND DANIEL B. SZYLD^{||}

5 **Abstract.** Parallel implementations of linear iterative solvers generally alternate between phases
6 of data exchange and phases of local computation. Increasingly large problem sizes and more hetero-
7 geneous compute architectures make load balancing and the design of low latency network intercon-
8 nects that are able to satisfy the communication requirements of linear solvers very challenging tasks.
9 In particular, global communication patterns such as inner products become increasingly limiting at
10 scale.

11 We explore the use of asynchronous communication based on one-sided MPI primitives in the
12 context of domain decomposition solvers. In particular, a scalable asynchronous two-level Schwarz
13 method is presented. We discuss practical issues encountered in the development of a scalable solver
14 and show experimental results obtained on a state-of-the-art supercomputer system that illustrate
15 the benefits of asynchronous solvers in load balanced as well as load imbalanced scenarios. Using the
16 novel method, we can observe speed-ups of up to 4x over its classical synchronous equivalent.

17 **Key words.** Asynchronous iteration, domain decomposition, Schwarz methods, chaotic relax-
18 ation

19 **AMS subject classifications.** 68W10, 65Y05, 68W15, 65N55

20 **1. Introduction.** Multilevel methods such as multigrid and domain decomposi-
21 tion are among the most efficient and scalable solvers for partial differential equations
22 developed to date. Adapting them to the next generation of supercomputers and
23 improving their performance and scalability is crucial in the push towards exascale.
24 Domain decomposition methods subdivide the global problem into subdomains, and
25 then alternate between local solves and boundary data exchange. This puts a signif-
26 icant stress on the network interconnect, since all processes try to communicate at
27 once. On the other hand, during the solve phase, the network is under-utilized. The
28 use of non-blocking communication can only alleviate this issue, but not fully resolve
29 it. In asynchronous methods, on the other hand, computation and communication
30 occur at the same time, with some processes performing computation while others
31 communicate, so that the network is consistently in use.

32 The term “asynchronous” can have several different meanings in the literature.
33 In computer science, it is sometimes used to describe communication patterns that
34 are non-blocking, so that computation and communication can be overlapped. Itera-
35 tive algorithms that use such “asynchronous” communication yield the same iterates
36 (results) up to round-off error, as they do not change the mathematical algorithm. In
37 applied mathematics, on the other hand, “asynchronous” denotes parallel algorithms
38 where each process (processor) proceeds at its own speed without synchronization.

* Part of this work has been accepted for publication in the form of a proceedings paper by the 25th International Domain Decomposition Conference.

[†]Center for Computing Research, Sandia National Laboratories, Albuquerque, New Mexico, USA (caglusa@sandia.gov).

[‡]Center for Computing Research, Sandia National Laboratories, Albuquerque, New Mexico, USA (egboman@sandia.gov).

[§]School of Computational Science and Engineering, College of Computing, Georgia Institute of Technology, Atlanta, Georgia, USA (echow@cc.gatech.edu).

[¶]Center for Computing Research, Sandia National Laboratories, Albuquerque, New Mexico, USA (srajama@sandia.gov).

^{||}Temple University, Philadelphia, Pennsylvania, USA (szyld@temple.edu).

39 Thus, asynchronous algorithms go beyond the widely used bulk-synchronous paral-
40 lel (BSP) model. More importantly, they are mathematically different than syn-
41 chronous methods and generate different iterates. The earliest work in this area was
42 called “chaotic relaxation” [11]. Both types of asynchronous approaches are expected
43 to play an important role on future supercomputers. In this paper, we focus on
44 asynchronous methods in the mathematical sense, and we will use the terms “asyn-
45 chronous” and “synchronous” to distinguish between methods that are asynchronous
46 and synchronous in the mathematical sense.

47 Domain decomposition solvers [16, 34, 33] are often used as preconditioners in
48 Krylov subspace iterations. Unfortunately, the computation of inner products and
49 norms widely used in Krylov methods requires global communication. Global com-
50 munication primitives, such as `MPI_Reduce`, asymptotically scale as the logarithm of
51 the number of processes involved. This can become a limiting factor when very large
52 process counts are used. The underlying domain decomposition method, however,
53 can do away with globally synchronous communication, assuming the coarse prob-
54 lem in multilevel methods can be solved in a parallel way. Therefore, we will focus
55 on using domain decomposition methods purely as iterative methods in the present
56 work. We will note, however, that the discussed algorithms could be coupled with ex-
57 isting pipelined methods [22] which alleviate the global synchronization requirement
58 of Krylov solvers.

59 Another issue that is crucial to good scaling behavior is load imbalance. Load
60 imbalance might occur due to heterogeneous hardware in the system, network noise,
61 dynamic power capping [1], or due to local, problem specific causes, such as iteration
62 counts for local solves that vary from subdomain to subdomain. The latter are espe-
63 cially difficult to predict, so that load balancing cannot occur before the actual solve.
64 Therefore, processes in a synchronous parallel program must be idle until its slowest
65 process has finished. In an asynchronous method, local computation can continue,
66 and potentially improve the quality of the global solution.

67 An added benefit of asynchronous methods is that, since the interdependence
68 between subdomains has been weakened, fault tolerance [9, 10] can be more easily
69 achieved. When one process must stop, be it for a hard or a soft fault, it can be
70 replaced without having to halt every other process.

71 The main drawback of asynchronous iterations is the fact that deterministic be-
72 havior is sacrificed. Consecutive runs do not produce the same result. (But one would
73 hope that they are at most a distance proportional to the convergence tolerance apart
74 from each other.) This also makes the mathematical analysis of asynchronous methods
75 significantly more difficult than for its synchronous counterparts. Analytical frame-
76 works for asynchronous linear (and nonlinear) iterations have long been available
77 [11, 4, 5, 18], but generally cannot produce sharp convergence bounds except in the
78 simplest cases.

79 The main contributions of our work are:

- 80 • A novel asynchronous two-level domain decomposition method, scalable to
81 thousands of processors.
- 82 • An empirical study of one-sided MPI performance in a scientific computing
83 setting.
- 84 • Empirical comparisons of synchronous and asynchronous variants of domain
85 decomposition solvers on a state-of-the-art parallel computer.

86 Our work demonstrates that asynchronous methods have the potential of outper-
87 forming conventional synchronous solvers and offer a viable alternative in the push
88 towards exascale.

89 The present work is structured as follows: In [Section 2](#), we present overlapping do-
 90 main decomposition methods, and explain their use in synchronous and asynchronous
 91 fashion. For a general introduction to domain decomposition methods we refer the
 92 reader to [\[16, 34, 33\]](#). The section concludes with a convergence analysis of the pre-
 93 sented one- and two-level methods. [Section 3](#) is dedicated to a description of the
 94 presently available mechanisms in MPI and hardware to achieve truly asynchronous
 95 communication. Numerical experiments exploring asynchronous communication and
 96 using the presented domain decomposition methods are given in [Section 4](#), where
 97 we compare the strong and weak scaling behavior of synchronous and asynchronous
 98 solvers with and without load imbalance.

99 **1.1. Related work.** An asynchronous one-level domain decomposition solver
 100 with optimized artificial boundary conditions was proposed in [\[30\]](#); see also [\[21, 20, 17\]](#)
 101 for its analysis in two different settings. An implementation of asynchronous optimized
 102 Schwarz is described in [\[36\]](#). An optimization package that leverages asynchronous
 103 coordinate updates is presented in [\[31\]](#). An asynchronous multigrid method for shared
 104 memory systems was proposed in [\[35\]](#). Synchronization reducing Krylov methods have
 105 a long history [\[14\]](#). However, preconditioning such methods is unresolved apart from
 106 some simple preconditioners [\[13\]](#). Recent work extends their applicability to one level
 107 domain decomposition preconditioning [\[37\]](#). Pipelined Krylov methods [\[22\]](#) reduce
 108 synchronization costs by overlapping inner products with matrix-vector products and
 109 preconditioner applications, and can be used with any preconditioner.

110 2. Domain decomposition methods.

111 **2.1. One-level Restricted Additive Schwarz (RAS).** We want to solve the
 112 global system

$$113 \quad \mathbf{A}u = f,$$

115 where $\mathbf{A} \in \mathbb{R}^{N \times N}$ arises from the finite element or finite difference discretization of
 116 a partial differential equation. Informally, one-level domain decomposition solvers
 117 break up the global system of equations into overlapping sub-problems that cover the
 118 whole global system. This requires that the matrix \mathbf{A} is sparse and couples unknowns
 119 only in a local manner.

120 The iteration then alternates between computation of the global residual, which
 121 involves communication, and local solves for solution corrections. Special attention
 122 needs to be paid to the unknowns in the overlap, in order to avoid over-correction.
 123 Below, we describe the different methods considered in this work in detail in order
 124 to understand what data is required to be exchanged and how the methods can be
 125 executed in asynchronous fashion.

126 Based on the graph of \mathbf{A} or geometric information for the underlying problem
 127 the unknowns are grouped into P overlapping sets \mathcal{N}_p of size N_p , $p = 1, \dots, P$. An
 128 example of such a partitioning is given in [Figure 2.1](#). We further split the sets \mathcal{N}_p
 129 into

$$130 \quad \mathcal{S}_p := \{j \in \mathcal{N}_p \mid \exists k \in \mathcal{N}_p^c : \mathbf{A}_{jk} \neq 0\},$$

132 i.e., unknowns that are on the boundary of the set \mathcal{N}_p , and interior unknowns $\mathcal{I}_p :=$
 133 $\mathcal{N}_p \setminus \mathcal{S}_p$.

134 The notation throughout this section is based on Dolean et al. [\[16\]](#). We call the
 135 restriction to the p -th set $\mathbf{R}_p \in \mathbb{R}^{N_p \times N}$. The entries of the matrices \mathbf{R}_p are all either

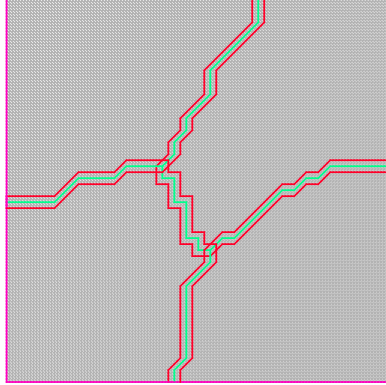


FIG. 2.1. Partitioning of a uniform triangular mesh of the unit square into 4 overlapping subdomains. The non-overlapping partitioning produced using METIS [25] is shown in green; the extended overlapping subdomains are shown in red.

136 one or zero, with exactly one entry per row and at most one entry per column being
 137 non-zero. The local parts of \mathbf{A} are given by

138
$$\mathbf{A}_p = \mathbf{R}_p \mathbf{A} \mathbf{R}_p^T \in \mathbb{R}^{N_p \times N_p}.$$

140 Furthermore, we require a partition of unity, represented by diagonal weighting
 141 matrices \mathbf{D}_p , such that the discrete partition of unity property holds

142 (2.1)
$$\mathbf{I} = \sum_{p=1}^P \mathbf{R}_p^T \mathbf{D}_p \mathbf{R}_p.$$

 143

144 In what follows, we will assume that \mathbf{D}_p are Boolean, i.e. their entries are either zero
 145 or one. This means that every (potentially shared) unknown has a special attachment
 146 with exactly one subdomain. We will furthermore require that $(\mathbf{D}_p)_{jj} = 0$ for all
 147 surface unknowns $j \in \mathcal{S}_p$. One way of satisfying these restrictions is to extend overlaps
 148 starting with a *non-overlapping* partition and then define the special attachment via
 149 the partition.

150 Consequently,

151 (2.2)
$$\mathbf{D}_p \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q = \mathbf{0} \quad \text{for } p \neq q$$

153 and

154 (2.3)
$$\mathbf{D}_p \mathbf{R}_p \mathbf{R}_p^T \mathbf{D}_p = \mathbf{D}_p.$$

 155

156 Moreover, the identity

157 (2.4)
$$\mathbf{R}_p \mathbf{A} \mathbf{R}_q^T \mathbf{D}_q = \mathbf{R}_p \mathbf{R}_q^T \mathbf{R}_q \mathbf{A} \mathbf{R}_q^T \mathbf{D}_q$$

 158

159 holds, since for any $u_q \in \mathbb{R}^{N_q}$, $\mathbf{D}_q u_q$ is supported on the interior unknowns \mathcal{I}_p , and
 160 hence $\mathbf{A} \mathbf{R}_q^T \mathbf{D}_q u_q$ is supported in \mathcal{N}_q . But on \mathcal{N}_q , $\mathbf{R}_q^T \mathbf{R}_q$ acts as the identity.

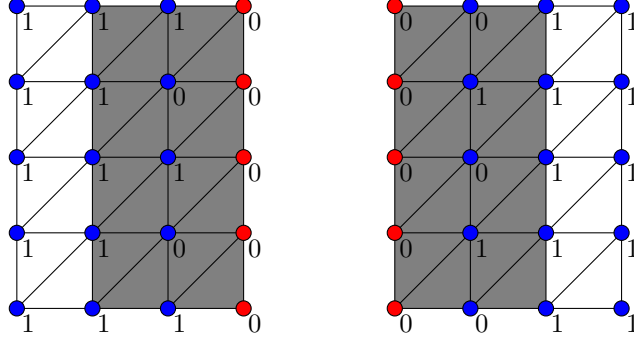


FIG. 2.2. Two overlapping subdomains. The overlap between the subdomains is shaded in gray; the respective surface sets S_\bullet are shown by red circles, the interior unknowns I_\bullet as blue circle. The diagonal values of the respective D_\bullet are shown next to the nodes.

161 A stationary iterative method based on the splitting $\mathbf{A} = \mathbf{M} - \mathbf{N}$ is given globally
 162 as

$$163 \quad u^{n+1} = u^n + \mathbf{M}^{-1} (f - \mathbf{A}u^n),$$

165 where \mathbf{M}^{-1} is a preconditioner for \mathbf{A} .

166 This means that we need to calculate the residual $r^n = f - \mathbf{A}u^n$. Its local part
 167 on node p is given by

$$\begin{aligned}
 168 \quad \mathbf{R}_p r^n &= \mathbf{R}_p f - \mathbf{R}_p \mathbf{A} u^n \\
 169 &= \mathbf{R}_p \left(\sum_{q=1}^P \mathbf{R}_q^T \mathbf{D}_q \mathbf{R}_q \right) f - \mathbf{R}_p \mathbf{A} \left(\sum_{q=1}^P \mathbf{R}_q^T \mathbf{D}_q \mathbf{R}_q \right) u^n \\
 170 &= \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q \mathbf{R}_q f - \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{A}_q \mathbf{D}_q \mathbf{R}_q u^n \\
 171 &= \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q \mathbf{R}_q u^n),
 \end{aligned}$$

173 where we used (2.1) and (2.4). This means that in order to obtain the local part of
 174 the global residual, we first compute locally $\mathbf{D}_p \mathbf{R}_p f - \mathbf{A}_p \mathbf{D}_p \mathbf{R}_p u^n$ on every node p ,
 175 and then communicate and accumulate the overlapping parts of these local residual
 176 vectors. The latter operation is represented by the operator $\sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T$.

177 The *restricted additive Schwarz (RAS) preconditioner* [8, 7] is given by

$$178 \quad \mathbf{M}_{RAS}^{-1} = \sum_{p=1}^P \mathbf{R}_p^T \mathbf{D}_p \mathbf{A}_p^{-1} \mathbf{R}_p.$$

180 RAS is widely used and is the default option for overlapping domain decomposition
 181 preconditioners in PETSc [3]. It can be thought of as a variant of the additive Schwarz
 182 preconditioner

$$183 \quad \mathbf{M}_{AS}^{-1} = \sum_{p=1}^P \mathbf{R}_p^T \mathbf{A}_p^{-1} \mathbf{R}_p$$

185 that is convergent as an iterative method, since the damping by \mathbf{D}_p in the overlapping
 186 parts avoids over-correction; see [18]. Note that for a natural choice of \mathbf{D}_p , the number
 187 of communication steps is cut in half as there is no communication associated with
 188 $\mathbf{R}_p^T \mathbf{D}_p$.

189 Now, the local part of the RAS iteration is given by

$$\begin{aligned}
 190 \quad \mathbf{R}_p u^{n+1} &= \mathbf{R}_p u^n + \mathbf{R}_p \mathbf{M}_{RAS}^{-1} r^n \\
 191 &= \mathbf{R}_p u^n + \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q \mathbf{A}_q^{-1} \mathbf{R}_q r^n.
 \end{aligned}$$

193 If we set $u_p^n = \mathbf{R}_p u^n$ and $r_p^n = \mathbf{R}_p r^n$ as the local parts of solution and residual
 194 respectively, the RAS iteration is

$$\begin{aligned}
 195 \quad r_p^n &= \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q u_q^n), \\
 196 \quad u_p^{n+1} &= u_p^n + \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q \mathbf{A}_q^{-1} r_q^n.
 \end{aligned}$$

198 This seems to suggest that the update step requires neighborhood communication as
 199 well. But in fact, in the next iteration, computation of the residual only requires
 200 $\mathbf{D}_p u_p^{n+1}$. From (2.2), (2.3), we see that the iterative scheme without the communica-
 201 tion step in the update

$$202 \quad (2.5) \quad r_p^n = \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q w_q^n),$$

$$203 \quad (2.6) \quad w_p^{n+1} = w_p^n + \mathbf{A}_p^{-1} r_p^n$$

205 is equivalent because $\mathbf{D}_p u_p^n = \mathbf{D}_p w_p^n$ for all n . The solution u_p^n can be recovered from
 206 w_p^n in the post-processing step

$$207 \quad u_p^n = \mathbf{R}_p u^n = \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q \mathbf{R}_q u^n = \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q w_q^n.$$

209 Finally, we use the norm of the residual in the stopping criterion. The norm can
 210 be computed from local quantities as

$$\begin{aligned}
 211 \quad \|r^n\|^2 &= r^n \cdot r^n = r^n \cdot \left(\sum_{p=1}^P \mathbf{R}_p^T \mathbf{D}_p \mathbf{R}_p r^n \right) \\
 212 &= \sum_{p=1}^P (\mathbf{R}_p r^n) \cdot (\mathbf{D}_p \mathbf{R}_p r^n) = \sum_{p=1}^P r_p^n \cdot (\mathbf{D}_p r_p^n).
 \end{aligned}$$

214 In conclusion, we can give the local form of RAS as in [Algorithm 2.1](#), where
 215 we have dropped the superscript n for the iteration number. In fact, [Algorithm 2.1](#)
 216 describes both the synchronous *and* the asynchronous version of RAS. In the syn-
 217 chronous version, line 4 is executed in lock step fashion by all subdomains using
 218 non-blocking two-sided communication primitives. This communication step could
 219 be overlapped by computation. However, in established frameworks such as Trilinos,
 220 such overlapping requires major changes to the framework¹. PETSc allows some over-

¹<https://github.com/trilinos/Trilinos/issues/767>

Algorithm 2.1 Restricted additive Schwarz (RAS) in local form, “ \leftrightarrow ” signifies communication.

```

1:  $w_p \leftarrow 0$ 
2: while not converged do
3:   Local residual:  $s_p \leftarrow \mathbf{D}_p \mathbf{R}_p f - \mathbf{A}_p \mathbf{D}_p w_p$ 
4:   Accumulate:  $r_p \leftarrow \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T s_q$   $\leftrightarrow$ 
5:   Solve:  $\mathbf{A}_p v_p = r_p$ 
6:   Update:  $w_p \leftarrow w_p + v_p$ 
7: end while
8: Post-process:  $u_p \leftarrow \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T \mathbf{D}_q w_q$   $\leftrightarrow$ 

```

221 lap of computation and communication with two-phase assembly [3]. It is possible to
222 modify such established libraries for the asynchronous iterations of this paper. How-
223 ever, in order to keep the focus on algorithmic development, we developed a library
224 that supports the one-sided communication primitives, and build the new solvers using
225 the communication primitives.

226 In the asynchronous variant, each subdomain exposes a memory region for remote
227 access. On execution of line 4, the relevant components of the current local residual
228 vector $s_p = \mathbf{D}_p \mathbf{R}_p f - \mathbf{A}_p \mathbf{D}_p w_p$ are written to the neighboring subdomains, and the
229 latest locally available data s_q from every neighbor q is used. We refer to Section 4.1
230 for a discussion of the options for actually achieving this neighborhood exchange in
231 practice. The implementation of a convergence check (as used on line 2) that does
232 not require synchronization is detailed in Section 4.4.

233 **2.2. Two-level synchronous RAS.** In order to improve the scalability of the
234 solver, a mechanism of global information exchange is required. Let $\mathbf{R}_0 \in \mathbb{R}^{n_0 \times n}$ be
235 the restriction from the fine grid problem to a coarser mesh, and let the coarse-grid
236 matrix \mathbf{A}_0 be given by the Galerkin relation $\mathbf{A}_0 = \mathbf{R}_0 \mathbf{A} \mathbf{R}_0^T$. The coarse-grid solve
237 can be incorporated in the RAS iteration either in additive fashion:

$$238 \quad (2.7) \quad u^{n+1} = u^n + \left(\frac{1}{2} \mathbf{M}_{RAS}^{-1} + \frac{1}{2} \mathbf{R}_0^T \mathbf{A}_0^{-1} \mathbf{R}_0 \right) (f - \mathbf{A} u^n),$$

239
240 or in multiplicative fashion:

$$241 \quad u^{n+1/2} = u^n + \mathbf{R}_0^T \mathbf{A}_0^{-1} \mathbf{R}_0 (f - \mathbf{A} u^n),$$

$$242 \quad u^{n+1} = u^{n+1/2} + \mathbf{M}_{RAS}^{-1} (f - \mathbf{A} u^{n+1/2}).$$

244 In what follows, we focus on the additive version, since it naturally lends itself to
245 asynchronous iterations: subdomain solves and coarse-grid solves are independent of
246 each other.

247 We now determine the local form of the global algorithm. It is understood that
248 the solve with \mathbf{A}_0 itself might be distributed over several processes. This internal
249 computation is not meant to be performed in an asynchronous manner, which is why
250 we do not need to further explore the local form of the coarse-grid solve. For simplicity
251 of exposition we therefore do not describe the solution of the coarse-grid problem itself
252 in local form, i.e. we will simply write \mathbf{A}_0^{-1} . The local part of the coarse-grid update

Algorithm 2.2 Synchronous RAS with additive coarse grid in local form, “ \leftrightarrow ” signifies communication.

```

1:  $w_p \leftarrow 0$ 
2: while not converged do
3:   On subdomains
4:     Local residual:  $s_p \leftarrow D_p R_p f - A_p D_p w_p$ 
5:     Send  $R_0 R_p^T s_p$  to coarse grid  $\leftrightarrow$ 
6:     Accumulate:  $r_p \leftarrow \sum_{q=1}^P R_p R_q^T s_q$   $\leftrightarrow$ 
7:     Solve:  $A_p v_p = r_p$ 
8:     Update:  $w_p \leftarrow w_p + \frac{1}{2} v_p$ 
9:     Receive  $c_p = R_p R_0^T v_0$  from coarse grid  $\leftrightarrow$ 
10:    Update:  $w_p \leftarrow w_p + \frac{1}{2} c_p$ 
11:   On coarse grid
12:     Receive  $R_0 R_p^T s_p$  from subdomains  $\leftrightarrow$ 
13:     Accumulate  $r_0 = \sum_{p=1}^P R_0 R_p^T s_p$ 
14:     Solve  $A_0 v_0 = r_0$ 
15:     Send  $c_p = R_p R_0^T v_0$ ,  $p = 1, \dots, P$  to subdomains  $\leftrightarrow$ 
16:   end while
17:   On subdomains
18:     Post-process  $u_p \leftarrow \sum_{q=1}^P R_p R_q^T D_q w_q$   $\leftrightarrow$ 

```

253 is

$$\begin{aligned}
254 \quad & \frac{1}{2} R_p R_0^T A_0^{-1} R_0 (f - A u^n) \\
255 \quad & = \frac{1}{2} \left(R_p R_0^T \right) A_0^{-1} \sum_{p=1}^P \left(R_0 R_p^T \right) (D_p R_p f - A_p R_p u^n). \\
256 \quad &
\end{aligned}$$

257 Here, the operators $\left(R_0 R_p^T \right)$ and $\left(R_p R_0^T \right)$ encode the communication from sub-
258 domain p to the coarse grid and vice versa. We notice that while the communica-
259 tion among subdomains consist in one neighborhood data exchange per iteration, the
260 coarse-grid solve involves sending data from the subdomains to the coarse grid, and
261 sending a solution from the coarse grid to the subdomains. In conclusion, the local
262 form of RAS with an additive coarse grid is given in [Algorithm 2.2](#). Again, we have
263 dropped the superscript for the iteration number. The communication between coarse
264 and fine grid can be implemented in multiple ways. Since we want to allow the coarse
265 grid solve to be distributed itself and the same coarse unknown can be owned by
266 several coarse grid ranks (just as is the case for the fine grid), we do not consider
267 options involving `MPI_Reduce/MPI_Bcast` or `MPI_Gather/MPI_Scatter` or their non-
268 blocking equivalents. Instead, we opted for use of `MPI_Isend` and `MPI_Irecv`. A future
269 improvement could involve the use of intercommunicators and `MPI_Iallgatherv` or
270 other collectives. The advantage of the current approach is that the changes between
271 synchronous and asynchronous implementation of the communication layer (described
272 in the next section) are minimal.

273 **2.3. Two-level asynchronous RAS.** From the mathematical description (2.7)
274 of two-level additive RAS, one might be tempted to see the coarse-grid problem simply
275 as an additional subdomain. From [Algorithm 2.2](#) the fundamental differences between

276 the subdomains and the coarse-grid problem become apparent. Subdomains deter-
277 mine the right-hand side for their local solve and correct it by transmitting boundary
278 data to their neighbors. The coarse grid, on the other hand, receives its entire right-
279 hand side from the subdomains, and hence it has to communicate with every single
280 one of them.

281 In order to perform asynchronous coarse-grid solves, we therefore need to make
282 sure that all the right-hand side data necessary for the solve has been received by the
283 processes responsible for the coarse grid. Moreover, corrections sent by the coarse
284 grid should be used exactly once by the subdomains. This is achieved by not only
285 allocating memory regions to hold the coarse-grid right-hand side on the coarse-grid
286 processes and the coarse-grid correction on the subdomains, but also Boolean variables
287 that are polled to determine whether writing or reading right-hand side or solution
288 data is permitted. More precisely, writing of the local subdomain residuals to the
289 coarse-grid memory region of r_0 is contingent upon the state of the Boolean variable
290 `canWriteRHSp`. (See [Algorithm 2.3](#).) When `canWriteRHSp` is `True`, right-hand side
291 data is written to the coarse grid, otherwise this operation is omitted. Here, the
292 subscripts are used to signify the MPI rank owning the accessed memory region. As
293 before, index 0 corresponds to the (potentially distributed) coarse grid and indices
294 $1, \dots, P$ correspond to the subdomains. To improve readability, we show access to a
295 memory region on the calling process in blue, while remote access is printed in red.

296 In a similar fashion, the coarse grid checks whether every subdomain has written a
297 right-hand side to r_0 by polling the state of the local Boolean array `RHSisReady0`. The
298 communication of the obtained coarse-grid solution back to the subdomains follows the
299 same pattern, using the variables `solutionIsReadyp`. The subdomains update their
300 current iterate using the local subdomain solution and the coarse-grid solution. If the
301 latter is not available, the subdomain solution is used unweighted. If both solutions
302 are available, then the same weighting (1/2, 1/2) as in the synchronous case (2.7)
303 is used. We note that the algorithm is asynchronous despite the data dependencies.
304 Coarse grid and subdomain solves do not wait for each other.

305 We determined by experiments that overall performance is adversely affected
306 if the coarse grid constantly polls the status variable `RHSisReady0`, waiting for all
307 subdomains to provide right-hand side information. Therefore, we added a sleep
308 statement into its work loop. If the sleep interval is too short, the sleep statement
309 is ineffective. If the sleep interval is too large, the coarse grid will be under-used.
310 Keeping the ratio of attempted coarse-grid solves (i.e. reads from `RHSisReady0`) to
311 actual performed coarse-grid solves at around 1/20 has been proven effective to us.
312 This can easily be achieved by an adaptive procedure that counts both successful
313 solves and solve attempts and then either increases or decreases the sleep interval
314 accordingly.

315 **2.4. Convergence Analysis of Asynchronous Iterations.** We present below
316 the mathematical framework used to describe and study asynchronous algorithms. We
317 modify the model introduced by Bertsekas [5], [6] to take into account the fact that
318 data available at a process p from another process q might have been produced during
319 different local iterations. This issue can arise when data is accessed on process p while
320 it is being overwritten by a new transmission from process q .

321 For a mathematical model of these asynchronous iterations on P processors, let
322 us denote by $\{\sigma_n\}_{n \in \mathbb{N}}$ the sequence of non-empty subsets of $\{1, \dots, P\}$, defining which
323 processes update their components at the “iteration” n , where here “iteration” can
324 be thought of as a time stamp. We call these *sets of update indices*. Define further for

Algorithm 2.3 Asynchronous RAS with additive coarse grid in local form. Variables printed in blue are exposed memory regions that are local to the calling process. Red variables are remote memory regions. Subscripts denote the owning process of the variable. Array access is denoted by “[.]”.

```

1: while not converged do
2:   On subdomains
3:     Local residual:  $s_p \leftarrow D_p R_p f - A_p D_p w_p$ 
4:     if canWriteRHSp then
5:        $r_0 \leftarrow r_0 + R_0 R_p^T s_p$ 
6:       canWriteRHSp  $\leftarrow$  False
7:       RHSisReady0[p]  $\leftarrow$  True
8:     end if
9:     Accumulate asynchronously:  $r_p \leftarrow \sum_{q=1}^P R_p R_q^T s_q$ 
10:    Solve:  $A_p v_p = r_p$ 
11:    if solutionIsReadyp then
12:      Update:  $w_p \leftarrow w_p + \frac{1}{2} v_p + \frac{1}{2} c_p$ 
13:      solutionIsReadyp  $\leftarrow$  False
14:    else
15:      Update:  $w_p \leftarrow w_p + v_p$ 
16:    end if
17:    On coarse grid
18:    if RHSisReady0[p]  $\forall p = 1, \dots, P$  then
19:      Solve  $A_0 v_0 = r_0$ 
20:      for  $p = 1, \dots, P$  do
21:        RHSisReady0[p]  $\leftarrow$  False
22:        canWriteRHSp  $\leftarrow$  True
23:         $c_p \leftarrow R_p R_0^T v_0$ 
24:        solutionIsReadyp  $\leftarrow$  True
25:      end for
26:    else
27:      Sleep (time adjusted adaptively)
28:    end if
29:  end while
30:  On subdomains
31:  Post-process synchronously  $u_p \leftarrow \sum_{q=1}^P R_p R_q^T D_q w_q$ 

```

325 $p, q \in \{1, \dots, P\}$, $\{\tau_{q,n}^{(p)}\}_{n \in \mathbb{N}}$ a sequence of integer vectors, where $\left(\tau_{q,n}^{(p)}\right)_i$, $1 \leq i \leq N_q$
326 represents the iteration number (or time stamp) of the i -th component of data coming
327 from process q and available on process p at the beginning of the computation of
328 the process which produces $u_{p,n}$ at time n . Thus, these are the time stamps of
329 previous computations that are used by process p , and thus, the quantities $n - \tau_{q,n,i}^{(p)}$
330 are sometimes called *delays*. We use the notation

$$331 \quad X_p = \mathbb{R}^{N_p}, \quad \text{and} \quad \tilde{X} = X_1 \times \dots \times X_P$$

333 to denote local and global solution spaces, and $\mathcal{T}_{p,n} : \tilde{X} \rightarrow X_p$ the rule that is used
334 to update the local iterate $u_{p,n}$ at iteration n . We can now define, for each process p ,

335 the asynchronous iterations as follows:

$$336 \quad (2.8) \quad u_{p,n} = \begin{cases} \mathcal{T}_{p,n} \left(u_{1,n}^{(p)}, \dots, u_{P,n}^{(p)} \right) & \text{if } p \in \sigma_n, \\ u_{p,n-1} & \text{if } p \notin \sigma_n. \end{cases}$$

338 The iteration is initialized using some initial guess for $u_{p,0}$, and we used the notation
 339 $u_{q,n}^{(p)} := u_{q,\tau_q^{(p)}(n)}$ to denote the data from process q that is available to process p at
 340 time n .

341 In other words, at time n , either $u_{p,\bullet}$ is not updated (if $p \notin \sigma_n$) or it is updated
 342 with the result of applying the (local) operator $\mathcal{T}_{p,n}$ to the variables computed at
 343 times $\tau_{\bullet}^{(p)}$. For comparison, the corresponding synchronous iteration is given by

$$344 \quad (2.9) \quad u_{p,n} = \mathcal{T}_{p,n} (u_{1,n-1}, \dots, u_{P,n-1}),$$

346 or, in compact form, as

$$347 \quad (2.10) \quad \tilde{u}_n = \tilde{\mathcal{T}}_n (\tilde{u}_{n-1}),$$

349 where

$$350 \quad \tilde{u}_n = (u_{1,n}, \dots, u_{P,n}), \quad \text{and} \quad \tilde{\mathcal{T}}_n = (\mathcal{T}_{1,n}, \mathcal{T}_{2,n}, \dots, \mathcal{T}_{P,n}).$$

352 We further assume that the three following conditions are satisfied

$$353 \quad (2.11) \quad \forall p, q \in \{1, \dots, P\}, 1 \leq i \leq N_q, \forall n \in \mathbb{N}^*, \left(\tau_{q,n}^{(p)} \right)_i \leq n,$$

$$354 \quad (2.12) \quad \forall p \in \{1, \dots, P\}, \text{card} \{n \in \mathbb{N}^* \mid p \in \sigma_n\} = \infty,$$

$$355 \quad (2.13) \quad \forall p, q \in \{1, \dots, P\}, 1 \leq i \leq N_q, \lim_{n \rightarrow +\infty} \left(\tau_{q,n}^{(p)} \right)_i = \infty.$$

356 Condition (2.11) indicates that data used at the time n must have been produced
 357 before time n , i.e., time does not flow backward. Condition (2.12) means that no
 358 process will ever stop updating its components. Condition (2.13) corresponds to the
 359 fact that new data will always be provided to the process. In other words, no process
 360 will have a piece of data that is never updated.

361 We note that these assumptions pose no significant restrictions on the iterations
 362 that we consider, but are necessary for the analysis.

363 Assume that each X_p is a normed linear space, equipped with a norm $\|\cdot\|_p$. Given
 364 a positive vector $w \in \mathbb{R}_{>0}^P$, the weighted norm $\|\cdot\|_w$ on the product space X is defined
 365 to be

$$366 \quad \|\tilde{u}\|_w = \max_{p=1, \dots, P} \frac{\|u_p\|_p}{w_p}.$$

368 We are ready to present a convergence theorem for asynchronous iterative algo-
 369 rithms, whose proof can be found in [18, Theorem 3.3].

370 **THEOREM 2.1.** *Assume that there exists $\tilde{u}^* \in X$ such that $\tilde{\mathcal{T}}_n (\tilde{u}^*) = \tilde{u}^*$ for all n .
 371 Moreover, assume that there exists $\gamma \in [0, 1)$ and $w \in \mathbb{R}_{>0}^P$ such that for all n we have*

$$372 \quad \left\| \tilde{\mathcal{T}}_n (\tilde{u}) - \tilde{u}^* \right\|_w \leq \gamma \|\tilde{u} - \tilde{u}^*\|_w.$$

374 *Then the asynchronous iterates \tilde{u}_n converge to \tilde{u}^* , the unique common fixed point of
 375 all $\tilde{\mathcal{T}}_n$.*

376 In view of equations (2.5) and (2.6), we have

$$377 \quad \mathcal{T}_{p,n}^{1L}(w_1, \dots, w_P) = w_p + \mathbf{A}_p^{-1} \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q w_q)$$

378

379 for the one-level method. We immediately observe that the mappings $\mathcal{T}_{p,\bullet}^{1L}$ do not
 380 depend on n , and that the iteration is stationary.

381 In order to tackle the two-level method, based on Algorithm 2.3 we set

$$382 \quad \mathcal{T}_{0,n}^{2L}(v_0, w_1, \dots, w_P) = \mathbf{A}_0^{-1} \sum_{q=1}^P \mathbf{R}_0 \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q w_q).$$

383

384 Moreover, for $p = 1, \dots, P$, we set

$$385 \quad \mathcal{T}_{p,n}^{2L}(v_0, w_1, \dots, w_P) = w_p + \frac{1}{2} \mathbf{R}_p \mathbf{R}_0^T v_0 + \frac{1}{2} \mathbf{A}_p^{-1} \sum_{q=1}^P \mathbf{R}_p \mathbf{R}_q^T (\mathbf{D}_q \mathbf{R}_q f - \mathbf{A}_q \mathbf{D}_q w_q)$$

386

387 for iteration numbers n that include coarse-grid updates, and

$$388 \quad \mathcal{T}_{p,n}^{2L}(v_0, w_1, \dots, w_P) = \mathcal{T}_{p,n}^{1L}(w_1, \dots, w_P)$$

390 for iterations n without coarse-grid update. Status variables such as `canWriteRHSp`
 391 act implicitly as constraints on the sets of update indices σ_n and do not appear in the
 392 definition of the mappings $\mathcal{T}_{p,n}^{2L}$.

393 It has been shown in [19] that both the one- and the two-level iterations are
 394 contracting in a weighted max-norm, provided that \mathbf{A} is a non-singular M-matrix, i.e.
 395 if \mathbf{A} has nonpositive off-diagonal elements and all entries of \mathbf{A}^{-1} are nonnegative.

396 Thus, we have the following result.

397 **THEOREM 2.2.** *The one-level method given in (2.5) and (2.6) and the two-level*
 398 *method given in Algorithm 2.3 converge, provided that \mathbf{A} is a non-singular M-matrix*
 399 *and that the conditions 2.11–2.13 hold.*

400 For further extensions of the theory, such as inexact sub-solves with \mathbf{A}_p^{-1} replaced
 401 by some (potentially nonstationary) $\mathbf{S}_{p,n} \approx \mathbf{A}_p^{-1}$, we refer the reader to [19].

402 **3. One-sided Message Passing Interface.** In order to drive the asynchronous
 403 method in a distributed memory setting, we use a one-sided approach wherein the re-
 404 mote process incurs minimal overhead for servicing received messages from the sender
 405 process. The one-sided approach is achieved in MPI using the Remote Memory Access
 406 (RMA) semantics, wherein every process exposes a part of its local memory window
 407 to remote processes for read as well as write operations. However, in reality, a syn-
 408 chronization between the source and the target process is required for progress of
 409 the underlying application. This active synchronization step, while still preserving
 410 the asynchronous nature of the algorithm, is expensive and might erode the natu-
 411 ral gains obtained from the asynchronous method. Therefore in order to extract the
 412 maximum gains from an asynchronous method, a passive approach is required. A
 413 passive approach entails transmission of messages which causes little to no interfer-
 414 ence to the target process. As a result, the target process does not need to yield its
 415 operating system time for servicing incoming message interrupts and therefore does
 416 not participate in the communication process. The RMA framework on MPI imple-
 417 ments passive target synchronization with the help of two sets of primitives `MPI_Win_`

418 `lock/MPI_Win_unlock` and `MPI_Win_lock_all/MPI_Win_unlock_all`. While the for-
419 mer involves opening and closing the exposure epoch on remote nodes for each access
420 operation, the latter only requires opening and closing of access epoch once during
421 the application lifetime incurring less target synchronization overhead.

422 RMA's passive one-sided communication can leverage a hardware mechanism
423 known as Remote Direct Memory Access (RDMA) [28] when available. It allows RMA
424 to directly map memory windows to the RDMA engine, allowing messages written by
425 remote processes to be directly read by each process. This leads to minimum distur-
426 bance to the remote process and achieving a truly passive, one-sided communication
427 scheme.

428 RDMA is usually a hardware characteristic that may not be supported by all
429 machines. Though we expect one-sided communication of RMA to be able to handle
430 progress of communication in an entirely asynchronous manner, it generally fails to do
431 so since MPI does not guarantee asynchronous progress. In such a case, asynchronous
432 progress may be enforced by allocating certain auxiliary cores to ghost processes that
433 solely perform the task of asynchronous progress control. As a consequence we obtain
434 an RDMA agnostic system while simultaneously obtaining the benefits of RDMA.
435 Even in the presence of RDMA, asynchronous progress control mechanism can be
436 complementary since the low level RDMA engine may not be capable to handle high
437 volumes of communication. Casper [32] and Intel Asynchronous Progress Control
438 (APC) are two such implementations that provide ghost processes for asynchronous
439 progress control.

440 **4. Implementation and numerical experiments.**

441 **4.1. Comparison of one-sided MPI communication options.** There are a
442 multitude of options for achieving asynchronous neighborhood exchange. Data that is
443 supposed to be moved from rank p to rank q could be held in MPI windows on either
444 p or q . In the first case, rank p will write the data to its local buffer using `MPI_Put`,
445 and rank q will retrieve it from the remote buffer using `MPI_Get`. In the second case,
446 rank p writes the data to the remote memory region using a `MPI_Put`, and q retrieves
447 using a local `MPI_Get`.

448 The second distinction comes from the type of locking mechanism used. Exclusive
449 or shared locks can either be applied for each individual memory access (`MPI_Win_`
450 `lock`), or windows can be locked in shared fashion for all subsequent access (`MPI_`
451 `WIN_lock_all`). In the latter case, windows can be flushed using any of the available
452 flush operations.

453 We benchmark the different available options in a simple test case in order to de-
454 termine which one should be used in the implementation of our domain decomposition
455 solvers. The performance of one-sided MPI communication depends on the support
456 provided by the MPI implementation as well as the network hardware. These experi-
457 ments are performed on the Haswell partition of Cori at the National Energy Research
458 Scientific Computing Center (NERSC), using the default Cray MPICH, version 7.7.3.
459 Since one-sided MPI has not been widely adopted, performance variations compared
460 to the classical two-sided routines can be expected to be much more significant. It
461 should be noted that different network hardware and better support in future MPI
462 versions could further improve timings for one-sided MPI routines.

463 64 MPI ranks are arranged in a three dimensional regular periodic grid (3D torus),
464 and each rank repeatedly exchanges a vector of doubles with its 26 neighbors. This
465 test mimics the communication pattern in the neighborhood exchange of the one-
466 level method. For each of the possible communication option as given in [Table 4.1](#),

467 we measure the time it takes to perform 50,000 exchanges of vectors of 500 doubles.
 468 By exchanging vectors that have a constant value corresponding to the exchange
 469 iteration, we can also measure how often inconsistent data is accessed (i.e. data that
 470 is accessed before it has been completely been transmitted). This phenomenon does
 471 not occur when using two-sided communication, since completion is guaranteed by
 472 the implementation. While the absolute number of accesses to incomplete writes
 473 is probably quite dependent on the ratio of computation to communication, we are
 474 interested in the susceptibility of the different communication options.

475 We make several observations. Unsurprisingly, the use of exclusive locks does not
 476 perform well in terms of time. However, the use of shared locks in every communica-
 477 tion phase performs equally poorly, which is why we decide to use global locking and
 478 unlocking (`MPI_Win_lock_all` / `MPI_Win_unlock_all`) in what follows. Using global
 479 locking, we see that using remote puts instead of remote gets is significantly faster.

480 We also observe that unless exclusive locks are used, we always experience access
 481 to inconsistent data. This might not be of too much importance within our
 482 application, since it amounts to using residual information that is only slightly more
 483 outdated. Finally, we observe that using global locking and puts results in faster
 484 communication than classical two-sided non-blocking communication.

485 Based on the above results, we choose to use global locking using `MPI_Win_lock_`
 486 `all` / `MPI_Win_unlock_all`, paired with remote `MPI_Puts` and local `MPI_Gets` and
 487 `MPI_flush_all`, since it appears to provide a good balance of speed and consistency.
 488 We note however that these results might depend significantly on characteristics of
 489 the system and the MPI implementation.

490 **4.2. Performance metrics.** The average contraction factor per iteration is de-
 491 fined as $\tilde{\rho} = (r_{\text{final}}/r_0)^{\frac{1}{K}}$, where r_0 is the norm of the initial residual vector, r_{final}
 492 the norm of the final residual vector, and K is the number of iterations that were taken
 493 to decrease the residual from r_0 to r_{final} . For an asynchronous method, the number
 494 of iterations varies from subdomain to subdomain, and hence $\tilde{\rho}$ is not well-defined.
 495 The following generalization permits us to compare synchronous methods with their
 496 asynchronous counterpart:

$$497 \quad \hat{\rho} = \left(\frac{r_{\text{final}}}{r_0} \right)^{\frac{\tau_{\text{sync}}}{T}} .$$

499 Here, T is the total iteration time, and τ_{sync} is the average time for a single iteration
 500 of the synchronous method. In the synchronous case, since $T = \tau_{\text{sync}}K$, $\hat{\rho}$ recovers $\tilde{\rho}$.
 501 The approximate contraction factor $\hat{\rho}$ can be interpreted as the average contraction
 502 of the residual norm in the time of a single synchronous iteration. As it will be
 503 visible in the results to follow, we note here that $\hat{\rho}$ for the asynchronous method
 504 obviously depends on the total iteration time for the synchronous method. Assume
 505 that the total iteration time for the synchronous method doubles, but the time taken
 506 by the asynchronous one stays constant. Consequently, the approximate contraction
 507 factor for the synchronous method stays constant, but the contraction factor for the
 508 asynchronous method gets squared and therefore decreases.

509 **4.3. Test problem.** As a test problem, we solve

$$510 \quad -\Delta u = f \quad \text{in } \Omega = [0, 1]^d, \quad u = 0 \quad \text{on } \partial\Omega,$$

512 where the right-hand side is $f = d\pi^2 \prod_{k=1}^d \sin(\pi x_k)$. The corresponding solutions is
 513 $u = \prod_{k=1}^d \sin(\pi x_k)$. We discretize Ω using a uniform simplicial mesh and approximate

global lock	per comm phase	per neighbor	time in seconds	inconsistency fraction
✗	MPI_Win_lock(EXCLUSIVE) MPI_Win_unlock	local MPI_Put, remote MPI_Get	34.6	0.0
✗	MPI_Win_lock(EXCLUSIVE) MPI_Win_unlock	remote MPI_Put, local MPI_Get	37.8	0.0
✗	MPI_Win_lock(SHARED) MPI_Win_unlock	local MPI_Put, remote MPI_Get	31.8	0.00151
✗	MPI_Win_lock(SHARED) MPI_Win_unlock	remote MPI_Put, local MPI_Get	33.0	0.00254
n/a	MPI_Wait_all	MPI_Isend, MPI_Irecv	9.59	0.0
✓	-	local MPI_Put, remote MPI_Get	25.8	0.123
✓	-	remote MPI_Put, local MPI_Get	8.42	0.00716
✓	MPI_flush_all	local MPI_Put, remote MPI_Get	22.1	0.117
✓	MPI_flush_all	remote MPI_Put, local MPI_Get	9.06	0.00491
✓	MPI_flush_local_all	local MPI_Put, remote MPI_Get	22.1	0.099
✓	MPI_flush_local_all	remote MPI_Put, local MPI_Get	9.02	0.00501
✓	MPI_flush_local	local MPI_Put, remote MPI_Get	24.1	0.172
✓	MPI_flush_local	remote MPI_Put, local MPI_Get	10.7	0.00198
✓	MPI_flush	local MPI_Put, remote MPI_Get	21.8	0.105
✓	MPI_flush	remote MPI_Put, local MPI_Get	11.2	0.00207

TABLE 4.1

Results of communication test described in Section 4.1 on 64 MPI ranks. The listed operations are either performed once per neighborhood communication phase, or for each individual neighborhood exchange. If `MPI_Win_lock_all/MPI_Win_unlock_all` is used, the column “global lock” has a ✓. We measured the time for 50,000 repetitions and the fraction of neighborhood exchanges leading to incompletely written data.

514 the solution using piece-wise linear finite elements. We note that the arising system
515 matrix \mathbf{A} is a non-singular M-matrix, and therefore [Theorem 2.2](#) applies. Further-
516 more, we mention that the generalization of the test problem to convection-diffusion
517 problems with non-constant diffusion coefficient is possible, but does not alter the
518 numerical results obtained below in a significant way, which is why we only present
519 the case of the standard Poisson problem.

520 **4.4. Convergence detection.** In classical synchronous iterative methods, a
521 stopping criterion of the form $\|r\| < \varepsilon$ is evaluated at every iteration. Here, r is the
522 residual vector, ε is a prescribed tolerance (that might be chosen as a function of the
523 discretization error), and $\|\cdot\|$ is an appropriate norm. The global quantity $\|r\|$ needs
524 to be computed as the sum of local contributions from all the subdomains. This
525 implies that convergence detection in asynchronous methods is not straightforward,
526 since collective communication primitives require synchronization. In the numerical
527 examples below, we use a simplistic convergence criterion, consisting in writing the
528 local contributions to a master rank, say rank 0. This master rank sums the contri-
529 butions, and determines if this approximation of the global residual norm is smaller
530 than the prescribed tolerance. If so, the master rank declares global convergence and
531 notifies the other ranks by sending a non-blocking message. This simplistic conver-
532 gence detection mechanism has several drawbacks. For one, the global residual is
533 updated by the master rank, which might not happen frequently enough. Hence it
534 is possible that the iteration continues despite the true global residual norm already
535 being smaller than the tolerance. Moreover, the mechanism puts an increased load
536 on the network connection to the master rank, since every subdomain writes to its
537 memory region. Finally, since the local contributions to the residual norm are not
538 necessarily monotonically decreasing, the criterion might actually detect convergence
539 when the true global residual is not yet smaller than the tolerance. The delicate topic
540 of asynchronous convergence detection has been treated in much detail in the liter-
541 ature, and we refer to [\[2, 29\]](#) for an overview of more elaborate approaches. While
542 these detection schemes mostly address the shortcomings of the above approach, their
543 correct implementation turns out to be quite involved. Since we are not observing
544 any major issues with our simplistic convergence detection scheme for the test prob-
545 lems that we consider, we have not implemented any of the schemes available in the
546 literature.

547 **4.5. Platform and implementation details.** All runs are performed on the
548 Haswell partition of the Cori supercomputer at NERSC. While all of the code was
549 written from scratch, the differences between the synchronous and the asynchronous
550 code path are limited, since only the communication layer and the stopping criterion
551 need to be changed. (E.g. compare [Algorithms 2.2](#) and [2.3](#).) We stress that the asyn-
552 chronous solver uses one-sided communication only in the solve phase. Therefore, we
553 record solve times only, since the time to set up the solver is unaffected by the type of
554 communication in the solution phase. Furthermore, all subdomains are synchronized
555 via a `MPI_Barrier` before entering the solve phase. One MPI rank is used per core,
556 i.e. 32 ranks per Haswell node. Moreover, one subdomain is assigned to each MPI
557 rank. The underlying mesh is partitioned either into uniformly sized rectangular sub-
558 domains or using the METIS library [\[25\]](#). In the latter case, the option to minimize
559 the overall communication volume is used. Our solvers handle general unstructured
560 matrices, and the structure of the mesh is not exploited. We either use

- 561 • direct solvers for subdomain and coarse-grid problems, provided by SuperLU
562 [\[27, 15\]](#), or

- conjugate gradient method preconditioned with an incomplete Cholesky factorization for the subdomain problems and a geometric multigrid solver for the coarse problem.

The latter option would allow for a distributed coarse-grid solve and is therefore in principle more scalable. In all numerical examples, we will use only a single core for the coarse-grid solve.

4.6. Comparison against HPDDM. We verify the performance of the synchronous version of our code against the HPDDM library [24, 23, 16] using the 2D and 3D test problems from Section 4.3. In all cases, we set up a GMRES solver and use a two-level additive RAS as right preconditioner. The reason for using a Krylov method here is that domain decomposition methods in general are commonly used as preconditioners, and HPDDM is most likely developed with that use case in mind. HPDDM was linked against the Intel Math Kernel Library, SuiteSparse [12] and ARPACK [26] and the option for coarse grid data exchange using MPI_Igather/MPI_Iscatter was enabled. We use the following parameters in HPDDM: `-hpddm_krylov_method=gmres -hpddm_schwarz_method=ras -hpddm_schwarz_coarse_correction=additive -hpddm_geneo_nu=NU`, where NU is chosen so that the size of the coarse grid matches our solver. In 2D the subdomains consist of roughly 20k unknowns, and the coarse grid contains about 16 unknowns per subdomain. In 3D the subdomains consist of roughly 40k unknowns, and the coarse grid contains about 1 unknown per subdomain. In Figures 4.1 and 4.2 we plot the results of weak scaling experiments: overall solve time, the reached residual norm and the time per iteration. We repeated each run 5 times. Mean values are given by solid lines, and individual runs as dots. We observe that while the time to convergence behaves quite differently for both implementations, the time per iteration follows the same trend. HPDDM behaves slightly better at large subdomain count which could be explained by the use of MPI collectives, but might also be an artifact of the difference in convergence behavior or the difference in coarse solvers (HPDDM uses Cholmod). We can therefore use our synchronous method as a base of comparison for the newly developed asynchronous solver.

4.7. One-level RAS, 2D test problem, strong scaling. We compare synchronous and asynchronous one-level RAS in a strong scaling experiment, where we fix the global problem size of a 2D test problem to about 261,000 unknowns, and vary the number of subdomains between 4 and 256. We cannot expect good scaling behavior for this one level method, since increasing the number of subdomains adversely affects the rate of convergence. The iteration is terminated based on the simplistic convergence criterion described in Section 4.4. In Figure 4.3 we show solve time, final residual norm and approximate rate of convergence. It can be observed that the synchronous method is faster for smaller numbers of subdomains, yet comparatively slower for larger number of subdomains. The crossover point between the two regimes appears to be at 64 subdomains.

An important question is whether the asynchronous method happens to converge because every subdomain performs the same number of local iterations, and hence the asynchronous method just mirrors the synchronous one, merely with a different communication method. The histogram in Figure 4.4 shows that this is not the case. The number of local iterations varies significantly. The slowest subdomain performs barely more than 11,000 iterations, whereas the fastest one almost reaches 16,000. The problem was load balanced by the number of degrees of freedom in each subdomain, thus the local solves are also approximately balanced but the communication is likely slightly imbalanced. This means that in this scenario, system and network noise

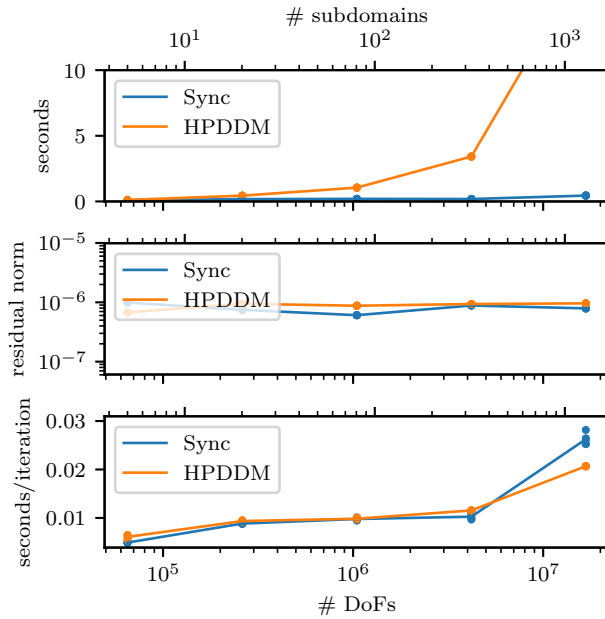


FIG. 4.1. Weak scaling of GMRES preconditioned by two-level additive RAS using the synchronous version of our code (Sync) and HPDDM for the 2D test problem in a load balanced case. From top to bottom: Total solution time, final residual norm, and time per iteration. Mean values are given by solid lines, and individual runs as dots.

612 are the main contributions to the observed variations in local iteration counts. For
 613 comparison, when only 4 subdomains were used, the local iteration counts were 1497,
 614 1500, 1504 and 1527.

615 The advantage of asynchronous RAS becomes even clearer when the experiment
 616 is repeated under load imbalance. We create an artificial load imbalance by choosing
 617 one of the subdomains to be 50% larger than the rest. In Figure 4.5 it is observed that
 618 the asynchronous method outperforms the synchronous one in all but the 4 subdomain
 619 case.

620 **4.8. Two-level RAS, 2D test problem.** In order to gauge the performance
 621 and scalability of the synchronous and asynchronous two-level RAS solvers, we per-
 622 form weak and strong scaling experiments.

623 **4.8.1. Weak scaling.** In the weak scaling experiment the number of subdomains
 624 P and the global number of degrees of freedom (DoFs) are increased proportionally.
 625 We use 16, 64, 256 and 1024 subdomains to solve the 2D test problem. The local
 626 number of unknowns on each subdomain is kept constant at almost 20,000. The
 627 coarse-grid problem increases in size proportionally to the number of subdomains,
 628 with approximately 16 unknowns per subdomain. Again, the iteration is terminated
 629 based on the simplistic convergence criterion described in Section 4.4.

630 In Figure 4.6 we plot the solution time, the achieved residual norm and the average
 631 contraction factor $\hat{\rho}$ depending on the global problem size. Both the synchronous and
 632 the asynchronous method reach the prescribed tolerance of 10^{-8} . Due to the lack
 633 of an efficient mechanism of convergence detection, the asynchronous method ends
 634 up iterating longer than necessary, so that the final residual norm often is smaller

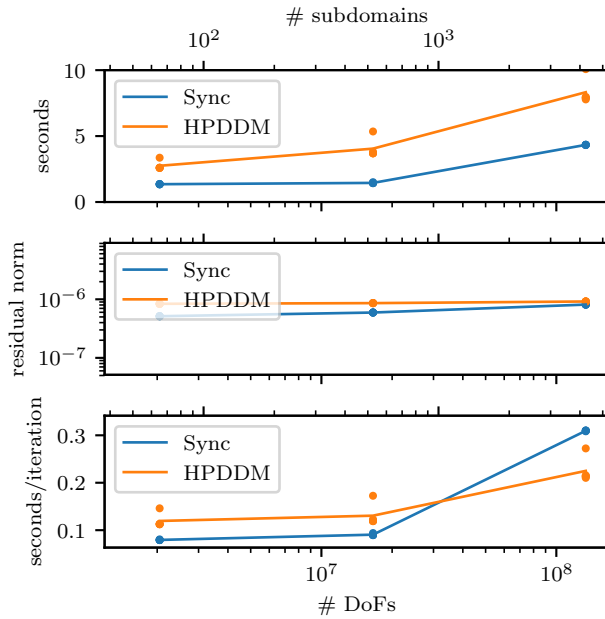


FIG. 4.2. Weak scaling of GMRES preconditioned by two-level additive RAS using the synchronous version of our code (Sync) and HPDDM for the 3D test problem in a load balanced case. From top to bottom: Total solution time, final residual norm, and time per iteration. Mean values are given by solid lines, and individual runs as dots.

635 than 10^{-9} . The number of iterations in the synchronous case is about 110, whereas
 636 the number of local iterations in the asynchronous case varies between 110 and 150.
 637 (See Figure 4.7.) The iteration counts are significantly lower than for the one-level
 638 methods. One can observe that for 16, 64 and 256 subdomains, the asynchronous
 639 and the synchronous methods take almost the same time for the solve. For 1024
 640 subdomains, however, the synchronous method is seen to take significantly more time.
 641 This can be explained by the fact that for 1024 subdomains, the size of the coarse grid
 642 is comparable to the size of the subdomains, and hence the coarse-grid solve which
 643 exchanges information with all the subdomains slows down the overall progress. For
 644 the asynchronous case this is not observed, since the subdomains do not have to
 645 wait for information from the coarse grid. This explains why we see better weak
 646 scalability for the asynchronous method than for the synchronous variant, and why
 647 we can observe a speedup of 2x of the asynchronous method over its synchronous
 648 counterpart. The third subplot of Figure 4.6 shows that the asynchronous method
 649 outperforms its synchronous equivalent in all but the smallest problem.

650 To further illustrate the effect of load imbalance, we repeat the previous experi-
 651 ment with one subdomain being 50% larger than the rest. The results are shown in
 652 Figure 4.8. While the results are mostly consistent with the previous case, it can be
 653 seen that, as expected, the performance advantage of the asynchronous method over
 654 the synchronous one has increased. Even before the size of the coarse-grid system is
 655 comparable to the size of the typical subdomain problem, the asynchronous method
 656 outperforms its synchronous counterpart.

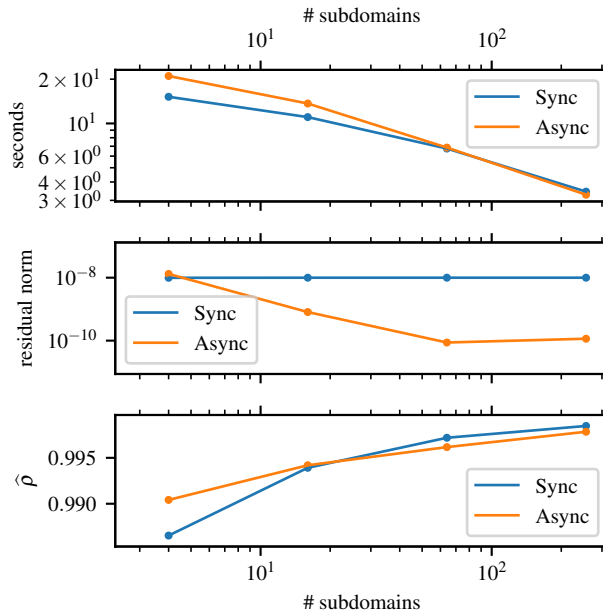


FIG. 4.3. Strong scaling of synchronous and asynchronous one-level RAS for the 2D test problem with system size of approximately 261,000 unknowns. The subdomains are load balanced. From top to bottom: Solution time, final residual norm, and the resulting approximate contraction factor $\hat{\rho}$. It can be observed that the synchronous method is significantly faster than for smaller numbers of subdomains (cores), yet comparatively slower for larger number of subdomains, as shown by the contraction factor.

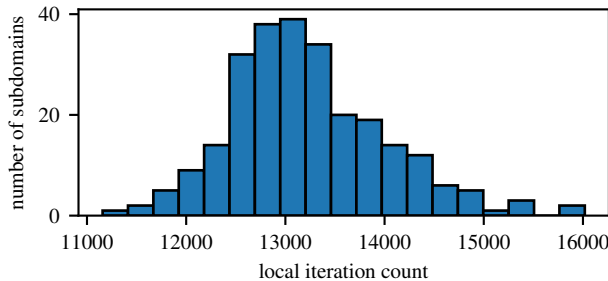


FIG. 4.4. Histogram of local iteration counts for asynchronous one-level RAS for the 2D test problem with 256 subdomains (load balanced case).

657 **4.8.2. Strong scaling.** For the strong scaling experiment the global number of
 658 degrees of freedom used to discretize the 2D test problem is fixed at about 4 million.
 659 The coarse-grid problem consists of approximately 4,000 unknowns. The number of
 660 subdomains used on the fine level takes values in $\{4, 16, 64, 256\}$. This means that
 661 the coarse-grid problem is always smaller than the typical subdomain problem, and
 662 no slowdown due to an imbalance of the computational cost of coarse and fine solve
 663 should arise.

664 The timing results are shown in the top of Figure 4.9. Both synchronous and
 665 asynchronous method display good strong scaling behavior. It is observed that the
 666 synchronous method is faster than the asynchronous method for smaller subdomain
 667 count. But already for 64 subdomains this behavior is reversed, and the asynchronous
 668 method outperforms the synchronous one. This suggests that synchronization is an

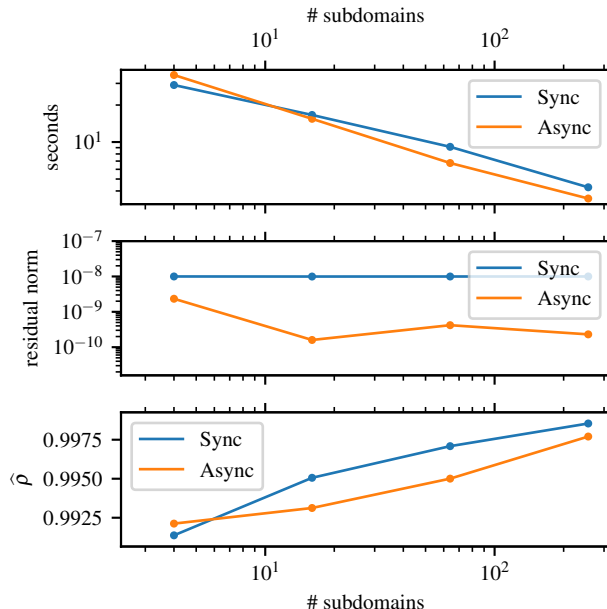


FIG. 4.5. Strong scaling of synchronous and asynchronous one-level RAS for the 2D test problem with a system size of approximately 261,000 unknowns under load imbalance: one subdomain is 50% larger than the rest. From top to bottom: Solution time, final residual norm, and approximate contraction factor $\hat{\rho}$. It can be observed that the asynchronous method outperforms the synchronous one in all but the 4 subdomain case, as shown by the contraction factor. The advantage of the asynchronous method over the synchronous one is increased, as compared to Figure 4.3.

669 important factor already at modest core count.

670 At the bottom of Figure 4.9, we show the timing results in the case of load
 671 imbalance. It can be seen that the asynchronous method is faster than the synchronous
 672 one independent of the number of subdomains, and that its performance advantage
 673 increases as more processes are used.

674 **4.9. Two-level RAS with iterative sub-solves, 3D test problem.** The
 675 density of the subdomain matrices A_p in 3D (about 15 entries per row) is higher than
 676 for the 2D test problem (about 7 entries per row). This means that direct factorization
 677 leads to more fill-in and thereby is more expensive. Therefore, we solve subdomain
 678 and coarse problem of the three dimensional test case using iterative solvers. For
 679 the subdomains, we use a conjugate gradient solver preconditioned by an incomplete
 680 Cholesky factorization. We employ a relative tolerance of 1/10 which has been de-
 681 termined experimentally to be sufficient. The coarse-grid problem is solved using a
 682 single V-cycle of a geometric multigrid solver with one step of Gauss-Seidel for pre-
 683 and post-smoothing. The use of multigrid allows us to solve the coarse-grid problem
 684 in a distributed fashion when it becomes too large for a single MPI rank. The global
 685 problem is partitioned into uniformly sized regular subdomains. The local number of
 686 unknowns on each subdomain is kept constant at about 40,000. The coarse-grid prob-
 687 lem increases in size proportionally to the number of subdomains, with approximately
 688 one unknown per subdomain.

689 The results of a weak scaling experiment are shown in Figure 4.10. We observe
 690 behavior that is similar to the 2D case. We notice however that the size of the coarse-
 691 grid problem (and hence the solution of the coarse-grid problem) are not the issue

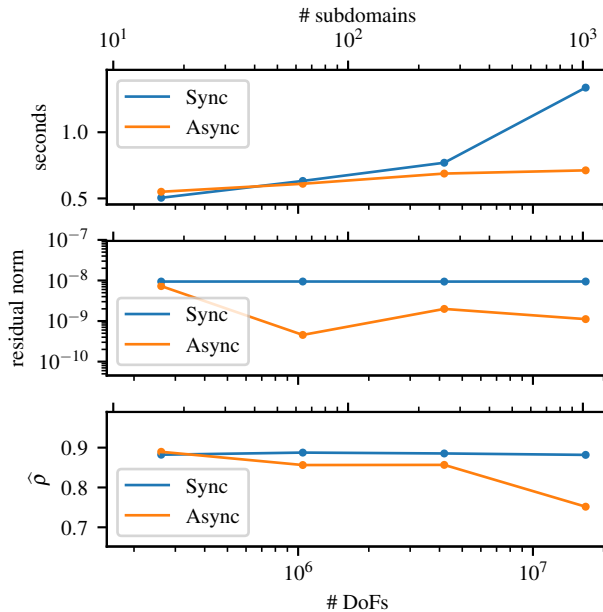


FIG. 4.6. Weak scaling of synchronous and asynchronous two-level additive RAS for the 2D test problem, load balanced case. From top to bottom: Total solution time, final residual norm, and approximate contraction factor $\hat{\rho}$. One can observe that for 16, 64 and 256 subdomains, the asynchronous and the synchronous method take almost the same time for the solve, with a slight advantage for the asynchronous method. For 1024 subdomains, however, the synchronous method is seen to take significantly more time, since the coarse grid, due to its size, starts to be the limiting factor. The asynchronous method is not affected by this.

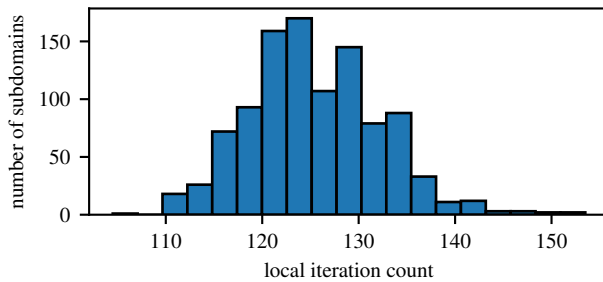


FIG. 4.7. Histogram of local iteration counts asynchronous two-level additive RAS for the 2D test problem with 1024 subdomains.

692 here. At 4096 ranks, the coarse-grid problem is an order of magnitude smaller than
 693 the typical subdomain problem. The apparent slowdown of the synchronous method
 694 is caused by the cost of exchanging information between the coarse grid and the
 695 subdomains. While a slowdown is also visible in the asynchronous method, it is much
 696 less pronounced, resulting in a speedup of 4x over the synchronous method.

697 We also observe that compared to the 2D case where direct solvers were used, both
 698 synchronous and asynchronous iterations terminate almost exactly once the prescribed
 699 tolerance has been achieved. The reason for this is that convergence checks occur much
 700 more frequently as the tolerance is reached, since the iterative sub-solves converge to
 701 their local tolerance typically within one iteration.

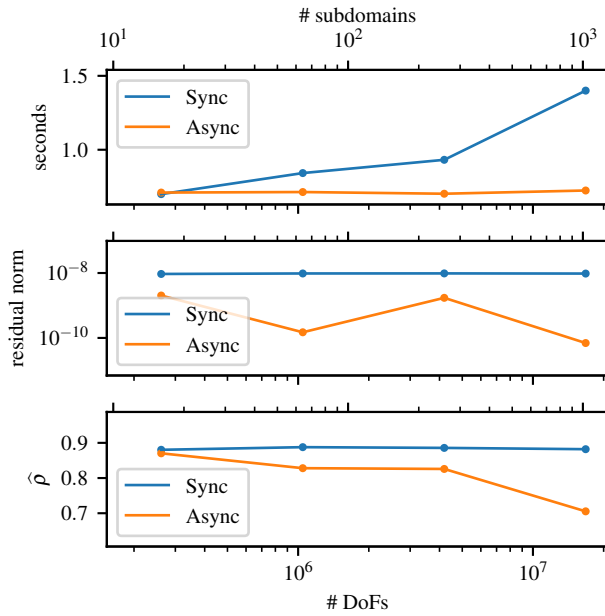


FIG. 4.8. Weak scaling of synchronous and asynchronous two-level additive RAS for the 2D test problem under load imbalance: one subdomain is 50% larger than all the other ones. From top to bottom: Total solution time, final residual norm, and approximate contraction factor $\hat{\rho}$. The advantage of the asynchronous method over the synchronous one is increased, as compared to Figure 4.6.

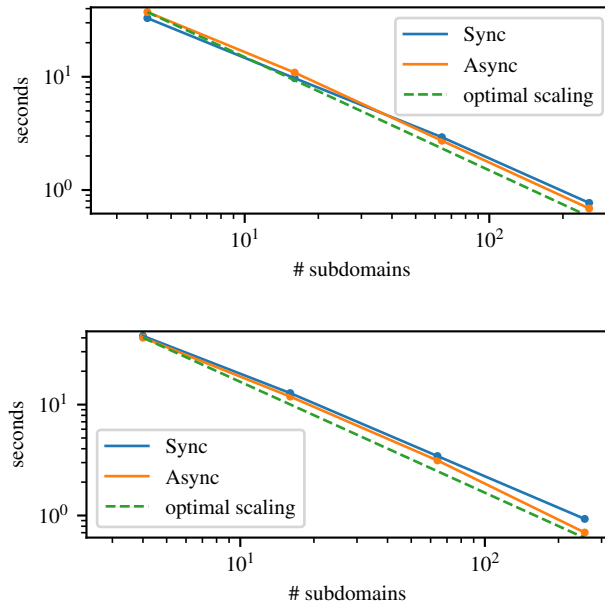


FIG. 4.9. Strong scaling of synchronous and asynchronous two-level additive RAS for the 2D test problem. On top: load balanced subdomains. At the bottom: load imbalance, one subdomain is 50% larger than the others.

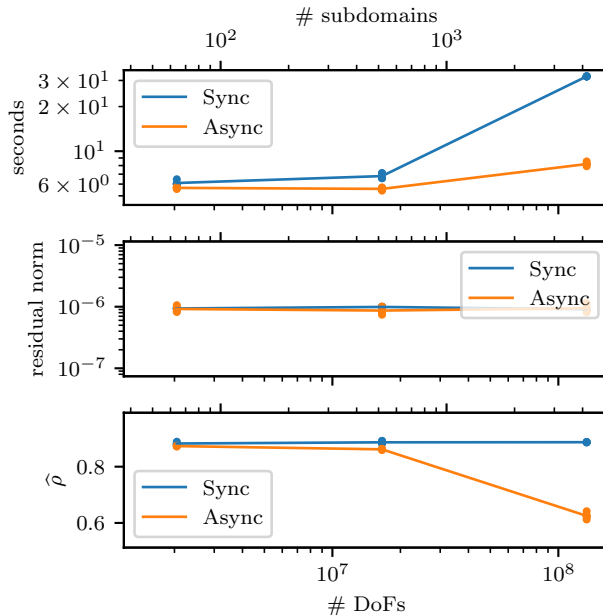


FIG. 4.10. Weak scaling of synchronous and asynchronous two-level additive RAS, load balanced 3D case. From top to bottom: Total solution time, final residual norm, and approximate contraction factor $\hat{\rho}$. One can observe that for 64 and 512 subdomains, the asynchronous and the synchronous method take almost the same time for the solve, with a slight advantage for the asynchronous method. For 4096 subdomains, however, the synchronous method is seen to take significantly more time. The reason for this is not the solution of the coarse-grid problem, as in 2D, but the cost of the data exchange. The effect on the asynchronous method is much less pronounced.

702 **5. Conclusion.** In the present work, we have explored the use of asynchronous
703 alternatives to conventional (synchronous) one-level and two-level domain decom-
704 position solvers. To the best of our knowledge, we proposed the first truly asyn-
705 chronous two-level method, where each processor can do different number of updates
706 (iterations). Several options to achieve asynchronous communication were tested,
707 and we found that our use case benefited most from using `MPI_Win_lock_all` /
708 `MPI_Win_unlock_all`, remote `MPI_Puts` and local `MPI_Gets`. The numerical results
709 presented demonstrate that asynchronous iterations can be considered a viable alter-
710 native to synchronous methods, despite partial availability of information from neigh-
711 bors. Asynchronous methods seem to be beneficial already at modest core count, even
712 for load balanced scenarios. In the presence of load imbalance, their performance ad-
713 vantage becomes even clearer, and we observed speedups up to 4x. While we focused
714 our attention on a particular Schwarz method, it is of inherent interest to explore asyn-
715 chronous variants of other, potentially more effective domain decomposition methods
716 involving deflation or non-overlapping decompositions (such as FETI and BDDC) or
717 more than two levels. The presented inclusion of a novel asynchronous coarse-grid
718 correction paves the way for asynchronous methods to be used in extremely scalable
719 parallel solvers.

720 **Acknowledgment.** Sandia National Laboratories is a multimission laboratory
721 managed and operated by National Technology and Engineering Solutions of Sand-
722 dia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S.
723 Department of Energy's National Nuclear Security Administration under contract

724 DE-NA0003525.

725 This paper describes objective technical results and analysis. Any subjective
726 views or opinions that might be expressed in the paper do not necessarily represent
727 the views of the U.S. Department of Energy or the United States Government.

728 SAND Number: SAND2020-4049 O

729 This material is based upon work supported by the U.S. Department of Energy,
730 Office of Science, Office of Advanced Scientific Computing Research, Applied Mathe-
731 matics program under Award Numbers DE-SC-0016564. This research used resources
732 of the National Energy Research Scientific Computing Center, a DOE Office of Science
733 User Facility supported by the Office of Science of the U.S. Department of Energy
734 under Contract No. DE-AC02-05CH11231.

735

REFERENCES

- 736 [1] V. AHLGREN, S. ANDERSSON, J. M. BRANDT, N. CARDO, S. CHUNDURI, P. FIELDS, A. C. GEN-
737 TILE, R. GERBER, J. GREENSEID, A. GREINER, ET AL., *Cray system monitoring: Successes*
738 *requirements and priorities.*, tech. report, Sandia National Lab.(SNL-NM), Albuquerque,
739 NM (United States); Sandia . . . , 2018.
- 740 [2] J. M. BAHJ, S. CONTASSOT-VIVIER, R. COUTURIER, AND F. VERNIER, *A decentralized conver-*
741 *gence detection algorithm for asynchronous parallel iterative algorithms*, IEEE Transac-
742 tions on Parallel and Distributed Systems, 16 (2005), pp. 4–13.
- 743 [3] S. BALAY, S. ABHYANKAR, M. F. ADAMS, J. BROWN, P. BRUNE, K. BUSCHELMAN, L. DALCIN,
744 A. DENER, V. EIJKHOUT, W. D. GROPP, D. KARPEYEV, D. KAUSHIK, M. G. KNEPLEY,
745 D. A. MAY, L. C. MCINNES, R. T. MILLS, T. MUNSON, K. RUPP, P. SANAN, B. F. SMITH,
746 S. ZAMPINI, H. ZHANG, AND H. ZHANG, *PETSc users manual*, Tech. Report ANL-95/11 -
747 Revision 3.13, Argonne National Laboratory, 2020, <https://www.mcs.anl.gov/petsc>.
- 748 [4] G. M. BAUDET, *Asynchronous iterative methods for multiprocessors*, Journal of the ACM
749 (JACM), 25 (1978), pp. 226–244.
- 750 [5] D. P. BERTSEKAS, *Distributed asynchronous computation of fixed points*, Mathematical Pro-
751 gramming, 27 (1983), pp. 107–120.
- 752 [6] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Parallel and distributed computation: numerical meth-*
753 *ods*, vol. 23, Prentice Hall Englewood Cliffs, NJ, 1989.
- 754 [7] X.-C. CAI, M. DRYJA, AND M. SARKIS, *Restricted additive Schwarz preconditioners with har-*
755 *monic overlap for symmetric positive definite linear systems*, SIAM Journal on Numerical
756 Analysis, 41 (2003), pp. 1209–1231.
- 757 [8] X.-C. CAI AND M. SARKIS, *A restricted additive Schwarz preconditioner for general sparse*
758 *linear systems*, SIAM Journal on Scientific Computing, 21 (1999), pp. 792–797.
- 759 [9] F. CAPPELLO, A. GEIST, B. GROPP, L. KALE, B. KRAMER, AND M. SNIR, *Toward exascale*
760 *resilience*, International Journal of High Performance Computing Applications, 23 (2009),
761 pp. 374–388, <https://doi.org/10.1177/1094342009347767>.
- 762 [10] F. CAPPELLO, A. GEIST, W. GROPP, S. KALE, B. KRAMER, AND M. SNIR, *Toward exascale*
763 *resilience: 2014 update*, Supercomputing frontiers and innovations, 1 (2014), pp. 5–28,
764 <https://doi.org/10.14529/jsfi140101>.
- 765 [11] D. CHAZAN AND W. MIRANKER, *Chaotic relaxation*, Linear algebra and its applications, 2
766 (1969), pp. 199–222.
- 767 [12] Y. CHEN, T. A. DAVIS, W. W. HAGER, AND S. RAJAMANICKAM, *Algorithm 887: Cholmod,*
768 *supernodal sparse cholesky factorization and update/downdate*, ACM Trans. Math. Softw.,
769 35 (2008), <https://doi.org/10.1145/1391989.1391995>, <https://doi.org/10.1145/1391989.1391995>.
- 770 [13] A. T. CHRONOPOULOS AND C. W. GEAR, *On the efficient implementation of preconditioned*
771 *s-step conjugate gradient methods on multiprocessors with memory hierarchy*, Parallel com-
772 puting, 11 (1989), pp. 37–53.
- 773 [14] A. T. CHRONOPOULOS AND C. W. GEAR, *s-step iterative methods for symmetric linear systems*,
774 J. Comput. Appl. Math., 25 (1989), pp. 153–168, [https://doi.org/http://dx.doi.org/10.1016/0377-0427\(89\)90045-9](https://doi.org/http://dx.doi.org/10.1016/0377-0427(89)90045-9).
- 775 [15] J. W. DEMMEL, S. C. EISENSTAT, J. R. GILBERT, X. S. LI, AND J. W. H. LIU, *A supernodal*
776 *approach to sparse partial pivoting*, SIAM J. Matrix Analysis and Applications, 20 (1999),
777 pp. 720–755.
- 778 [16] V. DOLEAN, P. JOLIVET, AND F. NATAF, *An introduction to domain decomposition methods*,

780

- 781 Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2015, [http://](http://dx.doi.org/10.1137/1.9781611974065.ch1)
782 dx.doi.org/10.1137/1.9781611974065.ch1. Algorithms, theory, and parallel implementation.
- 783 [17] M. EL HADDAD, J. C. GARAY, F. MAGOULÈS, AND D. B. SZYLD, *Synchronous and Asyn-*
784 *chronous optimized Schwarz Methods for one-way subdivision of bounded domains*, Nu-
785 *merical Linear Algebra and Applications*, 27 (2020), p. e2279. 30 pages.
- 786 [18] A. FROMMER AND D. B. SZYLD, *On asynchronous iterations*, *Journal of Computational and*
787 *Applied Mathematics*, 123 (2000), pp. 201–216.
- 788 [19] A. FROMMER AND D. B. SZYLD, *An algebraic convergence theory for restricted additive Schwarz*
789 *methods using weighted max norms*, *SIAM Journal on Numerical Analysis*, 39 (2001),
790 pp. 463–479.
- 791 [20] J. C. GARAY, F. MAGOULÈS, AND D. B. SZYLD, *Synchronous and asynchronous optimized*
792 *Schwarz method for Poisson’s equation in rectangular domains*, Tech. Report 17-10-18,
793 Department of Mathematics, Temple University, Oct. 2017. Revised April 2018.
- 794 [21] J. C. GARAY, F. MAGOULÈS, AND D. B. SZYLD, *Convergence of asynchronous optimized*
795 *Schwarz methods in the plane*, in *Domain Decomposition Methods in Science and En-*
796 *gineering XXIV*, P. E. B. stard, S. C. Brenner, L. Halpern, H. H. Kim, R. Kornhuber,
797 T. Rahman, and O. B. Widlund, eds., *Lecture Notes in Computer Science and Engineer-*
798 *ing*, Berlin and Heidelberg, 2018, Springer, pp. 333–341.
- 799 [22] P. GHYSELS, T. J. ASHBY, K. MEERBERGEN, AND W. VANROOSE, *Hiding global communica-*
800 *tion latency in the GMRES algorithm on massively parallel machines*, *SIAM Journal on*
801 *Scientific Computing*, 35 (2013), pp. C48–C71.
- 802 [23] P. JOLIVET, F. HECHT, F. NATAF, AND C. PRUD’HOMME, *Scalable domain decomposition pre-*
803 *conditioners for heterogeneous elliptic problems*, in *Proceedings of the International Con-*
804 *ference on High Performance Computing, Networking, Storage and Analysis, SC ’13*, New
805 *York, NY, USA, 2013*, Association for Computing Machinery, [https://doi.org/10.1145/](https://doi.org/10.1145/2503210.2503212)
806 [2503210.2503212](https://doi.org/10.1145/2503210.2503212), <https://doi.org/10.1145/2503210.2503212>.
- 807 [24] P. JOLIVET AND F. NATAF, *HPDDM – high-performance unified framework for domain decom-*
808 *position methods*, 2020, <https://github.com/hpddm/hpddm> (accessed March 27, 2020).
- 809 [25] G. KARYPIS AND V. KUMAR, *A Fast and High Quality Multilevel Scheme for Partitioning*
810 *Irregular Graphs*, *SIAM Journal on Scientific Computing*, 20 (1998), pp. 359–392, <https://doi.org/10.1137/S1064827595287997>.
- 811 [26] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK users’ guide: solution of large-scale*
812 *eigenvalue problems with implicitly restarted Arnoldi methods*, vol. 6, Siam, 1998.
- 813 [27] X. LI, J. DEMMEL, J. GILBERT, I. L. GRIGORI, M. SHAO, AND I. YAMAZAKI, *SuperLU Users’*
814 *Guide*, Tech. Report LBNL-44289, Lawrence Berkeley National Laboratory, September
815 1999. <http://crd.lbl.gov/~xiaoye/SuperLU/>. Last update: August 2011.
- 816 [28] J. LIU, J. WU, AND D. K. PANDA, *High performance RDMA-based MPI implementation over*
817 *InfiniBand*, *International Journal of Parallel Programming*, 32 (2004), pp. 167–198.
- 818 [29] F. MAGOULÈS AND G. GBIKPI-BENISSAN, *Distributed convergence detection based on global*
819 *residual error under asynchronous iterations*, *IEEE Transactions on Parallel and Dis-*
820 *tributed Systems*, (2017).
- 821 [30] F. MAGOULÈS, D. B. SZYLD, AND C. VENET, *Asynchronous optimized Schwarz methods with*
822 *and without overlap*, *Numerische Mathematik*, (2017), pp. 1–29, [https://doi.org/10.1007/](https://doi.org/10.1007/s00211-017-0872-z)
823 [s00211-017-0872-z](https://doi.org/10.1007/s00211-017-0872-z).
- 824 [31] Z. PENG, Y. XU, M. YAN, AND W. YIN, *A Rock: An Algorithmic Framework for Asynchronous*
825 *Parallel Coordinate Updates*, *SIAM Journal on Scientific Computing*, 38 (2016), pp. A2851–
826 A2879, <https://doi.org/10.1137/15M1024950>.
- 827 [32] M. SI, A. J. PENA, J. HAMMOND, P. BALAJI, M. TAKAGI, AND Y. ISHIKAWA, *Casper: An asyn-*
828 *chronous progress model for MPI RMA on many-core architectures*, in *Parallel and Dis-*
829 *tributed Processing Symposium (IPDPS), 2015 IEEE International*, IEEE, 2015, pp. 665–
830 676.
- 831 [33] B. SMITH, P. BJORSTAD, AND W. GROPP, *Domain decomposition: parallel multilevel methods*
832 *for elliptic partial differential equations*, Cambridge University Press, 2004.
- 833 [34] A. TOSELLI AND O. WIDLUND, *Domain decomposition methods: algorithms and theory*, vol. 34,
834 Springer Science & Business Media, 2006.
- 835 [35] J. WOLFSON-POU AND E. CHOW, *Asynchronous multigrid methods*, in *33rd IEEE International*
836 *Parallel and Distributed Processing Symposium (IPDPS), IEEE Computer Society, 2019*,
837 pp. 101–110.
- 838 [36] I. YAMAZAKI, E. CHOW, A. BOUTELLER, AND J. DONGARRA, *Performance of Asynchronous*
839 *Optimized Schwarz with One-sided Communication*, *Parallel Computing*, 86 (2019), pp. 66–
840 81.
- 841 [37] I. YAMAZAKI, S. RAJAMANICKAM, E. G. BOMAN, M. HOEMMEN, M. A. HEROUX, AND S. TO-
- 842

843 MOV, *Domain decomposition preconditioners for communication-avoiding Krylov methods*
844 *on a hybrid CPU/GPU cluster*, in Proceedings of the International Conference for High
845 Performance Computing, Networking, Storage and Analysis, IEEE Press, 2014, pp. 933–
846 944.