

Petal Tool for Analyzing and Transforming Legacy MPI Applications

Hadia Ahmed¹(✉), Anthony Skjellum², and Peter Pirkelbauer¹

¹ University of Alabama at Birmingham, Birmingham, AL 35294, USA
{hadia,pirkelbauer}@uab.edu

² Auburn University, Auburn, AL 36830, USA
skjellum@auburn.edu

Abstract. Legacy MPI applications are an important and economically valuable category of parallel software that rely on the MPI-1, MPI-2 (and, more recently, MPI-3) standards to achieve performance and portability. Many of these applications have been developed or ported to MPI over the past two decades, with the implicit (dual) goal of achieving acceptably high performance and scalability, and a high level of portability between diverse parallel architectures. However they were often created implicitly using MPI in ways that exploited how a particular underlying MPI behaved at the time (such as those with polling progress and poor implementation of some operations). Thus, they did not necessarily take advantage of the full potential for describing latent concurrency or for loosening the coupling of the application thread from the message scheduling and transfer.

This paper presents a first transformation tool, Petal, that identifies calls to legacy MPI primitives. Petal is implemented on top of the ROSE source-to-source infrastructure and automates the analysis and transformation of existing codes to utilize non-blocking MPI and persistent MPI primitives. We use control flow and pointer alias analysis to overlap communication and computation. The transformed code is capable of supporting better application bypass, yielding better overlapping of communication, computation, and I/O. We present the design of the tool and its evaluation on available benchmarks.

1 Introduction

The Message Passing Interface (MPI) describes a library that enables the development of portable parallel software for large-scale systems. The first MPI standard [12] focused on providing a basic framework for point-to-point and collective communication. MPI-2 [8] introduced one-sided communication, added support for parallel file access, and dynamic process management, and extended the usefulness of two-group (inter-communicator) operations. MPI offers a small set of core functions that are sufficient for the development of many applications, and also offers functionality that helps experts optimize applications [10].

MPI bindings exist for C++, Fortran, and many other languages, making MPI one of the most prevalent programming models for high-performance computing. MPI is supported on many platforms, which makes applications developed with MPI portable to many large-scale systems. Building high-performance computing systems constitutes a large investment in human resources. As the communication infrastructure advances and the MPI standards and library implementations follow suite, legacy codes becomes a potential liability. Code that does not utilize more recent MPI primitives will not scale well on newer architectures. This effect will become more marked over time.

With Exascale systems on the horizon, the cost of communication is becoming a major concern. Compared to older architectures, communication incurs relatively more overhead. Legacy software written for older architectures often utilizes `MPI_Send` and `MPI_Recv` for the communication of point-to-point messages. These two primitives block until the data exchange completes (or at least till the send buffer can be reused by the calling thread). While this makes it easy for programmers to reason about communication, such methods fail to utilize computing resources efficiently. On next generation hardware, the implied cost of sending data using a polling and/or blocking mode of communication significantly rises and it is expected that software relying on blocking communication will have too much overhead. In order to take advantage of the architectural changes in Exascale, existing code needs to be transformed to use better primitives, some of which are only available in MPI-3 or higher. Non-blocking primitives allow overlap of communication with local computation¹. A paired, non-blocking communication uses two MPI routines, one to start (`MPI_Isend`, `MPI_Irecv`) and one to complete (`MPI_wait`). After a communication has been initiated, code can compute, and only waits at the `MPI_wait` to synchronize with the communication operation. In addition to the benefits of non-blocking, applications that exhibit fixed point-to-point communication patterns can further utilize persistent operations introduced in MPI-1 and being extended in MPI-3.x. Persistent MPI primitives reduce communication overhead in applications that exhibit fixed patterns. Persistent MPI operations minimize the overhead incurred from redundant message setup.

Rewriting legacy MPI programs by hand is both tedious and error prone. To relieve programmers of the task of manually rewriting applications, the authors have developed tool support to replace uses of MPI primitives that are known to perform slowly on modern hardware (or may have better alternatives, especially on next-generation architectures) with better alternatives in the MPI standard. We have implemented a source code rejuvenation tool [16] called Petal using the ROSE source-to-source infrastructure [3, 17]. We chose ROSE for its support of many languages relevant for high-performance computing. Petal analyzes existing source code and finds calls to `MPI_Send` and `MPI_Recv`. It replaces these primitives with their non-blocking counterparts and uses data-dependency

¹ Provided the underlying MPI does not poll excessively to make progress or for message completion, the messages are long enough, and there is sufficient memory bandwidth for both communication and computation.

and control-flow information to find code locations where corresponding calls to `MPI.Wait` need to be inserted. If Petal can determine that the communication partners, message buffer, and message length do not change, persistent communication primitives will be used in lieu of non-persistent functions.

Overall, this paper offers the following contributions:

- program analysis and transformation to replace blocking MPI calls with non-blocking calls;
- program analysis and transformation to introduce persistent MPI calls; and,
- analysis of persistent MPI implementations.

The remainder of this paper is organized as follows. Section 2 presents more detailed information on MPI and ROSE. Section 3 describes our implementations and Sect. 4 discusses our evaluation and findings. Section 5 gives an overview of related work on MPI transformations, and Sect. 6 offers conclusions and an outlook on possible future work.

2 Background

This section provides background information on MPI and the ROSE compiler infrastructure.

2.1 MPI Primitives

MPI offers several modes of operation for point-to-point communication. Many programs employ `MPI.Send` and `MPI.Recv`, two blocking MPI primitives. `MPI.Send` takes the following arguments: base pointer to message data, the number of elements to send, a type descriptor, the destination, and a communicator. The base pointer to data typically points to a send buffer, but it could also point to data described by a type descriptor. Blocking means that the MPI primitive waits until the message buffer containing the data being sent/received is safe to be used again by the calling process. Only then is control returned to the caller. On send, actual implementations of `MPI.Send` may either block until all data has been transmitted or copy the data to an intermediate internal buffer. The use of blocking primitives may be prone to deadlocks, if programmers do not carefully consider send and receive order [13]

`MPI.Isend` and `MPI.Irecv` are non-blocking versions for point-to-point message communication. Compared to `MPI.Send`'s arguments, `MPI.Isend` adds an additional argument for a request handle. The handle is used in calls to `MPI.Wait` to identify which send to wait for. Non-blocking calls return immediately after initiating the communication and the user thread can execute more operations, eventually followed by a completion operation (a wait or test) on the request. The communication is considered complete after a successful call to `MPI.Wait` (or `MPI.Test`, etc.). Non-blocking is used to help promote overlap communication and computation, resulting in communicating cost hiding and yielding overall

better performance on systems that support it. To avoid tampering with the data, programmers must ensure that the message data is not modified before the communication is completed.

Another mode is offered by persistent communication primitives. If a program exhibits regular communication patterns (static arguments), where the same communication partners exchange fixed size messages, utilization of persistent MPI enables exploitation of faster communication paths. Provided MPI implementations efficiently implement these operations, persistence supports reduced overhead by eliminating cost associated with repeated operations and streamlined processing of derived datatypes. Persistence also can reduce jitter and allow for preplanned choice of algorithms, such as for MPI collectives. Since persistence in MPI offers many benefits (potential and long observed), it is likely that future MPI standards will enhance support for persistent primitives, for example by supporting variable length messages between the same communication partners.

Note that all three modes can be used interchangeably. It is possible that one side uses persistent MPI, while the other side does not. That is why the functions are sometimes referred to as providing half-channels.

Figure 1 shows the use of blocking, non-blocking, and persistent operations for a simple 1D heat transfer code. The basic design of the heat-transfer code is depicted in Fig. 1d. The code uses two arrays, containing cells with temperature information. The initial temperatures are located in the even array. In odd numbered timesteps the odd array is computed from the even array and in even numbered timesteps vice versa. Red cells are computed by neighbors and dark blue cells are needed by neighbors for the next iteration. Figure 1a shows a blocking implementation. The order of sends and receives is important to avoid deadlock. Even-numbered MPI processes send first, odd numbered processes receive first. `D` stands for `MPI_DOUBLE`, and `n` is the rank of this node. For simplicity, the codes assume that each process has two neighbors and ignores send and receive status. Figure 1b demonstrates the overlap of communication and computation in non-blocking mode. The key idea is that the inner (light blue) cells can be computed before the data from neighbors are received. The code starts two receive operations to receive both neighbor's data from the last iteration. Then it starts two send operations to communicate its values from the previous iteration to its neighbors. While the communication is ongoing, the inner cells are computed. Before cells depending on neighbors' data can be computed, the code waits until the data have been received (Line 10). After computing the outer cells, the wait in Line 13 blocks until the data have been sent. This is necessary in order not to overwrite the data in the next iteration. Figure 1c shows the persistent version of the code. Since the communication patterns, buffer, and buffer size do not change, we can set up the communication for sends and receives at the beginning of the program, and reuse this pattern in every iteration.

```

1  double b[4]; // send/receive buffer
3  for (int i = 0; i<MAX; ++i) {
4      data_to_buf(prev, b+2);
5      if (n%2 == 0) {
6          MPI_Recv(b+0, 1, D, n-1, 0, com);
7          MPI_Recv(b+1, 1, D, n+1, 0, com);
8      }
9      MPI_Send(b+2, 1, D, n-1, 0, com);
10     MPI_Send(b+3, 1, D, n+1, 0, com);
11     if (n%2 == 1) {
12         MPI_Recv(b+0, 1, D, n-1, 0, com);
13         MPI_Recv(b+1, 1, D, n+1, 0, com);
14     }
15     buf_to_data(b, prev);
16     compute_all(prev, curr);
17     swap(curr, prev);
18 }

```

(a) Blocking operations

```

1  MPI_request r[4]; // request handler
2  double b[4]; // send/receive buffer
3
4  MPI_Recv_init(b+0, 1, D, n-1, 0, com, r+0);
5  MPI_Recv_init(b+1, 1, D, n+1, 0, com, r+1);
6  MPI_Send_init(b+2, 1, D, n-1, 0, com, r+2);
7  MPI_Send_init(b+3, 1, D, n+1, 0, com, r+3);
8  for (int i = 0; i<MAX; ++i) {
9      data_to_buf(prev, b+2);
10     for (int j = 0; j < 4; ++j)
11         MPI_Start(r+j);
12     compute_inner(prev, curr);
13     MPI_Wait(2, r+0, IGNORE);
14     buf_to_data(b, prev);
15     compute_outer(prev, curr);
16     MPI_Wait(2, r+2, IGNORE);
17     swap(curr, prev);
18 }

```

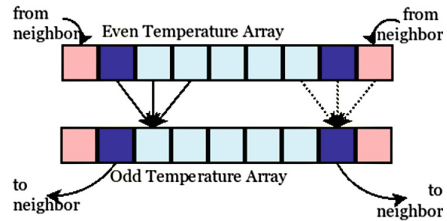
(c) Persistent operations

```

1  MPI_request r[4]; // request handler
2  double b[4]; // send/receive buffer
3
4  for (int i = 0; i<MAX; ++i) {
5      data_to_buf(prev, b+2);
6      MPI_Irecv(b+0, 1, D, n-1, 0, com, r+0);
7      MPI_Irecv(b+1, 1, D, n+1, 0, com, r+1);
8      MPI_Isend(b+2, 1, D, n-1, 0, com, r+2);
9      MPI_Isend(b+3, 1, D, n+1, 0, com, r+3);
10     compute_inner(prev, curr);
11     MPI_Wait(2, req+0, IGNORE);
12     buf_to_data(b, prev);
13     compute_outer(prev, curr);
14     MPI_Wait(2, req+2, IGNORE);
15     swap(curr, prev);
16 }

```

(b) Non-blocking operations



(d) Design Overview

Fig. 1. 1D heat transfer

2.2 The ROSE Compiler Infrastructure

The ROSE source-to-source translation infrastructure is under active development currently at the Lawrence Livermore National Laboratory (LLNL). ROSE provides front ends for many languages, including C/C++, Fortran 77/95/2003, Java, and UPC. ROSE also supports several parallel extensions, such as OpenMP and CUDA. ROSE generates an Abstract Syntax Tree (AST) for the source code. The ASTs are uniformly built for all input languages. ROSE offers many specific analyses (e.g., pointer alias analysis) and makes these available through an API. Users can write their own analyses by utilizing frameworks that ROSE provides. These include attribute evaluation traversals, call graph analysis, control flow graphs, class hierarchies, SSA representation, and dataflow analysis. The Fuse framework [4], is an object-oriented dataflow analysis framework that affords users with the ability to create their own inter- and intra-procedural dataflow analyses by implementing standard dataflow components. ROSE has been used for building custom tools for static analysis, program optimization, arbitrary

program transformation, domain-specific optimizations, performance analysis, and cyber-security. With the representation of the code as an AST and using the static analysis provided from the ROSE libraries, one can explore the code and determine how to improve it by looking for certain code style, inserting new code, changing and/or removing old code, hence generating modified source code while preserving the semantics of the original code.

3 Implementation

In this section, we describe Petal’s implementation of a mechanism to transform applications from using blocking MPI point-to-point routines to using non-blocking versions. We also describe the analysis and transformations to introduce persistent routines.

3.1 Design

Petal transforms code to use non-blocking MPI operations to reveal a better potential overlap of computation and communication and adds persistent operations, whenever possible, to eliminate much of the overhead of repeatedly communicating with a partner node.

Figure 2 shows an overview of our transformation framework. The tool takes MPI source files, for which ROSE compiles and generates the Abstract Syntax Tree (AST), then function calls are inlined if the function implementation should be available. Once inlined, ROSE’s query and builder libraries are used to find and replace blocking with non-blocking calls and to identify where to insert corresponding calls to `MPI.Wait`. If some or all of these non-blocking calls are used repeatedly with the same arguments, they are replaced with persistent communication operations. At the end, Petal generates a new transformed source file as its output, using either non-blocking or persistent communications (which are always non-blocking).

The idea of following this approach is based on trying to maximize the overlap between communication and computation without compromising the semantics of the original application. Inlining eliminates the need to use inter-procedural analysis and simplifies moving `MPI.Wait` downward, crossing its original function boundaries if no unsafe access to the message buffer is found across the function calls. MPI uses pointers to the message buffers that they use in their communication. This fact allowed us to simplify the analysis used by the tool and focus only on using pointer alias analysis. ROSE’s pointer alias analysis implements Steensgaard’s algorithm, which has linear time complexity [19]. This allows our tool to scale well with large applications.

3.2 Blocking to Non-blocking Transformation

Petal allows changing the blocking function call `MPI.Send/MPI.Recv` to the corresponding `MPI.Isend/MPI.Irecv` while ensuring proper access to the message

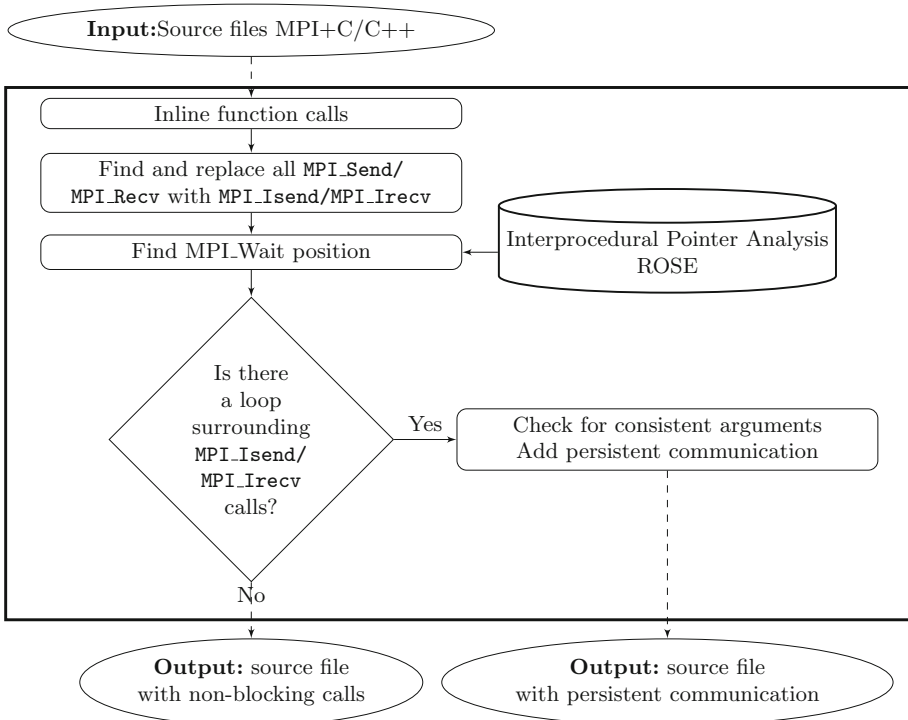


Fig. 2. Transformation framework

buffers, and once an operation that access the message buffer is encountered, `MPI.Wait` is inserted before it to ensure the safety of the data.

Calling `MPI.Send/MPI.Recv` is in effect the same as calling `MPI.Isend/MPI.Irecv` immediately followed by `MPI.Wait`. Our tool moves calls to `MPI.Wait` downward along forward control flow edges as long as the operations are safe with respect to the MPI operation and buffer access. Any write to a message buffer that is used in a send operation, and any access to a message buffer that is used in a receive operation is considered an unsafe access and `MPI.Wait` must be called before that to maintain the correctness of the code.

For each blocking call, to be replaced by the corresponding non-blocking, three variables are created, two of which are handlers for `MPI.Request` and `MPI.Status` plus a flag introduced to ensure the execution of `MPI.Wait` if and only if its corresponding non-blocking call is executed. Each blocking call is replaced with the corresponding `MPI.Isend/MPI.Irecv`. After finding and replacing blocking calls, control flow analysis is used to find subsequent statements, extract the variables used in these statements and use pointer analysis to test for aliasing between the message buffer used and the variables in hand. For the send operation, we identify potential update operations, such as a variable occurring on the left hand side of an assignment. We use pointer alias analysis to check whether an update could

modify some data. For the receive operation, all expressions that read values from a variable are tested. Variable extraction includes subscripts of an array, arguments in non-inlined function calls, variables used in conditions of control statements, initial and increment statements of for loop, and operands of binary and unary operations. Our tool uses ROSE's pointer alias analysis to test whether the extracted variables and the communication buffer could alias. If there could be an alias, the tool inserts the corresponding `MPI_Wait` before the statement using this variable.

Because of inlining, Petal is able to bypass the end of the function and keep searching for potential usage of the message buffer outside the function containing the original MPI calls. If no alias is found in all the statements following the block call, the tool identifies where this statement is located. If it is in `main()`, that means that no alias is found and the `MPI_Wait` is inserted before the `MPI_Finalize`. Because of the complexity of loop-carried data dependencies, currently the tool does not support moving `MPI_Wait` outside the loop body. Hence, if it is in a loop statement (for, while, do-while) `MPI_Wait` is inserted as the last statement in the loop. Otherwise the statement following the block that has the blocking call is examined for alias analysis. To ensure that the `MPI_Wait` in its new position gets executed only if its corresponding non-blocking call is executed, a flag is set to true with each non-blocking call and then based on its value, the corresponding `MPI_Wait` is executed.

Figure 3 shows an example of a snippet of code before and after transformation. Figure 3a shows the original blocking code and Fig. 3b shows how the code looks after the transformation. Lines 3–5 shows the declaration of the `MPI_Request`, `MPI_Status` and the flag variables. Line 10 sets the flag to 1 where Line 21 tests for the flag's value before executing the `MPI_Wait` on Line 22. Since this is a send call, the `printf` function call is a safe read access and the wait call is inserted after it.

3.3 Non-persistent to Persistent Transformation

If a program exhibits regular communication patterns, where the same communication partners exchange fixed size messages, utilization of persistent MPI enables exploitation of faster communication paths². In Shao et al. [18] work to identify communication patterns for MPI programs, they discovered that many programs that are considered dynamic can use persistent communication. This means that changing these programs to use persistence will result in better performance. The difficulty of persistent communications is that possible uses in real world codes are hard to determine statically. To overcome this limitation, we use dynamic analysis. Petal transforms code to persistent mode and inserts guards that test that the arguments did not change. Persistent communication is a four-step process. First, a persistent request is created. Then, data transmission is initiated. After that, wait routines must be called to ensure proper completion. Lastly, the persistent request handlers must be explicitly deallocated.

² At least on high quality implementations of MPI.


```

1  int *buffer;
   int x;
3  ... //code for main,initialization,...
5  for(int i=0;i<1000;i++)
6  {
7      if (myid == source) {
8          *buffer = 123;
9          MPI_Send(buffer,count,MPI_INT,
10             dest,tag,MPI_COMM_WORLD);
11             x = 0;
12         }
13         else {
14             *buffer = 456;
15             x = 1;
16         }
17         printf("%d\n",*buffer);
18     }
19 }
20 }

```

(a) Before

```

   int *buffer;
2  int x;
   MPI_Request reqs[1];
   MPI_Status stats[1];
   int flags[1];
3  ... //code for main,initialization,...
4  for(int i=0;i<1000;i++)
5  {
6      if (myid == source) {
7          flags[0]=1;
8          *buffer = 123;
9          MPI_Isend(buffer,count,MPI_INT,
10             dest,tag,MPI_COMM_WORLD,&reqs[0]);
11             x = 0;
12         }
13         else {
14             *buffer = 456;
15             x = 1;
16         }
17         printf("%d\n",*buffer);
18         if (flags[0] == 1)
19             MPI_Start(&reqs[0]);
20     }
21 }

```

(b) After

Fig. 3. Non-blocking transformation example

Changing to persistent mode is best suited for non-blocking calls in a loop. Petal does such transformations from non-blocking non-persistent to persistent automatically. A structure is created to hold initial values for non-blocking call arguments as its members. Using ROSE queries, the tool identifies `MPI_Isend/MPI_Irecv` and checks to see which one is enclosed in a loop. If no call is in a loop, no transformations are performed. If one or more are found inside a loop, the tool initiates a persistent request with the same arguments as the corresponding non-blocking call and places this initiation process before the loop (`MPI_Send/Recv_Init`). In addition, it stores the values of the `MPI_Isend/MPI_Irecv` arguments in a struct variable for comparing the values across iterations. Then inside the loop, it inserts an if statement to check if the current values are the same as the persistent request argument values, if the outcome is yes, it uses this persistent request using `MPI_Start(&request)`, otherwise it uses the normal `MPI_Isend/MPI_Irecv` call. After the loop, all the created persistent requests are freed.

Following the output from Figs. 3b and 4 shows the result of applying the persistence transformation. On the left side, line 6 shows the persistent request handler and line 7–16 shows the struct definition and its instance declaration. Line 20 initiates the persistent communication passing it all the non-blocking arguments and lines 23–29 represents the copying of the arguments values to the struct instance. On the right side, line 6–11 represents the test against the current values with the values stored in the persistent request. If they are the same `MPI_Start` on line 13 is executed, otherwise the original `MPI_Isend` is executed on line 16–17. Line 29 shows the deallocation of the persistent request.

```

1  int *buffer;
2  int x;
3  MPI_Request reqs[1];
4  MPI_Status stats[1];
5  int flags[1];
6  MPI_Request preqs[1];
7  struct buf_data
8  {
9      void *buf;
10     int count;
11     MPI_Datatype datatype;
12     int dest;
13     int tag;
14     MPI_Comm comm;
15 }
16 struct buf_data temp_data[1];
17 ... //code for main,initialization,...
18 MPI_Send_init(buffer, count, MPI_INT,
19 dest, tag, MPI_COMM_WORLD, &preqs[0]);
20
21     temp_data[0] . buf = buffer;
22     temp_data[0] . count = count;
23     temp_data[0] . datatype = MPI_INT;
24     temp_data[0] . dest = dest;
25     temp_data[0] . tag = tag;
26     temp_data[0] . comm =
27     MPI_COMM_WORLD;
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

Fig. 4. Persistent transformation example

3.4 Discussion

Even though the tool can detect any unsafe access to the message buffers correctly, the applied analysis has limitations in two cases. First, it treats any access to a part of the array as an access to the whole array. For example if MPI sends the first 10 elements of a 100-element array, an assignment to the 20th element will be considered unsafe even though it is in a different place and can be safely used. The second case is that Steensgaard algorithm treats a struct member access as an access to the whole struct [19]. These two cases might lead to placing the `MPI_Wait` in overly conservative positions in some applications. We plan to improve our tool to handles these cases better, since identifying these cases could result into achieving better communication-computation overlap.

Currently, Petal cannot combine multiple consecutive calls to `MPI_Wait`, if found together, into a single `MPI_Waitall` call. This is because different calls to `MPI_Isend`/`MPI_Irecv` may originate in alternative blocks. For example, two calls are part of the then and else branch of an if statement. We hope to find a better solution instead of using flags and if-statement, to ensure the semantics of the code and being able to take advantage of using `MPI_Waitall`.

4 Evaluation

In this section, we present the preliminary evaluation of using Petal and the effect of its transformations on overall application performance. The experiments were

performed on the TACC Stampede system. Stampede is a 10 Petaflop (PF) Dell Linux Cluster with 6400+ Dell PowerEdge server nodes each with 32 GB memory, 2 Intel Xeon E5 (8-core Sandy Bridge) processors and an additional Intel Xeon Phi Coprocessor (61-core Knights Corner) (MIC Architecture) [20]. We used the mvapich2 MPI library. Petal was tested with the 1D heat decomposition described earlier, 2D heat [7] and DT from the NAS NPB 3.3 benchmark [1].

We tested the performance of the application while varying the number of MPI processes. For 1D heat, we varied the number of MPI processes in each case ranging from 6 to 200 tasks. For 2D heat and DT with classes W and A, the number of MPI processors varied between 16 and 256. Figure 5 shows the execution time speedup ($S = T_{\text{original}}/T_{\text{transformed}}$) after applying non-blocking transformation, and adding persistent communication. Figure 5a shows the effect when running applications with only 16 MPI processes, while Fig. 5b shows the effect on applications with 200 and more processes. As shown in the figures, we experienced good improvement with larger number of processes while flat to minor slowdown was observed with fewer numbers of processes. However, in both cases we experienced minor slowdown when adding persistence³.

4.1 Discussion of Results

Petal was able successfully to transform applications from blocking to non-blocking while pushing `MPI.Wait` as far as possible, while also preserving the correctness of the code output. The results shows that with smaller programs and few tasks, the non-blocking improvement is negligible and sometimes hurts

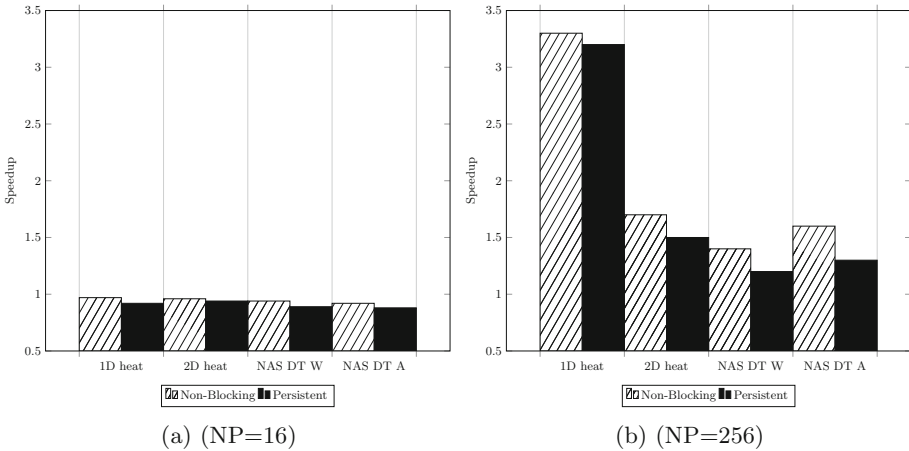


Fig. 5. Execution time speedup

³ This indicates that mvapich may not optimize the code path for persistent send and/or receive.

the application performance. However, with increasing problem size and number of MPI tasks, non-blocking enhanced the performance by up to 30 %.

Unfortunately, even though Petal was able to transform code to persistent mode, the results of persistent performance showed a flat improvement and sometimes a slowdown.

To gain more insight into the usage of persistent communications, we applied the persistent transformation on the LULESH code from LLNL [11] on Stampede and on a Debian 7.6 amd64 computer with 1 Xeon E5410 @ 2.33 GHz using the Open MPI 1.6.5 library. LULESH already exploits non-blocking operations. Since it has some communications that are fixed for most of the program's execution time, persistent communication should be beneficial. However, upon transforming to persistent no gain was seen and with increasing number of tasks we saw a minor slowdown. Since Open MPI is open source, we investigated how it implements its non-blocking and persistent communications. We found that they optimize the code by creating persistent requests and using them whenever possible. Hence, changing the applications' code to persistent will not give a speedup as Open MPI already uses similar optimization techniques internally. The slowdown might be because of the overhead of checking the arguments on each iteration.

According to the MPI Forum [2], persistent requests are considered to be half-channels, which makes the connection faster by reducing the overhead of communication processing within each of the sender and receiver. Our results suggest that the performance improvement is dependent not only on the standard definition of how code should work but it also depends on the actual MPI implementation and architecture. While the tested systems did not show any performance improvements, the transformation may be beneficial on other systems.

5 Related Work

The idea of overlapping communication and computation code is of interest to many researchers because of the promising results in better performance it can give when applied efficiently. In this section, we describe previous research work done to produce overlapped communication and computation in MPI.

Several methods were studied and implemented to handle the communication computation overlap approach. Das et al. [6] represents the closest work to our tool in which they developed an algorithm for pushing wait downward in a segment of code. However, they use Static Single Assignment (SSA) use-def analysis to determine the statements that access the message buffer. Even though they describe a method for moving a `MPI_Wait` out of its current scope interval possibility of going to the parent, they did not implement their method and currently their compiler tool only detects MPI calls and finds `MPI_Wait`'s final position; however, insertion is done by hand. Haque et al. [9] developed a similar tool for transforming blocking to non-blocking; however, it does not use any compiler analysis techniques and relies heavily on the programmer annotation to identify where to move the corresponding non-blocking call and its corresponding wait.

Another work is presented by Nguyen et al. in [14] in which they developed Bamboo, a transformation system that transforms MPI C code into a data-driven application that overlaps computation and communication. It was implemented with the ROSE compiler framework and runtime support using the Tarragon runtime library. Their approach is to determine task precedence. It relies on programmer annotations to mark parallel loops and data packing/unpacking plus calls to communication routines. Other approaches were developed using different techniques to achieve the same goal of maximizing communication and computation overlap. Danalis et al. developed the ASPhALT tool [5] within Open64. Their idea is based on automatically detecting where data is available and applying the pre-pushing transformation to send data as soon as possible. They focused on specific a type of applications that does its communication in two parts where at first, it computes the data in a loop with minimum dependencies across iterations, and then uses communication call(s) after the loop to exchange the data generated by the loop. Pellegrini et al. [15] offer a different approach in which they use the polyhedral model to determine exact dependencies and automatically detect potential overlap on a finer grain. To simplify the analysis, they normalize the code by changing non-blocking to blocking. Their work is limited by polyhedral model requirements of using only affine expressions.

Even though MPI included persistent communication since MPI-1 and these calls emphasize the benefits of using persistent, to our knowledge, no available work offers a tool that automatically transforms non-persistent to persistent communication, when such patterns can be identified.

6 Conclusions and Future Work

In this paper, we described our development of Petal, a tool that supports transforming a blocking MPI code to non-blocking version and introduces persistent communication if possible. We have described the approach used in order to push `MPI.Wait` as far as possible from the corresponding communication call in order to improve the potential for overlap of communication and computation code and also to use persistent communication whenever two points communicate the same type and amount of data over multiple iterations. Petal is based on the ROSE framework and uses ROSE's alias analysis to apply transformation required and to preserve correctness of the code. Preliminary results showed that we can improve performance by using non-blocking. In some cases we found that persistent communication does not improve performance even with code that is proved to have fixed communication for most of the execution time. It does not only depend on having fixed arguments but the MPI library used has an effect too. Further detailed analyses of persistent performance on different architectures with different libraries will be explored.

In addition to analyzing data dependency within loop iterations and moving `MPI.Wait` outside the loop body, if no dependency found, techniques to eliminate loop-carried dependencies on send and receive buffers and perhaps unrolling

the loops will also be explored. This will provide another opportunity to move `MPI_Wait(s)` outside loops boundaries. Another future step is to work on cases where we have 3-D data models and to explore how they can be safely overlapped in communication.

We are also extending the Petal tool to do other automatic translation and/or refactoring that will allow a smooth transition for legacy MPI systems to Exascale systems, such as the use of one-sided communications and changing further to use non-blocking and persistent collective operations (being proposed at present in MPI-3.x).

Acknowledgements. This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1053575.

References

1. The NAS parallel benchmarks. <https://www.nas.nasa.gov/publications/npb.html>
2. Persistent MPI communication. <http://www.mpi-forum.org/docs/mpi-1.1/mpi-11-html/node51.html>. Accessed on 28 June 2015
3. The ROSE source-to-source compiler. <http://rosecompiler.org>
4. Aananthakrishnan, S., Bronevetsky, G., Gopalakrishnan, G.: Hybrid approach for data-flow analysis of mpi programs. In: Proceedings of the 27th International ACM Conference on International Conference on Supercomputing, ICS 2013, pp. 455–456. ACM, New York (2013)
5. Danalis, A., Pollock, L., Swamy, M.: Automatic MPI application transformation with asphalt. In: Parallel and Distributed Processing Symposium, IPDpPS 2007, pp. 1–8, IEEE International, March 2007
6. Das, D., Gupta, M., Ravindran, R., Shivani, W., Sivakeshava, P., Uppal, R.: Compiler-controlled extraction of computation-communication overlap in MPI applications. In: IEEE International Symposium on Parallel and Distributed Processing, IPDpPS 2008, pp. 1–8, April 2008
7. Frexus: mpi-2d-plate (2013). <http://project.github.com>. Accessed on 1 September 2015
8. Gropp, W., Lusk, E., Thakur, R.: Using MPI-2: Advanced Features of the Message-Passing Interface. MIT Press, Cambridge (1999)
9. Haque, M., Yi, Q., Dinan, J., Balaji, P.: Enhancing performance portability of MPI applications through annotation-based transformations. In: 2013 42nd International Conference on Parallel Processing (ICPP), pp. 631–640, October 2013
10. Hoefler, T.: New and old features in MPI-3.0: The past, the standard, and the future, April 2012
11. Karlin, I., Keasler, J., Neely, R.: Lulesh 2.0 updates and changes. Technical report LLNL-TR-641973, August 2013
12. Forum, Message Passing Interface: MPI: A message-passing interface standard. Technical report, Knoxville (1994)
13. Message Passing Interface Forum: MPI: A message-passing interface standard version 3.1, June 2015

14. Nguyen, T., Cicotti, P., Bylaska, E., Quinlan, D., Baden, S.B.: Bamboo: translating MPI applications to a latency-tolerant, data-driven form. In: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC 2012, pp. 39:1–39:11. IEEE Computer Society Press, Los Alamitos (2012)
15. Pellegrini, S., Hoefler, T., Fahringer, T.: Exact dependence analysis for increased communication overlap. In: Träff, J.L., Benkner, S., Dongarra, J.J. (eds.) EuroMPI 2012. LNCS, vol. 7490, pp. 89–99. Springer, Heidelberg (2012)
16. Pirkelbauer, P., Dechev, D., Stroustrup, B.: Source code rejuvenation is not refactoring. In: van Leeuwen, J., Muscholl, A., Peleg, D., Pokorný, J., Rumpe, B. (eds.) SOFSEM 2010. LNCS, vol. 5901, pp. 639–650. Springer, Heidelberg (2010)
17. Schordan, M., Quinlan, D.: A source-to-source architecture for user-defined optimizations. In: Böszörményi, L., Schojer, P. (eds.) JMLC 2003. LNCS, vol. 2789, pp. 214–223. Springer, Heidelberg (2003)
18. Shao, S., Jones, A., Melhem, R.: A compiler-based communication analysis approach for multiprocessor systems. In: 20th International Parallel and Distributed Processing Symposium, IPDpPS 2006, p. 10, April 2006
19. Steensgaard, B.: Points-to analysis in almost linear time. In: Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 1996, pp. 32–41. ACM, New York (1996)
20. Towns, J., Cockerill, T., Dahan, M., Foster, I., Gathier, K., Grimshaw, A., Hazlewood, V., Lathrop, S., Lifka, D., Peterson, G., Roskies, R., Scott, J., Wilkins-Diehr, N.: XSEDE: accelerating scientific discovery. *Comput. Sci. Eng.* **16**(5), 62–74 (2014)