

Compiler-Assisted Overlapping of Communication and Computation in MPI Applications

Jichi Guo, Qing Yi
 University of Colorado, Colorado Springs
 Colorado Springs, CO, USA
 {jguo2, qyi}@uccs.edu

Jiayuan Meng
 Google Inc.
 Mountain View, CA, USA
 meng.jiayuan@gmail.com

Junchao Zhang, Pavan Balaji
 Argonne National Laboratory
 Lemont, IL, USA
 {jczhang, balaji}@anl.gov

Abstract—The performance of distributed-memory applications, many of which are written in MPI, critically depends on how well the applications can ameliorate the long latency of data movement by overlapping them with ongoing computations, thereby minimizing wait time. This paper aims to enable such overlapping in large MPI applications and presents a framework that uses an analytical performance model and an optimizing compiler to systematically enable the optimizations. In particular, we first generate an analytical performance model of the application execution flow to automatically identify potential communication hot spots that may induce long wait time. Next, for each communication hot spot, we search the execution flow graph to find surrounding loops that include sufficient local computation to overlap with the communication. Then, blocking MPI communications are decoupled into nonblocking operations when necessary, and their surrounding loops are transformed to hide the communication latencies behind local computations. We evaluated our framework using 7 MPI applications from the NPB NAS benchmark suite. Our optimizations can attain 3-88% speedup over the original implementations.

Index Terms—Computer Applications; Computer performance; Parallel machines; Automatic programming

I. INTRODUCTION

As computing platforms migrate to clusters of increasingly larger scale of microprocessors, applications need to manage the distributed memories of the processors via explicit message-passing runtimes, for example, MPI, to attain high performance. The relative latency and bandwidth of the communication network in relation to the compute capacity of processors are often hard to predict a priori and may change dramatically from one system to the next. Even on supercomputers comprising homogeneous nodes, system noise is increasing on each node because of aspects such as power management, deeper memory hierarchies, and sharing of hardware such as caches and network. The “equal work means equal time” paradigm is no longer relevant on most systems, and load imbalance increasingly becomes the common scenario even on applications that are symmetrically structured. Consequently, bulk-synchronous communication, where all processes synchronize frequently, is no longer a valid option for high-performance MPI applications. Application performance is often critically determined by its ability to flexibly overlap communications with local computations, thereby minimizing wait time.

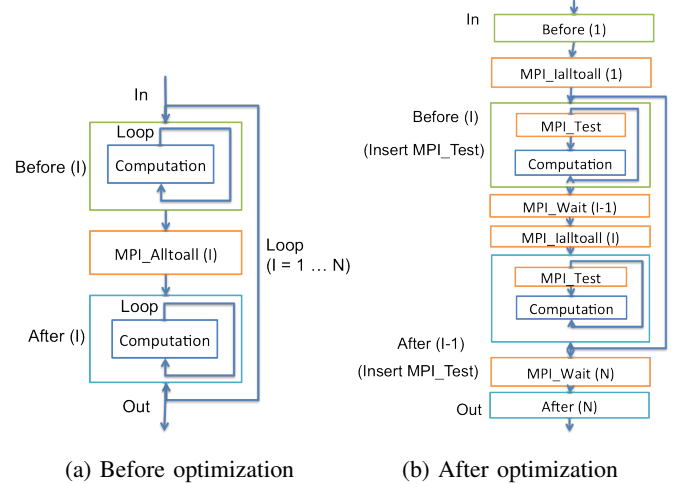


Fig. 1: Optimizing NAS FT by overlapping computation with communication (Before(i)/After(i) are computations before/after the MPI communication of the i th loop iteration)

This paper aims to automatically enable the use of non-blocking latency-hiding techniques to overlap local computation with remote communication in MPI applications, thereby enhancing their overall efficiency and performance portability. To illustrate the optimization, Figure 1a shows the structure of the NAS FT benchmark [4], which applies fast Fourier transform (FFT) to a 3D matrix through a loop that interleaves the computation of scaling the input matrix with a collective communication of MPI_Alltoall to exchange data among the processes. This is then followed by a final transposition of the resulting matrix. The clear separation of computation and communication phases makes the algorithm design easy to implement and maintain. Additionally, the communication buffers can be reused across different loop iterations, saving memory. However, the blocking MPI communication requires that all processes wait while the MPI_Alltoall operation is in progress. Consequently, unless the application is executed on a platform with the fastest network connections, its performance is likely to suffer because of the excessive wait time.

Figure 1b illustrates how the structure in 1a may be modified to better overlap computation with communication. In particular, the MPI_Alltoall operation is decoupled into

two finer-grained operations: a nonblocking `MPI_Ialltoall` and a blocking `MPI_Wait`. Then, the loop is modified so that `Before(i)`, which multiplies a local matrix with a time-evolution array and then saves a transpose of the matrix into a local buffer to be communicated to other processors, and `MPI_Ialltoall(i)`, which exchanges the local transposes among different processes, are essentially moved so that they are evaluated before `MPI_Wait(i-1)`, which waits for the completion of MPI communication of the previous iteration, and `After(i-1)`, which processes the just received remote data (of the $i - 1$ th iteration) and then prints the result into an output file. By using two distinct buffers to store the data used in consecutive MPI communications, the output dependence between `After(i-1)` and `Before(i)/MPI_Ialltoall(i)` can be eliminated, guaranteeing the correctness of optimization.

`MPI_Test` operations then are inserted into the local computation to ensure the progress of the nonblocking communication.¹ By overlapping the MPI communication with the local computation, the transformed code allows the application to perform well even on systems with slow network connections, although nonblocking communications generally take longer than their blocking counterparts do, and more memory may be needed to hold the data during communication.

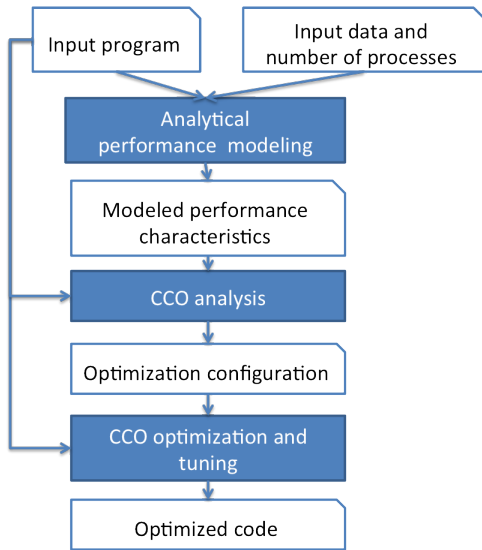


Fig. 2: Optimization workflow.

Figure 2 shows our workflow for systematically enabling communication-computation overlapping (CCO) in MPI applications to enhance their performance portability. The workflow contains three key components: (1) the *performance modeling* component, which analyzes the runtime statistics of an MPI application to extract a *Bayesian execution tree* [16] representation of its execution flow, including the frequencies of various runtime code paths and their performance characteristics such as computation intensities, working set sizes, and

¹Although MPI communications do not need full usage of the CPUs, they need some CPU time, e.g., to manage communication progress, which is supplied only when operations such as `MPI_Test` and `MPI_Wait` are invoked.

communication characteristics of MPI operations; (2) the *CCO analysis* component, which identifies *hot* computation and communication regions and performs profitability and safety analysis to determine whether to optimize these regions; and (3) the *CCO tuning* component, which applies the appropriate program transformations and inserts `MPI_Test` operations with a frequency determined by empirical tuning of the optimized code. For each optimization to be profitable, any communication slowdown from the use of the nonblocking operations must be fully overlapped with the local computation, and the insertion of `MPI_Test` operations should cause only marginal slowdown of the computation. Our framework currently uses empirical tuning to select appropriate optimization configurations and to skip nonprofitable optimizations.

The idea of overlapping computation and communication in MPI applications has been well studied [8], [13], including both using analytical performance models [21] and using compiler analysis [9] to assist the optimization. Similar to other existing work, we also manually applied the optimizing transformations for each application. However, our work is a step closer to complete automation than existing work in that our work fully integrates analytical performance modeling and compiler dependence analysis to automatically determine (with optional developer guidance) the profitability and safety of the overlapping optimization. Furthermore, while the special interprocedural pattern of loop-based communication-computation overlapping addressed by our work is common in scientific applications, it has not been addressed previously in the literature. While our example in Figure 1 contains only a single communication inside the loop body, the optimization works similarly when the body contains a chain of dependent communications, by overlapping them with independent computations of the previous iterations.

Our main technical contributions are the following:

- Our framework tackles a common case of enhancing the overlapping of computation-communication that has not been previously addressed for MPI applications and is able to fully automate the profitability and safety analysis of the optimization, by using advanced analytical performance modeling (which collectively considers all the dynamic paths through each code block) and compiler dependence analysis (which supports automated semantic inlining of developer-supplied knowledge). Developer guidance is required only to improve the accuracy of the analysis for large scientific applications, because not all source code of these applications is available and many low-level implementation details are impossible to fully automatically decipher [17]. We currently manually apply the necessary program transformations, because code need to be carefully moved across procedural boundaries. We believe this process can be automated with some developer guidance, which is our future work.
- We applied our approach to optimize 7 NAS Parallel Benchmarks (NPB) applications on both a high-speed and a slow network-connected cluster environment and achieved 3–88% speedup on both platforms.

The remainder of the paper is organized as follows. Section II presents our analytical performance modeling component for automatically identifying communication and computation hot spots in MPI applications. Section III discusses how to automatically determine the safety of the optimization through compiler analysis. Section IV summarizes strategies we used to perform the actual optimizations and the tuning of their configurations. Section V evaluates our framework using 7 NAS application benchmarks [4]. Section VI discusses related work, and Section VII presents our conclusions.

II. ANALYTICAL MODELING OF MPI APPLICATIONS

To effectively reduce the overhead of network communications in MPI applications, one must understand when and where it becomes beneficial to enhance the overlap of communications with local computations in these applications. Through an analytical approach, our framework aims to model the runtime execution flow of an input application in terms of its relative amount of time spent in local computations and network communications. This information then is used to automatically identify potential communication bottlenecks as candidates for optimization in the later steps.

To represent and estimate the time required to execute the local computation of each path, we use the Bayesian Execution Tree (BET) from the Skope analytical performance modeling framework [16]. Each BET essentially represents possible runtime code paths of an application together with their execution frequency and expected execution time. We use the Skope framework to automatically generate a BET representation of each application from the application source code combined with some sample input data and code-coverage profiling of the application execution. We then extend the Skope framework to additionally estimate the overhead of each MPI communication through the following steps.

- 1) Use a LogGP-based communication model for the MPI runtime to estimate the communication time for each individual MPI call.
- 2) Statistically estimate the expected average communication time for each code path by combining the individual communication with the execution frequencies.

Finally, the balance between the time required for each MPI communication and the expected execution time of its surrounding local computation is used to project optimization opportunities. The following first illustrates the BET representation that we inherit from [16] and then explains our extensions for modeling MPI communications.

A. Bayesian Execution Tree

Figure 3 shows an example BET for one of the MPI processes of the NAS FT benchmark in Figure 1a. Each node of the BET represents a code block (a sequence of statements in the user program) together with its *runtime execution frequency*, defined as the expected average number of times that statements in the code block will be executed at runtime. A *depth-first-traversal* (DFS) of each subtree of the BET corresponds to a possible runtime execution

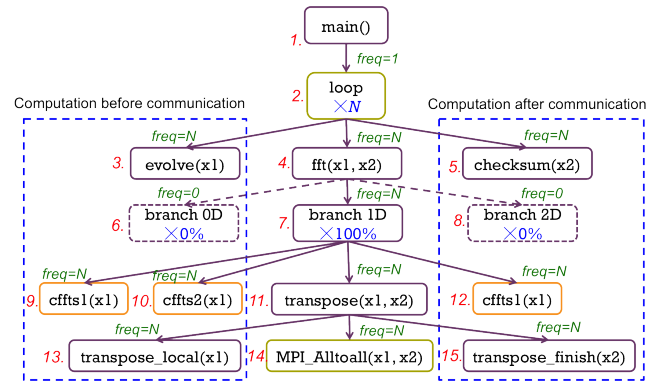


Fig. 3: Simplified Bayesian Execution Tree for NAS 1D FFT before overlapping computation and communications (*only important branches, loops, and function calls are shown*)

path of the statements. For example, in Figure 3, Node#2 is a loop of N iterations, so the frequency of its loop body is N . Node#7 is a branch inside the *fft* function. If the application is to perform 1D FFT, this branch is taken 100% of time during execution, so its frequency is N , while the frequency of the alternative branches (Node#6 and Node#8) are set to 0.

In order to derive execution frequencies of each code block, the Skope framework requires a description of the application input data, manually provided by the user or in our case, automatically collected via instrumented runs of the application. The input data description characterizes the possible values of data that the application may obtain from external sources, for example, command-line arguments, environment variables, or files. For array variables, only their dimensions and the size of each dimension are required. For MPI applications, the total number of MPI processes (`MPI_Comm_size`) and the rank of the process to model (`MPI_Rank`) are additionally required. Based on the input data description, the Skope framework applies constant propagation to derive possible values of the expressions that control the directions of branch and loop controls. A fall-through probability is assumed to be 50% if the values cannot be accurately determined. For this paper, we used `gcov` to profile applications with sample input data,

B. Modeling MPI Communications

To predict MPI communication overhead, we have extended the Skope framework with the LogGP model [2] to additionally model the latency (elapsed time) of each MPI operation using the following four parameters:

- 1) P : number of processes involved in the communication
- 2) n : size (in bytes) of the message being transferred
- 3) α : overhead of starting each message and time interval required between transmitting each pair of messages
- 4) β : communication time per byte for large messages, determined by the underlying network bandwidth.

Among the four parameters, α and β can be calculated ahead of time from characteristics of the underlying network. We compute β as the reciprocal of the network

bandwidth and $alpha$ by using microbenchmarks to measure the latency of `MPI_Send` and `MPI_Recv` operations on the target platform. The other two parameters, P and n , are determined through instrumented runs of the user application. In particular, P equals to `MPI_Comm_size`; and n is obtained from the values used to invoke the MPI operations.

Following the LogGP model, we model the cost (latency) of each MPI point-to-point communication as:

$$cost_{p2p}(n; alpha, beta) = alpha + n \cdot beta. \quad (1)$$

To model the `MPI_Alltoall` operation, we use the following two formulas.

$$cost_{short} = logP \cdot alpha + \frac{n}{2} \cdot logP \cdot beta \quad (2)$$

$$cost_{long} = (P - 1) \cdot alpha + n \cdot beta \quad (3)$$

The first formula models the latency of short messages and the second that of long messages. We use values of *control variables* from the MPI runtime library, for example, `MPIR_CVAR_ALLTOALL_SHORT_MSG_SIZE` for MPI alltoall, to determine whether a message should be categorized as *short* or *long* and thereby select the proper formulas to use.

After estimating the latency of each individual MPI operation, the overall communication time of a code path in BET can be calculated by adding the communication time of all code blocks along the path, using the following formula.

$$cost_m = \sum_i^m cost(i) * freq(i) \quad (4)$$

Specifically, the total communication time of a path of m nodes in the BET can be computed as the sum of the latency of each individual MPI operation multiplied by its execution frequency $freq(i)$. Here, $freq(i)$ is calculated as the same of the frequency of the BET node that contains the MPI operation, and $cost(i)$ is calculated as indicated above using LogGP formulas instantiated with the expected parameter values of the MPI operations. For example, the total communication time of `MPI_Alltoall` in Figure 3 can be computed by multiplying the average communication time of `MPI_Alltoall` by the number of iterations of the loop node#2 ($\times N$) and the fall-through probability of the 1D FFT branch node#7 ($\times 100\%$).

III. OPTIMIZATION ANALYSIS

The objective of our optimization analysis is to automatically identify which MPI communications to optimize and what local computations can be safely overlapped with the communication, through the following three steps.

- 1) Analytically identify MPI operations that are potential performance bottlenecks based on the modeling of communication cost and the execution flow modeling of the entire application described in Section II. In particular, based on the BET representation of the user application,

this step identifies from the top N most time-consuming MPI calls those that collectively take $\geq P\%$ of the overall communication time, where both N and P are user-configurable parameters and were set by default with $N = 10$ and $P = 80$. The selection is accomplished by simply sorting the pre-estimated communication time of all MPI calls in the BET and then selecting the top ones. For example, for the NAS FT application shown in Figure 3, a single MPI call, the `MPI_Alltoall` at the bottom of the BET, is selected since it takes more than 95% of the overall communication time.

- 2) For each identified MPI communication to optimize, locate the closest enclosing loops of the MPI communication in the BET—for example, node#2 in Figure 3 for the NAS FT application—to potentially overlap with the communication. If the enclosing loop does not exist, the communication is given up as an optimization target.
- 3) Suppose `comm(I)` is the MPI communication being considered for optimization at each loop iteration. Apply loop dependence analysis to identify the statements (`Before(I)`) that compute data to be transferred in `comm(I)` (that is, those that have dependence into `comm(I)`), and the statements (`After(I)`) that use the data transferred from `comm(i)` (that is, those that have dependence from `comm(I)`). Determine whether `comm(I)` is independent of `After(I-1)` and `Before(I+1)`; that is, whether it is safe to overlap `comm(I)` across loop iterations. If yes, the loop is selected for optimization.

A key challenge in optimizing large applications is that MPI communications are often scattered across procedural boundaries, and the computation that can be overlapped with them is often some distance away and similarly across abstraction boundaries. By using the Skope framework and through the BET representation of the whole user application, we are able to inter-procedurally select MPI communications as well as their surrounding loops as potential optimization targets. Then, an optimization pragma, `#pragma cco do`, illustrated at line 1 of Figure 5, is inserted automatically before each selected code region to instruct the compiler to perform additional analysis to determine the safety of optimization.

We use loop dependence analysis within the ROSE C/C++/Fortran compiler [33] to automatically determine the safety of the reordering optimization to each selected code region. By making the compiler inline all function calls within the region, the entire optimization analysis can be fully automated, if the compiler can find the source code of all the functions invoked inside the loop being optimized; when the source code of some of these function are not available for analysis, the compiler would opt to be conservative and deem the optimization unsafe. If such conservativeness is not desirable, the developer can insert the following pragmas to provide additional guidance to the loop dependence analysis.

- 1) `#pragma cco ignore`, illustrated at line 3 in Figure 5, which can be inserted before each function call that can be safely ignored when performing dependence

analysis. That is, these function calls will not implicate the safety of any reordering optimization. Examples of such function calls include the `timer_start()` and `timer_stop()` in Figure 5.

- 2) `#pragma cco override`, illustrated at the first line of Figure 4 and 6, which defines the memory side effects of the following function call. The override definitions, when manually specified, allow dependence analysis to proceed across procedural boundaries without requiring the source code or the inlining of the functions invoked.

For this paper, we manually inserted the above annotations to overcome situations where the source code of the callee is unavailable, or the low-level implementation details of the callee are too complex to be accurately deciphered by traditional compiler dependence analysis (e.g., the underlying implementations of MPI operations). Figure 4 shows an example override definition for `MPI_Alltoall`, where we use the `read` and `write` pseudo statements to indicate read and write memory accesses. Here based on the domain knowledge of the application that the data being sent/received have atomic types instead of user-defined types, the memory side effect of the operation can be expressed as consecutive read and write memory references to the communication buffers (`sendbuf` and `recvbuf`). We similarly composed memory side effect definitions for the other MPI operations.

Traditional loop dependence analysis in compilers is based on the disambiguation of subscripted array references, by solving diophantine equations parameterized by the surrounding loop index variables, to determine whether each pair of array references may refer to the same memory location at two arbitrary loop iterations [3]. The analysis therefore becomes ineffective when memory references are not expressed using array notations or when the subscript of an array reference cannot be expressed as a linear combination of the surrounding loop index variables. Additionally, since the compiler does not have any information about the runtime control flow of a program, it assumes all control paths can happen at runtime, and no unknown system calls (e.g., `timer_start` and `timer_stop` in Figure 5) can be reordered. Our annotation mechanisms are provided to developers to optionally overcome such conservativeness when desired. In our experience of optimizing the NAS application benchmarks, we have used these annotations to serve the following purposes.

- Define the memory side effects of MPI operations and system calls that can be ignored (e.g., Figure 4 and Figure 5), the source code of which is not available to the compiler. These annotations can be reused across MPI applications without additional work by the developers.
- Specialize the memory side effect of a function call when only a single runtime code path is known to be executed by the function call, to prevent the compiler considering all possible code paths through the default inlining mechanism. For example, in NAS FT, the procedure `fft` has 6 branches for solving different dimensions of the FFT problem (0D, 1D, or 2D), while only one branch will be

taken at each invocation. By manually overriding the default inlining, we can eliminate the unreachable branches from being considered. Figure 6 shows the annotation we used to override the `fft()` function in NAS FT. Here the original function have several branches for different data layout, while the override definition has only 1D layout that is the target code path to optimize. These annotations can be automated with developer approval by having the compiler directly interpret the modeling output of Skope.

```
$cco override
subroutine MPI_Alltoall(sendbuf, sendcount, sendtype,
> recvbuf, recvcount, recvtype, comm, ierror)
do i = 1, sendcount
  read sendbuf(i)
end do
do i = 1, recvcount
  write recvbuf(i)
end do
end subroutine
```

Fig. 4: Example: describe the memory side effect of `MPI_Alltoall()` for simple data types

```
1 !$cco do
2 do iter = 1, niter
3   !$cco ignore
4   if (timers_enabled) call timer_start(T_evolve)
5   call evolve(u0,u1,twiddle,dims(1,1),dims(2,1),dims(3,1))
6   !$cco ignore
7   if (timers_enabled) call timer_stop(T_evolve)
8   !$cco ignore
9   if (timers_enabled) call timer_start(T_fft)
10  call fft(-1,u1,u2)
11  !$cco ignore
12  if (timers_enabled) call timer_stop(T_fft)
13  !$cco ignore
14  if (timers_enabled) call timer_start(T_checksum)
15  call checksum(iter,u2,dims(1,1),dims(2,1),dims(3,1))
16  !$cco ignore
17  if (timers_enabled) call timer_stop(T_checksum)
18 end do
```

Fig. 5: Example: annotating a loop to optimize for NAS FT

```
$cco override
subroutine fft(dir, x1, x2)
  cffts1(-1,dims(1,3),dims(2,3),dims(3,3),x1,x1,scratch)
  transpose_xyz(3, 2, x1, x2)
  cffts2(-1,dims(1,2),dims(2,2),dims(3,2),x2,x2,scratch)
  cffts1(-1,dims(1,1),dims(2,1),dims(3,1),x2,x2,scratch)
end subroutine
```

Fig. 6: Example: using 1D layout code path to override the default inlining for NAS FT

IV. PROGRAM TRANSFORMATION

After being selected by our automatic compiler analysis component, each loop to be optimized includes at least one MPI communication, c , inside the loop body, and two sets of statements, $Before(c)$ and $After(c)$, that are safe to be overlapped with c across loop iterations. If the Loop body contains multiple MPI communications, say $c1$ followed by $c2$, both of which are selected for optimization, then $c1$ and $c2$ must be independent across different loop iterations, as $c1 \in Before(c2)$ and $c2 \in After(c1)$ must hold, and to be

both selected for optimization, we must have $c1(I)$ (the MPI communication $c1$ at loop iteration I) and $After(c1, I-1)$ (the statements after $c1$ at loop iteration $I-1$) being independent, and similarly $c2(I)$ and $Before(c2, I+1)$ being independent. Consequently, $c1$ and $c2$ must be independent of each other across loop iterations. So they can be optimized one after another without violating any dependence constraints, with each optimization considering only a single MPI communication and its surrounding loop, discussed in the following.

For each MPI communication and its surrounding loop to optimize, we currently manually transform the source code to enable the overlapping of computation and communication. The main difficulty of fully automating the optimization is to identify the actual statements corresponding to $Comm(I)$, $Before(I)$, and $After(I)$ respectively, as shown in Figure 7(a), when they are scattered across multiple procedures, and some of these procedures cannot be inlined in the compiler analysis phase, e.g., due to the need for code path specialization as demonstrated in Figure 6. Otherwise, our future work will automate the following currently manually applied steps.

A. Converting MPI communications

Convert each blocking MPI operation, for example, *alltoall* collectives and point-to-point send-receives, in $Comm(I)$ to an equivalent nonblocking communication combined with a blocking wait, as illustrated in converting Figure 7a to 7b.

B. Reordering computation and communication

Tag each statement in the loop body as belonging to $Before(I)$, $After(I)$, $Icomm(I)$, or $Wait(I)$, where I is the index variable of the surrounding loop, as shown in Figure 7(b). Then, *interleave* $Icomm(I)$ with $Before(I+1)$ and $After(I-1)$, as illustrated in Figure 7d, in two steps:

- 1) Move $Before(1)$ and $Icomm(1)$ to the outside before the first iteration of the loop starts, and move $Wait(N)$ and $After(N)$ outside after the last loop iteration as shown in Figure 7c.
- 2) Move $Before(I)$ and $Icomm(I)$ above $Wait(I-1)$ and $After(I-1)$ as shown in Figure 7d.

Figure 8 shows the new scheduling of the various computation and communication components after the reordering. Note that for each iteration I , $Before(I)$, $Icomm(I)$, $Wait(I)$ and $After(I)$ are evaluated in the same order as in the original computation. The main difference is that $After(I-1)$ and $Before(I+1)$ are now placed in between $Icomm(I)$ and $Wait(I)$, so that the non-blocking communication $Icomm(I)$ can now be processed in parallel with $After(I-1)$ and $Before(I+1)$, thereby facilitating overlapping of computation and communication. If $Before(I+1)$ or $After(I-1)$ contains other MPI communications that are independent of $comm(I)$, the same rescheduling modification can be applied to these communications in the same fashion, where $Icomm(I)$ and $wait(I)$ would belong to the $Before(I)$ or $After(I)$ components of these new communications.

```
DO I = 1 .. N
  Before(I)
  Comm(I)
  After(I)
END DO
```

(a) Input loop

```
DO I = 1 .. N
  Before(I)
  Icomm(I)
  Wait(I)
  After(I)
END DO
```

(b) Decouple blocking comm

```
Before(1)
Icomm(1)
DO I = 2 .. N
  Wait(I - 1)
  After(I - 1)
  Before(I)
  Icomm(I)
END DO
Wait(N)
After(N)
```

(c) Move first and last iterations

```
Before(1)
Icomm(1)
DO I = 2 .. N
  Before(I)
  Wait(I - 1)
  Icomm(I)
  After(I - 1)
END DO
Wait(N)
After(N)
```

(d) Interleave consecutive iterations

Fig. 7: Steps to reorder communication and computation

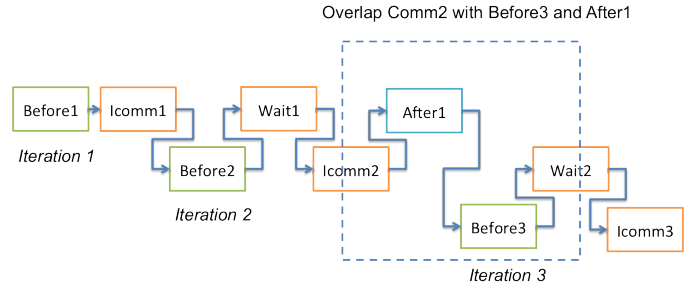


Fig. 8: Overlapped computation and communications

C. Replicating the communication buffer

Each MPI operation needs a dedicated buffer to hold the data being communicated. Applications typically first allocate the necessary communication buffers at the initialization stage and then reuse the same buffers in the MPI operations across different loop iterations. After applying our optimization, as illustrated in Figure 8, the communication ($Icomm(i)$ and $Wait(i)$) at each i th iteration, where $i \geq 2$, is overlapped with computation $Before(i+1)$ and $After(i-1)$. Assuming that two distinct buffers, *InBuf* and *OutBuf*, are used for sending and receiving each message, respectively, each buffer needs to be replicated into a pair of equal size to ensure that distinct buffers are used across the overlapping iterations, as illustrated in Figure 9. In particular, we replicate each buffer by allocating

```

Before(1, InBuf)
Icomm(1, InBuf, OutBuf)
DO I = 2 .. N
  Before(I, InBuf)
  Wait(I - 1)
  Icomm(I, InBuf, OutBuf)
  After(I - 1, OutBuf)
END DO
Wait(N, OutBuf)
After(N, OutBuf)

```

(a) Original code using one pair of input/output buffers

```

Before(1, InBuf)
Icomm(1, InBuf, OutBuf)
DO I = 2 .. N
  Before(I, I % 2 == 1 ? InBuf : InBuf2)
  Wait(I - 1)
  Icomm(I, I % 2 == 0 ? InBuf : InBuf2
        , I % 2 == 0 ? OutBuf : OutBuf2)
  After(I - 1, , I % 2 == 1 ? InBuf : InBuf2)
END DO
Wait(N, N % 2 == 1 ? OutBuf : OutBuf2)
After(N, N % 2 == 1 ? OutBuf : OutBuf2)

```

(b) After replicating the input/output buffers

Fig. 9: Replicate buffers for nonblocking communication

additional memory outside the loop and then alternately use a distinct buffer in every pair of consecutive loop iterations.

D. Inserting MPI_Tests

When using nonblocking MPI operations, some CPU time needs to be allocated, by embedding *MPI_Test* calls in the local computation, to ensure continuous progress of the communications. If the local computation is not inside a loop, we insert one or more *MPI_Test* calls evenly distributed into the computation. On the other hand, if the local computation is inside a loop, we insert *MPI_Test* into the beginning of the loop body and use a conditional variable to adjust its frequency. The inserted code is illustrated in Figure 10. In both cases, the frequency of *MPI_Test* is empirically adjusted as the application is ported to each architecture.

```

DO I = 1 ... L
  If I % Freq == 0
    MPI_Test
  Original_computation_statements
END DO

```

Fig. 10: Insert *MPI_Test* into a loop at specific frequency $Freq$

V. EXPERIMENTAL RESULTS

To evaluate the accuracy of our analytical modeling of MPI applications and the performance implications of our optimizations in automatically overlapping computation with communication, we applied our approach to model and optimize 7 MPI applications from the NAS NPB [4] on two clusters shown in Table I. The first Intel platform is a high-performance computing cluster with very fast internode communication through InfiniBand. The second platform is in a small data center where the internode communication is through relatively slow Ethernet. Both clusters use MPICH 3.1.1 [15] as the underlying MPI library. Because our current

TABLE I: Experiment platforms

Server	Intel	HP ProLiant BL460c Gen6
Instruction set	Intel Xeon x86	Intel Xeon x64
Frequency	2.6 GHz	3.2 GHz
Compiler	ICC/Ifort 13.1	GCC/Gfortran 4.4.7
Network	InfiniBand Qlogic QDR	1 Gbps Ethernet
Total nodes	301	24 on 3 racks
Max memory	64 GB	48 GB

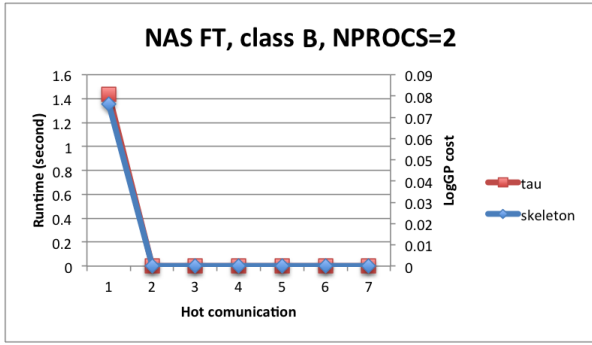
analytical performance model cannot estimate intra-node MPI communication time, we have allocated a single MPI process per node on both clusters.

For each MPI application, we first used our extended Skope modeling framework to find the most time-consuming MPI communication in the application. Then, after using the ROSE compiler to determine the safety of overlapping each communication with its surrounding loop, we manually applied the optimization to the enclosing loop. We measured the performance improvements from the optimizations using input data provided by the NPB benchmarks and using a range of 2 to 9 nodes for each application. Besides using the built-in timers within the NPB applications to collect their overall performance, we manually instrumented the source code of the applications to report the performance of individual communications. We used the class B input of the NAS benchmarks in most of our studies, as the bigger input sizes (class C and class D) took too long to run using the relatively small number of nodes available on our platforms.

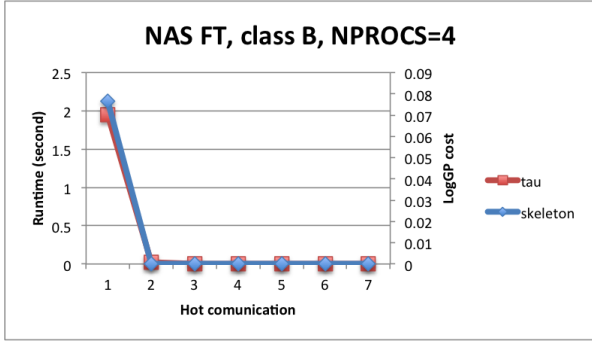
A. Accuracy of Hot Communication Prediction

To evaluate the accuracy of our modeling of MPI communications, Table II shows the differences in hot communication selections for the NPB applications, by comparing the set of communication hot spots selected using our model with those found by profiling the NPB applications, using class B input on 4 nodes. When selecting the top time-consuming MPI communications, we required that their overall time be at least 80% of the application’s overall communication time. In this case, our predictive modeling selected the same set of hot communications as found using application profiling. When asked to select a given number of the most time-consuming communications, the output by our predictive modeling differs by at most 2 selections compared with using profiling, for the NAS LU benchmark. Here the most time-consuming communications are pairs of sends/receives at four symmetric directions, which were estimated to take the same time by our predictive modeling. However, their actual runtime collected through profiling differ by 37%, because the execution of the processes is unbalanced, resulting in extra wait time to synchronize the corresponding *MPI_Send*/*MPI_Recv* operations.

Figure 11 compares our projected communication time with the actual measured time for NAS FT, using two and four processors. Here in spite of the small error rates in projecting the absolute values of the communication time, our



(a) Communication on 2 nodes



(b) Communication on 4 nodes

Fig. 11: Profiled runtime and modeled cost of NAS FT with middle-sized input (B) on x86 cluster

TABLE II: Differences between the projected hot-spot selection and the measured hot-spot selection based on profiling with 80% threshold for class B data on 4 nodes. Zero means the set of N hot spots equals the top N hot spots.

	Selected number of hot MPI communications							
	1	2	3	4	5	6	7	8
FT	0							
IS	0	0						
CG	0							
LU	0	1	2	2	1	1	0	0
MG	1	1	0	1	1	0		

modeling framework was able to accurately capture the relative importances of the various communication operations.

B. Impact of Optimizations

Figures 12 and 13 show the speedups we attained by enabling better computation-communication overlapping for the NPB applications. For each benchmark, the performance of the original and the optimized code is measured by using input class B on 2, 4, 8, and 9 nodes, with one MPI process bound to each node, with the exception of NAS BT and SP, which require the number of processes to be N^2 , where N is the number of nodes, and so are evaluated using 4 and 9 processes only. The overall elapsed time of each application is measured by using NPB’s built-in timer.

Our optimization attained 3–88% speedup for the NPB applications, with an average speedup of 18.8% on the InfiniBand cluster and 16.6% on the Ethernet. Significant speedups are

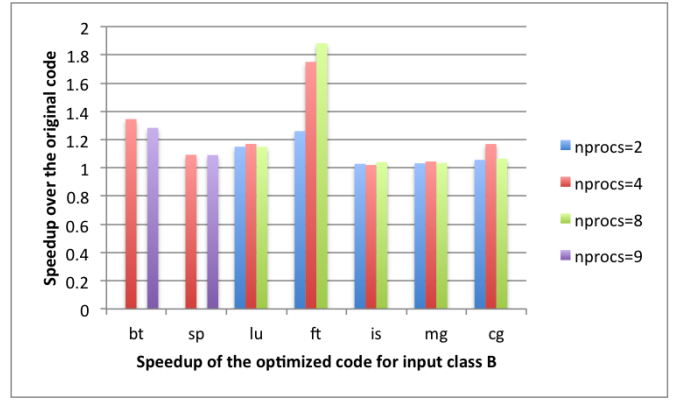


Fig. 12: Optimization speedups on the InfiniBand cluster.

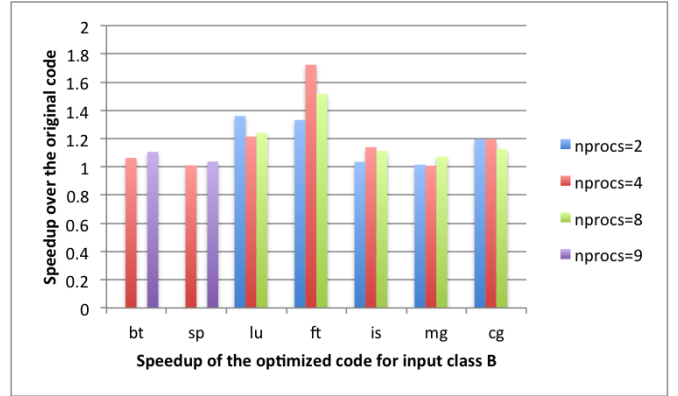


Fig. 13: Optimization speedups on the Ethernet cluster.

attained for FT and IS, which use *alltoall* collectives as the main communication operation, and less significant speedups for the other benchmarks, which mostly use point-to-point send/receives. The highest 88% improvement is observed with FT, which uses a time-consuming *MPI_Alltoall* operation, enclosed inside the outermost loop of the application, to exchange a large amount of data. The lowest speedup (3%) is observed with MG, which does not have sufficient local computation in the surrounding loop of the MPI communication to overlap with communication.

While our optimization was able to attain a consistent level of performance improvement on both platforms, the best speedups are observed when using different runtime configurations on these platforms. For example, the best speedup for NAS FT was attained when using 8 processors on the infiniband cluster but when using two processors on the Ethernet cluster, as the communication latency on the Ethernet is much longer than that of the Infiniband network, which in turn affects the amount of local computation time required to overlap with the communication. Overall, the possible speedup attained is bounded by the latency of the communication being optimized and the amount of available local computation to overlap with the communication time. A larger amount of local independent computation is generally required to fully compensate the latencies of a slower network (e.g., Ethernet) than that of a faster network (e.g., infiniband).

VI. RELATED WORK

Our work focuses on application-level performance enhancement by enabling automated overlapping of MPI communications with independent local computations. It therefore serves a different purpose than existing work on enhancing the communication efficiency of the many MPI operations, for example, alternative protocols for point-to-point communications [5], [10], collective operations [14], [26], [37], [38], remote direct memory access (RDMA) [18], [24], [41], load balancing [23], [28], and the elimination of redundant communications through software caching and the exploitation of data locality [7], [29], [39].

The Bamboo source-to-source translator [36] also aims to automatically translate MPI applications to latency-tolerant forms but relies on developer annotations to identify MPI regions to optimize. In contrast, our work automatically determines the profitability of the optimizations from analytical performance modeling. Iancu et al. [21] automatically selected message sizes and schedules for MPI communications through an analytical model of system scale and load. Danalis et al. [9] investigated compiler optimizations to potentially automate the overlapping of computations and MPI communications, by formulating the side effects of key MPI operations so that an MPI-aware compiler can automatically assess the safety of several optimizations, which were then manually applied in their paper. Various patterns of computation-communication overlapping and automated optimization schemes have also been discussed [8], [13]. Our work particularly focuses on automatically enabling a special form of loop-based communication-computation overlapping in scientific applications, a form that has not been addressed by existing work.

Sancho et al. [34] combined empirical tuning with networking models to quantify the potential benefit of overlapping communication and computation. Potluri et al. [31] empirically quantified the overlapping of MPI-2 operations in a seismic modeling application. Hu et al. [20], [35] identified the consumer-producer model from the control flow graph of the application to guide optimization decisions for overlapping Alltoall communication in a 3-D FFT. Didelot et al. [11], [12] developed a message progression model based on collaborative polling that allows an efficient auto-adaptive overlapping of communication phases with computation. In this paper, we predicted the most time-consuming code paths that contain MPI communications to optimize, using existing analytical models of communications [2].

Preissl et al. [32] summarized common communication patterns in MPI applications to enable automated optimization. Pellegrini et al. [30] proposed an exact dependence analysis approach for increasing the overlapping of computation and communication. Subotic et al. [40] speculatively extracted runtime data-flow to understand the dynamic dependence of the application. Aananthakrishnan et al. [1] used a hybrid static and runtime data-flow analysis of MPI programs. We also use dependence analysis to determine the correctness of optimization, enhanced with additional knowledge from

developers about their applications.

In order to find the optimal placement of nonblocking MPI operations within the computation control flow, accurate modeling of the underlying computation and communication is required [6]. Hoefer et al. [19] presented an analytical approach to model MPI barriers. Ino et al. [22] presented a parallel computational model for synchronization analysis in MPI. Martinez et al. [25] developed an analytical model extending LogGP [2] for accurate estimation of individual MPI communication. Moritz and Frank [27] modeled network contention in MPI applications. In our optimization, we first reposition each pair of local computation and nonblocking communication as far apart as safety allows across different loop iterations and then insert MPI_Test with empirically tuned frequencies into the local computation to ensure proper progress of the nonblocking communication.

VII. CONCLUSION

This paper presents a systematic approach to automate the overlapping of communications with independent computations in large MPI applications, thereby enhancing their performance portability. Our optimization workflow starts with analytical performance modeling of the overall application execution flow to identify long-lasting MPI communications to overlap. Next, we conduct automated safety and profitability analysis (with optional developer assistance) to find optimization opportunities. We complete the optimization by manually applying the necessary transformations in a systematic fashion that can be potentially automated. We applied our approach to optimize 7 NAS NPB applications on both a high-speed and a slow network-connected cluster environment. We achieved 3–88% speedup on both platforms.

REFERENCES

- [1] Sriram Aananthakrishnan, Greg Bronevetsky, and Ganesh Gopalakrishnan. Hybrid approach for data-flow analysis of mpi programs. In *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing, ICS '13*, pages 455–456, New York, NY, USA, 2013. ACM.
- [2] Albert Alexandrov, Mihai F. Ionescu, Klaus E. Schauer, and Chris Scheiman. LogGP: Incorporating long messages into the LogP model :one step closer towards a realistic model for parallel computation. In *Proceedings of the Seventh Annual ACM Symposium on Parallel Algorithms and Architectures, SPAA '95*, pages 95–105, New York, NY, USA, 1995. ACM.
- [3] R. Allen and Ken Kennedy. *Optimizing Compilers for Modern Architectures*. Morgan Kaufmann, San Francisco, October 2001.
- [4] D. H. Bailey, E. Barszcz, J. T. Barton, D. S. Browning, R. L. Carter, R. A. Fatoohi, P. O. Frederickson, T. A. Lasinski, H. D. Simon, V. Venkatakrisnan, and S. K. Weeratunga. The NAS parallel benchmarks. Technical report, The International Journal of Supercomputer Applications, 1991.
- [5] Ron Brightwell and Keith Underwood. Evaluation of an eager protocol optimization for MPI. In *Proceedings of EuroPVM/MPI*, pages 327–334, 2003.
- [6] Ron Brightwell and Keith D. Underwood. An analysis of the impact of MPI overlap and independent progress. In *Proceedings of the 18th Annual International Conference on Supercomputing, ICS '04*, pages 298–305, New York, NY, USA, 2004. ACM.
- [7] Darius Buntinas, Brice Goglin, David Goodell, Guillaume Mercier, and Stéphanie Moreaud. Cache-efficient, intranode, large-message MPI communication with MPICH2-Nemesis. In *Proceedings of the 2009 International Conference on Parallel Processing, ICPP '09*, pages 462–469, Washington, DC, USA, 2009. IEEE Computer Society.

- [8] Anthony Danalis, Ki-Yong Kim, Lori Pollock, and Martin Swany. Transformations to parallel codes for communication-computation overlap. In *Proceedings of the 2005 ACM/IEEE Conference on Supercomputing, SC '05*, pages 58–, Washington, DC, USA, 2005. IEEE Computer Society.
- [9] Anthony Danalis, Lori Pollock, Martin Swany, and John Cavazos. MPI-aware compiler optimizations for improving communication-computation overlap. In *Proceedings of the 23rd International Conference on Supercomputing, ICS '09*, pages 316–325, New York, NY, USA, 2009. ACM.
- [10] Alexandre Denis. A high performance superpipeline protocol for InfiniBand. In *Proceedings of the 17th International Conference on Parallel Processing - Volume Part II, Euro-Par'11*, pages 276–287, Berlin, Heidelberg, 2011. Springer-Verlag.
- [11] Sylvain Didelot, Patrick Carribault, Marc Pérache, and William Jalby. Improving MPI communication overlap with collaborative polling. In *Proceedings of the 19th European Conference on Recent Advances in the Message Passing Interface, EuroMPI'12*, pages 37–46, Berlin, Heidelberg, 2012. Springer-Verlag.
- [12] Sylvain Didelot, Patrick Carribault, Marc Pérache, and William Jalby. Improving MPI communication overlap with collaborative polling. *Computing*, 96(4):263–278, April 2014.
- [13] Lewis Fishgold, Anthony Danalis, Lori Pollock, and Martin Swany. An automated approach to improve communication-computation overlap in clusters. In *Proceedings of the 20th International Conference on Parallel and Distributed Processing, IPDPS'06*, pages 290–290, Washington, DC, USA, 2006. IEEE Computer Society.
- [14] Richard L. Graham and Galen Shipman. MPI support for multi-core architectures: Optimized shared memory collectives. In *Proceedings of the 15th European PVM/MPI Users' Group Meeting on Recent Advances in Parallel Virtual Machine and Message Passing Interface*, pages 130–140, Berlin, Heidelberg, 2008. Springer-Verlag.
- [15] William Gropp. MPICH2: A new start for MPI implementations. In *Proceedings of the 9th European PVM/MPI Users' Group Meeting on Recent Advances in Parallel Virtual Machine and Message Passing Interface*, pages 7–, London, UK, 2002. Springer-Verlag.
- [16] Jichi Guo, Jiayuan Meng, Qing Yi, Vitali A. Morozov, and Kalyan Kumaran. Analytically modeling application execution for software-hardware co-design. In *2014 IEEE 28th International Parallel and Distributed Processing Symposium*, pages 468–477, 2014.
- [17] Jichi Guo, Mike Stiles, Qing Yi, and Kleantlis Psarris. Enhancing the role of inlining in effective interprocedural parallelization. In *ICPP'11: International Conference On Parallel Processing*, Taipei, Taiwan, September 2011.
- [18] Masayuki Hatanaka, Atsushi Hori, and Yutaka Ishikawa. Optimization of MPI persistent communication. In *Proceedings of the 20th European MPI Users' Group Meeting, EuroMPI '13*, pages 79–84, New York, NY, USA, 2013. ACM.
- [19] Torsten Hoefler, Lavinio Cerquetti, and Frank Mietke. A practical approach to the rating of barrier algorithms using the logp model and open mpi. In *Proceedings of the 2005 International Conference on Parallel Processing Workshops, ICPPW '05*, pages 562–569, Washington, DC, USA, 2005. IEEE Computer Society.
- [20] Changjun Hu, Yewei Shao, Jue Wang, and Jianjiang Li. Automatic transformation for overlapping communication and computation. In *Proceedings of the IFIP International Conference on Network and Parallel Computing, NPC '08*, pages 210–220, Berlin, Heidelberg, 2008. Springer-Verlag.
- [21] Costin C. Iancu and Erich Strohmaier. Optimizing communication overlap for high-speed networks. In *Proceedings of the 12th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPoPP '07*, pages 35–45, New York, NY, USA, 2007. ACM.
- [22] Fumihiko Ino, Noriyuki Fujimoto, and Kenichi Hagihara. LogGPS: A parallel computational model for synchronization analysis. In *Proceedings of the Eighth ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming, PPoPP '01*, pages 133–142, New York, NY, USA, 2001. ACM.
- [23] Vivek Kale, Amanda Randles, and William D. Gropp. Locality-optimized mixed static/dynamic scheduling for improving load balancing on SMPs. In *Proceedings of the 21st European MPI Users' Group Meeting, EuroMPI/ASIA '14*, pages 115:115–115:116, New York, NY, USA, 2014. ACM.
- [24] Jiuxing Liu, Jiesheng Wu, Sushmitha P. Kini, Pete Wyckoff, and Dhabaleswar K. Panda. High performance RDMA-based MPI implementation over InfiniBand. In *Proceedings of the 17th Annual International Conference on Supercomputing, ICS '03*, pages 295–304, New York, NY, USA, 2003. ACM.
- [25] D. R. Martinez, J. C. Cabaleiro, T. F. Pena, F. F. Rivera, and V. Blanco. Accurate analytical performance model of communications in mpi applications. In *Proceedings of the 2009 IEEE International Symposium on Parallel & Distributed Processing, IPDPS '09*, pages 1–8, Washington, DC, USA, 2009. IEEE Computer Society.
- [26] Anshul Mittal, Nikhil Jain, Thomas George, Yogish Sabharwal, and Sameer Kumar. Collective algorithms for sub-communicators. *SIGPLAN Not.*, 47(8):315–316, February 2012.
- [27] Csaba Andras Moritz and Matthew I. Frank. LoGPC: Modeling network contention in message-passing programs. *IEEE Trans. Parallel Distrib. Syst.*, 12(4):404–415, April 2001.
- [28] Sun Nian and Liang Guangmin. Dynamic load balancing algorithm for MPI parallel computing. In *Proceedings of the 2009 International Conference on New Trends in Information and Service Science, NISS '09*, pages 95–99, Washington, DC, USA, 2009. IEEE Computer Society.
- [29] David Ozog, Sameer Shende, Allen Malony, Jeff R. Hammond, James Dinan, and Pavan Balaji. Inspector/executor load balancing algorithms for block-sparse tensor contractions. In *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing, ICS '13*, pages 483–484, New York, NY, USA, 2013. ACM.
- [30] Simone Pellegrini, Torsten Hoefler, and Thomas Fahringer. Exact dependence analysis for increased communication overlap. In *Proceedings of the 19th European Conference on Recent Advances in the Message Passing Interface, EuroMPI'12*, pages 89–99, Berlin, Heidelberg, 2012. Springer-Verlag.
- [31] Sreeram Potluri, Ping Lai, Karen Tomko, Sayantan Sur, Yifeng Cui, Mahidhar Tatineni, Karl W. Schulz, William L. Barth, Amitava Majumdar, and Dhabaleswar K. Panda. Quantifying performance benefits of overlap using MPI-2 in a seismic modeling application. In *Proceedings of the 24th ACM International Conference on Supercomputing, ICS '10*, pages 17–25, New York, NY, USA, 2010. ACM.
- [32] Robert Preissl, Martin Schulz, Dieter Kranzlmüller, Bronis R. de Supinski, and Daniel J. Quinlan. Transforming MPI source code based on communication patterns. *Future Gener. Comput. Syst.*, 26(1):147–154, January 2010.
- [33] Daniel J. Quinlan. Rose: Compiler support for object-oriented frameworks. *Parallel Processing Letters*, 10(2/3):215–226, 2000.
- [34] José Carlos Sancho, Kevin J. Barker, Darren J. Kerbyson, and Kei Davis. Quantifying the potential benefit of overlapping communication and computation in large-scale scientific applications. In *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing, SC '06*, New York, NY, USA, 2006. ACM.
- [35] Sukhyun Song and Jeffrey K. Hollingsworth. Designing and auto-tuning parallel 3-D FFT for computation-communication overlap. *SIGPLAN Not.*, 49(8):181–192, February 2014.
- [36] Eric Bylaska Dan Quinlan Tan Nguyen, Pietro Cicotti and Scott Baden. Bamboo - translating mpi applications to a latency-tolerant, data-driven form. In *Proceedings of the 2012 ACM/IEEE conference on Supercomputing (SC12)*, Salt Lake City, Utah, United States, 2012.
- [37] Jesper Larsson Träff. Optimal MPI datatype normalization for vector and index-block types. In *Proceedings of the 21st European MPI Users' Group Meeting, EuroMPI/ASIA '14*, pages 33:33–33:38, New York, NY, USA, 2014. ACM.
- [38] Jesper Larsson Träff and Antoine Rougier. MPI collectives and datatypes for hierarchical all-to-all communication. In *Proceedings of the 21st European MPI Users' Group Meeting, EuroMPI/ASIA '14*, pages 27:27–27:32, New York, NY, USA, 2014. ACM.
- [39] Yuichi Tsujita, Atsushi Hori, and Yutaka Ishikawa. Locality-aware process mapping for high performance collective MPI-IO on FEFS with Tofu interconnect. In *Proceedings of the 21st European MPI Users' Group Meeting, EuroMPI/ASIA '14*, pages 157:157–157:162, New York, NY, USA, 2014. ACM.
- [40] J. Labarta V. Subotic and M. Valero. Overlapping MPI computation and communication by enforcing speculative dataflow. 2008.
- [41] Tim S. Woodall, Galen M. Shipman, George Bosilca, Richard L. Graham, and Arthur B. Maccabe. High performance RDMA protocols in HPC. In *Proceedings of the 13th European PVM/MPI User's Group Conference on Recent Advances in Parallel Virtual Machine and Message Passing Interface, EuroPVM/MPI'06*, pages 76–85, Berlin, Heidelberg, 2006. Springer-Verlag.