

# Dynamic Adaptable Asynchronous Progress Model for MPI RMA Multiphase Applications

Min Si, Antonio J. Peña, Jeff Hammond, Pavan Balaji, Masamichi Takagi, Yutaka Ishikawa

**Abstract**—Casper is a process-based asynchronous progress model for MPI one-sided communication on multi- and many-core architectures. The one-sided communication is not yet truly one-sided in most MPI implementations: the target process still relies on software progress to complete incoming operations. Casper allows the user to specify an arbitrary number of cores dedicated to background ghost processes and transparently redirects the user RMA operations to ghost processes by utilizing the PMPI redirection and MPI-3 shared-memory technologies. Although Casper is effective and efficient for applications suffering from lack of asynchronous progress, permanently redirecting operations to a small number of ghost processes might not support complex multiphase applications effectively, which often involve dynamically changing communication density and computing workloads.

In this paper, we present an adaptive mechanism in Casper to address the limitation of static asynchronous progress in multiphase applications. We exploit two adaptive strategies, a precise user-guided strategy and a fully transparent and automatic strategy based on self-profiling and prediction, to dynamically reconfigure the asynchronous progress in Casper according to real-time performance characteristics during multiphase execution. We evaluate the adaptive approaches in both microbenchmarks and a real quantum chemistry application suite, NWChem, on the Cray XC30 supercomputer and an Intel Omni-Path cluster.

**Index Terms**—MPI; multiphase; one-sided; RMA; adaptation; asynchronous progress;

## 1 INTRODUCTION

Advances in high-end computing systems enable scientists to solve complex and large-scale problems with the integration of various fundamental solvers and algorithm modules. Tuning the configuration of runtime systems is a nontrivial task for obtaining highly efficient application performance. This task can be particularly challenging in multiphase applications because of their dynamically changing characteristics of communication and computation during the execution of multiple internal phases, especially when some of the internal phases prefer exactly opposite runtime configurations.

MPI is the dominant parallel programming model on distributed-memory systems. The one-sided communication model (also known as RMA) allows one process to specify all communication parameters for both the sending and receiving sides. Thus, a process can access a memory region of another process without the target process explicitly needing to receive the message. This asynchronous feature of RMA potentially provides a natural model for some applications that rely on irregular data movement [1], [2], [3], [4]. In practice, however, the RMA communication is not truly asynchronous in most MPI implementations. The reason is that even on RDMA-supported networks such as InfiniBand and the Cray Aries interconnect, most

ACCUMULATE operations still have to be handled in MPI software because of the limitations in hardware and MPI semantics. The completion of software-handled operations relies on explicit progress polling on the target process (e.g., making MPI calls). Consequently, an arbitrarily long delay in communication can occur if the target is busy in computation outside MPI.

The traditional approaches to ensure asynchronous progress for MPI communication have relied on two models. One is to utilize the background threads dedicated to each MPI process in order to handle incoming messages from other processes [5]. This model is widely provided in mainstream MPI implementations [6], [7], [8], [9]. However, the fundamental limitation of using this model in real applications is that the thread-based concept requires as many background threads as MPI processes on the system node. Thus, the user has to choose either to lose half of the computing cores or to enable expensive core oversubscription. In addition, this model requires MPI multithreading safety, that is known to be expensive because of the internal thread synchronization [10]. The other model of asynchronous progress is to utilize hardware interrupts to awaken a kernel thread on the target side and process the incoming RMA data within the interrupt context. The interrupt-based model can be found in Cray MPI [11] and IBM MPI on Blue Gene/P [12, Chapter 7]. Using this model, however, is limited in that an interrupt has to be generated on the target side to process every incoming data, which can be expensive.

In our previous work, we proposed an alternate model for asynchronous progress in MPI one-sided communication, called Casper [13]. Casper provides users the ability to specify a small number of cores as background “ghost processes;” it then transparently intercepts all RMA communication calls to the application processes and redirects them

- Min Si and Pavan Balaji are with the Mathematics and Computer Science Division at Argonne National Laboratory, USA.  
Email: {msi,balaji}@anl.gov
- Antonio J. Peña is with the Barcelona Supercomputing Center (BSC).  
Email: antonio.pena@bsc.es
- Jeff Hammond is with Intel Corporation, USA.  
Email: jeff.hammond@acm.org
- Masamichi Takagi and Yutaka Ishikawa are with the Riken Advanced Institute for Computational Science, Japan.  
Email: {masamichi.takagi, yutaka.ishikawa}@riken.jp

to the hidden ghost processes. Thus, the data movement can be completed by the ghost processes asynchronously while the application processes are tied by computation outside MPI. Unlike the traditional models, such a process-based model can avoid expensive overheads from either multi-threading safety or system interrupts. More importantly, it allows the user to control the number of cores being utilized for asynchronous progress, which we believe is a more suitable solution for applications running on multi- and many-core architectures.

The core concept of Casper is to redirect the RMA communication to a few number of ghost processes. This is suitable for common cases that need only a few cores to trigger data movements while the other resources are used to accelerate computing tasks. However, this design raises a potential bottleneck in that a large number of software-handled operations, which were handled by multiple cores on the node, are aggregated to only one or a few “asynchronous cores” in Casper. Such a small amount of progress resources might not be able to complete intensive operations quickly. Especially, when the application does not involve heavy computation and communication becomes dominant, this bottleneck might even counteract the benefit of asynchronous progress and result in poor performance. For single-phase applications, the user of Casper can adjust the number of cores according to the workloads of communication and computation. This method becomes impractical, however, when the application comprises multiple phases, some of which even indicate opposite performance patterns. That is, *the computation-dominant phases can benefit from asynchronous progress with only a small number of asynchronous cores, but some other communication-intensive phases might suffer performance degradation because of the overaggregated operations in Casper*. Thus, there is no way to deliver optimal performance for overall execution.

To address such complications, we present an adaptive mechanism in this paper that allows Casper to dynamically reconfigure the asynchronous progress during the execution of an application’s multiple phases. We analyze the performance tradeoff with regard to RMA progress, and design the adaptation to pursue the optimal performance for the overall execution of multiphase applications. We exploit two strategies. One is an effective user-guided strategy where the user can trigger reconfiguration in every application phase through MPI info hints. The other is a fully transparent strategy that involves automatic self-profiling and performance prediction at application runtime.

We carefully design the framework to ensure strict correctness in accordance with MPI-3 semantics. We evaluate the proposed adaptive approaches through both a set of microbenchmarks and a real application suite on a Cray XC30 supercomputer and an Intel Omni-Path Fabric based cluster. Throughout this work, we conclude that the process-based asynchronous progress model is a highly efficient, flexible, and portable approach for MPI RMA.

## 2 BACKGROUND

In this section, we briefly introduce the MPI RMA communication model and its implementation limitations on

modern network architectures. Here we highlight the important semantics on which this work highly relies. The comprehensive description of RMA semantics can be found in the MPI Standard [14] and past papers (e.g., [15]).

**MPI RMA Semantics:** MPI RMA is introduced in the MPI-2 and MPI-3 standards. To initialize an “RMA conversation,” every process in the communicator collectively creates a *window* as the exposure of their local memory regions, and a data transferring phase (called *epoch*) is opened and closed by a set of synchronization calls. During the epoch, a process can access the memory region on a remote process by issuing an RMA operation.

MPI defines an active-target synchronization mode including the *fence* and *post-start-complete-wait* (PSCW) epochs, and a passive-target mode that includes the *lock* and *lockall* epochs. A fence epoch requires all processes in the window to make the `MPI_WIN_FENCE` synchronization call; A PSCW epoch requires the processes in the origin group and those in the target group to make the `MPI_WIN_START|COMPLETE` and `MPI_WIN_POST|WAIT` calls, respectively; the lock or lockall epoch requires only the origin process to make the `MPI_WIN_LOCK|UNLOCK{ALL}` calls.

The data movement in RMA is defined through the operation calls including PUT, GET, and a set of ACCUMULATE operations (i.e., ACCUMULATE, GET\_ACCUMULATE, FETCH\_AND\_OP, and COMPARE\_AND\_SWAP). The ACCUMULATES guarantee strict *ordering* and *atomicity* for element-wise atomic access to remote memory locations (see page 461 in [14]). A closing synchronization call or an `MPI_WIN_FLUSH{ALL}` call in the passive target mode ensures the completion of operations.

**RMA Implementation:** The one-sided semantics provide MPI runtime developers the opportunity to offload the data movement to the hardware of Remote Direct Memory Access (RDMA) supported networks such as InfiniBand, Cray Aries, and Fujitsu Tofu. However, the state-of-the-art implementations are usually limited by two reasons. Firstly, most RDMA networks are able to process only simple data formats due to the limited processing power on the network interface controller (NIC). Complex operations such as computing a multi-dimensional non-contiguous double array still have to be handled by CPUs in the MPI software. Secondly, the RMA semantics force the runtime to guarantee ordering and atomicity among ACCUMULATE operations. Thus none of the ACCUMULATES can be offloaded as long as a data format is unsupported in hardware. Consequently, the MPI implementations for RDMA networks (e.g., MVA-PICH, Cray MPI) usually only offload PUTs and GETs with simple data formats to the hardware and keep the handling of other operations in MPI stack.

## 3 CASPER OVERVIEW AND CHALLENGE

In our previous work, we proposed Casper, a process-based asynchronous progress model for MPI one-sided communication [13], [16]. In this section, we present a brief overview of the Casper framework and discuss the challenge we observed in multiphase applications.

### 3.1 Overview

Casper is a process-based asynchronous progress model for MPI RMA on multi- and manycore architectures [13], [16].

It allows a few user-defined cores to be kept aside as background “ghost processes,” which are dedicated to helping the asynchronous progress for the application processes on the same node.

Casper is designed as an external library through the PMPI name-shifted profiling interface of MPI and transparently provides asynchronous progress for any RMA communication by overloading the necessary MPI functions. This design allows Casper to be platform and MPI implementation independent and enables the user to easily link Casper into the application binary.

When an application process tries to allocate an RMA window, Casper intercepts this call and internally allocates a shared-memory window for all application processes and the ghost processes that are on the same node, using the portable MPI-3 `MPI_WIN_ALLOCATE_SHARED` function. Thus, the ghost processes are able to access the window region located in the memories of the application processes. Then whenever a process tries to issue an RMA operation to the target process, Casper intercepts the call through PMPI and transparently redirects this message to a ghost process on the target node. The ghost processes simply wait in an `MPI_RECV` loop. Therefore, the MPI runtime can quickly make progress for any operations that are handled in the software stack of those ghost processes, and RMA operations that are offloaded to hardware see no difference in the way they behave.

The optimal number of ghost processes is *platform specific and application specific*. Choosing the optimal configuration is essentially a tradeoff between the amount of resources assigned for computation and that assigned for RMA progress. In practice, using one ghost process per node or one per socket is sufficient for most scientific applications running on CPU cores. This allows the remaining cores to be used to fulfill the heavy computing tasks.

### 3.2 Challenge in Multiphase Applications

Although using only one or a few asynchronous cores is suitable for most applications, such a small amount of progress resources might also lead to a performance bottleneck in some cases. That is, intensive software-handled RMA operations that were completed by a number of application cores on a node, are redirected to a few cores in Casper. This processing of overaggregated operations can be slow, and might eventually degrade the overall performance of an applications if the following two conditions are met: (1) the portion of computation is less significant than the data movement cost, and (2) the number of dedicated cores is much smaller than that of the remaining computing cores.

In most single phase applications, the user of Casper can empirically adjust the number of cores for ghost processes in order to workaround the above situation (avoid condition (2)). Such a method becomes impractical, however, when the application involves multiple internal phases and especially some of the heavy phases perform opposite performance characteristics. For instance, an internal phase may perform extremely expensive computation with a few data movement but the other phase may be dominated by enormous communication. It can be optimal to the former if redirecting communication to only a single asynchronous

core in Casper since the majority of core resources are still used to accelerate the computation; but such setting can cause severe overaggregation bottleneck in the latter phase. Unfortunately, a performance trade-off must be made for overall execution.

In order to provide the optimal overall performance, it is straightforward that we need an adaptive mechanism in Casper to dynamically update the message redirection for different application phases. To be specific, we need to address the following three questions.

Q-1. When does an adaptation become necessary?

Q-2. How can we make the adaptation?

Q-3. Where can the adaptation be taken?

## 4 DECOMPOSING RMA PROGRESS

To answer the question Q-1, we need first understand the MPI internal overhead for RMA communication. Because the asynchronous progress is needed only when the target process cannot make progress (e.g., computing outside MPI), we consider the simple scenario where the origin process initializes and completes the RMA conversation (e.g., issuing an `ACCUMULATE` and waiting in a flush), and the target process does computation. Figure 1 demonstrates the life time of such an RMA progress.

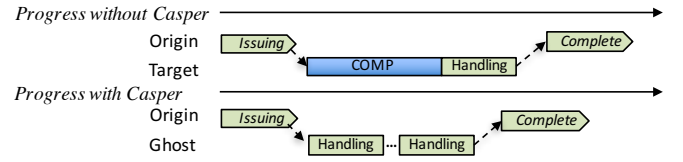


Fig. 1. Decomposing RMA progress.

We consider the completion of an RMA conversation can be decomposed into four portions: the operation issuing taken by the origin process, the network transfer between the origin and the target nodes, the operation handling on the target side, and the local completion on the origin side (e.g., receiving an ACK message in the software handled operation). We abbreviate the cost of each portion as  $T_{is}$ ,  $T_{nt}$ ,  $T_{hd}$ , and  $T_{dn}$ , respectively. Moreover, we consider the worst case that the message arrives on the target just when the target process joined a computation task which takes  $E$  time, thus the message cannot be handled until the target computation finishes. Consequently, we can formulate the execution time of the original case as  $T^{original} = T_{is} + T_{nt} + T_{hd} + E + T_{dn}$ .

We then formulate the cost when Casper is involved. Unlike the original case, Casper redirects the message to a ghost process, thus the message can be immediately handled. Therefore, we give the execution time  $T^{csp} = T_{is} + T_{nt} + T_{hd} + T_{dn}$ .

It seems that  $T^{csp}$  should be always smaller than  $T^{original}$ . In practice, the relationship is actually more complicated. An important factor is that, we always have a number of processes and take away only a few of them as ghost processes. Therefore, the situation becomes that a ghost process might receive messages instead of multiple target processes, which lead to the aggregation of  $T_{hd}$ . Let us set the ratio of target processes to a ghost process as  $r$ , thus we finalize the cost with Casper as  $T^{csp} = T_{is} + T_{nt} + rT_{hd} + T_{dn}$ .

Now we can conclude that *when the aggregated message handling cost in Casper ( $rT_{hd}$ ) is less significant than the target computation ( $E$ ), the Casper redirection should improve performance; in contrast, when  $rT_{hd}$  becomes more expensive than  $E$ , then the overaggregation bottleneck appears.*

As an ideal approach to catch this trend, we can measure and compare the overhead of each portion during the execution. However, the MPI standard does not expose a portable interface that allows the external user to insert timers and query the information for an arbitrary internal step. Obtaining the cost of the RMA message handling has to rely on the implementation-specific support of MPI tools interface (MPI\_T).

## 5 ADAPTABLE ASYNCHRONOUS PROGRESS

This section introduces the design and implementation of the adaptive mechanism in Casper.

### 5.1 Discussing Principle of Adaptation

When the overaggregation risk appears, one possible way to make adaptation in Casper is to dynamically assign more resources to asynchronous progress, thus the aggregated handling cost can be eliminated (i.e., reduced  $r$ , the ratio of target processes to a ghost process). However, the obvious drawback is that we need to dynamically transfer some application processes to the “ghost group,” which will result in heavy data re-partitioning, or we have to make expensive process oversubscription through MPI dynamic process<sup>1</sup>.

To minimize the overhead of adaptation, we choose a more lightweight approach. That is, we *disable the redirection when the overaggregation issue becomes significant, so that the operation messages can be handled by sufficient processes; when the overaggregation issue disappears and the delay caused by computation becomes dominant, we re-enable the redirection.*

We notice that this strategy has two potential issues when the redirection becomes disabled. The first is that the messages may suffer from the lack of asynchronous progress again, since they are now handled by the application processes. Given that the disabled setting is needed only when communication becomes dominant, we believe the application processes should be able to make sufficient progress by themselves. The second issue is, the cores dedicated to ghost processes are unutilized. We consider number of dedicated cores is usually small, this limitation is acceptable.

Following the basic adaptive approach, we then implement the mechanism in two directions. We first study a strategy based on *user guidance*. The assumption is that the user knows the performance characteristics of each internal phase; thus the user can request Casper to enable or disable redirection for each particular phase by passing hints to Casper at the beginning of that phase. The simplified solution allows us to concentrate on the important semantics correctness according to the MPI standard. As the second direction, we design a fully automatic strategy based on the idea of *self-profiling*. In the following sections, we describe the design of each strategy separately.

1. The MPI dynamic process concept allows a program to spawn additional MPI processes during execution. However, the support of dynamic process is very limited on HPC systems, for example, the MPICH implementation only supports it on TCP networks.

### 5.2 User-Guided Adaptation

Casper is required to maintain the consistency of message redirection over all application processes in an RMA window. This is because any simultaneous operations issued to the same target in that window must always be handled only by a single process, in order to ensure the *ordering* and *atomicity* of ACCUMULATES. Therefore, we allow the user to reconfigure through an MPI call only when the call guarantees that *all window-wide outstanding operations are completed and all application processes in the window can collectively apply the same change*. Specifically, the reconfiguration can be done either at a window allocation or at a window-wide synchronization call that guarantees both conditions. Therefore, we consider three levels of granularity.

**Global Configuration:** The user can specify a global configuration applied to the entire execution through the environment variable `CSP_ASYNC_CONFIG` with two possible values, `ON` or `OFF`, to enable or disable the redirection in Casper.

**Window Configuration:** Whenever a process allocates a window, the user can pass the MPI info hint `async_config=ON|OFF` to reconfigure for the communication performed through that window.

**Sync-Phase Configuration:** Epochs are the natural synchronization phases. However, not all the epochs can make adaptation. For instance, the synchronization calls in PSCW and the passive target epochs involve only partial processes of the window. Thus, updating in those calls can break the correctness. We can safely reconfigure only in fence. Users can pass the `async_config` info hint for a fence epoch by inserting `MPI_WIN_SET_INFO` before the starting fence call. We require the value of this info to be identical across all processes. Additionally, we propose a new “collective” info hint `symmetric=true|false` that users can pass in `MPI_WIN_SET_INFO`. The `symmetric=true` hint is parsed in Casper meaning that the user ensures all outstanding operations in the window have been completed and all processes have arrived at this call. This allows Casper to trigger an adaptation in `MPI_WIN_SET_INFO`; thus it is useful especially within the passive target epochs. For instance, it can be used after a *flush\_all-barrier*, which commonly exists in passive target programs.

### 5.3 Transparent Profiling-Based Adaptation

We then introduce techniques to enable the fully automatic adaptation. Instead of user guidance, we want to dynamically predict the impact of redirection in Casper for an application phase. This is based on the notion that the application usually performs a similar communication/computation pattern at a certain period of execution time (e.g., in the same solver). Therefore, we can study the performance of a recent period of execution and assume the pattern continues for the upcoming period, thus we can trigger adaptation in the Casper runtime.

The key challenge of this approach is that, in order to ensure portability of Casper, we need obtain the performance information under the constraint that we utilize only the PMPI layer resources. Moreover, we also need to address a second challenge that the dynamically predicted results are applied to all processes in the window consistently. Similar

to the user-guided approach, we relies on a collective synchronization at the window allocation or synchronization calls for updating the redirection. Although this solution does not show problem in the user-guided approach, it can result in failure of adaptation in the profiling-based solution if the timing of the synchronization in the application code does not match the change of performance. For instance, the computation can become heavy after a window is allocated, but there may not be any MPI call that allows Casper to perform the synchronization.

In the remainder of this section, we describe our solution as separated into three key components: the self-profiling and local prediction, a basic window-wide synchronization framework, and a special ghost-offloaded synchronization.

### 5.3.1 Self-Profiling and Local Prediction

Through only the timers inserted in PMPI layer, it is impractical to measure the time of RMA internal portions because they can be processed internally at arbitrary MPI calls. Therefore, instead of focusing on only the message handling cost  $rT_{hd}$ , we try to obtain an approximate relationship between the computation time  $E$  and the overall communication time  $T_{comm}$ . Theoretically, in a specific pattern (e.g., the same operations with the same data size and format), the proportion of  $T_{hd}$  related to the other internal portions should be the same on a system with the same MPI environment. Therefore, if  $rT_{hd} > E$ , we should be able to obtain  $xT_{comm} > E$  where the value of parameter  $x$  is approximately identical for a specific pattern.

This notion allows us to build an approximate prediction model. We define a *communication percentage rate* criterion  $CR = T_{comm}/(T_{comm}+E)$  to indicate the proportion of communication time  $T_{comm}$  in the overall execution time  $(T_{comm}+E)$ . We note that the portion  $E$  includes any time that is spent outside MPI (e.g., computation, I/O). We employ an offline preprocessing step to learn the reference values of  $CR$  that indicate the communication with asynchronous progress redirection takes the same amount of time as redirection disabled for different communication patterns and for different system deployments. We store the  $CR$  reference values for a system and use them as the threshold of real-time adaptation. When the user executes an application, the Casper runtime can perform online profiling and periodically predict the setting based on a corresponding threshold. Below, we describe each step.

**Offline Preprocessing:** In this step, we design benchmarks to simulate various communication patterns and estimate the  $CR$  rate when the condition ( $T_{comm} = T^{original} = T^{csp}$ ) is met. The RMA overhead construction can vary depending on several factors as listed in Table 1. Thus our preprocessing experiments must cover many different sets of those values to reduce the deviation. The second column shows the input matrices generated in our benchmark. We describe the benchmark details and study the results obtained on our test platforms in Section 7.2.

**Online Profiling:** During the application execution, we periodically measure the real-time  $CR$  rate for every period of execution. We insert timers in every MPI function through the PMPI wrapper in Casper to accumulate the time spent in MPI communication  $T_{comm}$  and the overall execution time

TABLE 1  
Important factors for RMA overhead construction.

Factor	Sample inputs used in offline preprocessing
Data size	Data size in bytes.
Datatype	Contiguous: double; Strided: 3D subarray (double).
Operation type	PUT; GET; ACCUMULATE; GET_ACCUMULATE; FOP; CAS.
Blocking pattern	Blocking: flush for every OP; Nonblocking: flush for multiple OPs.
Target pattern (t)	All-to-1-node: Everyone issues OP to the procs on one node; All-to-all: Everyone issues OP to all procs.
Num of procs (n)	Total number of processes.
Num of ghosts (g)	Number of ghost processes on a node.

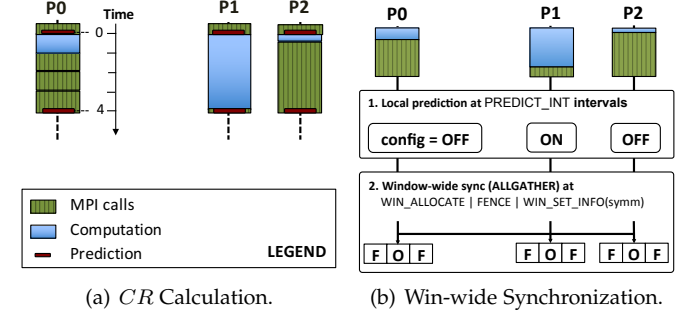


Fig. 2. Self-profiling adaptation.

( $T_{comm}+E$ ). As shown in Figure 2(a), we locally calculate the rate of  $CR$  on every process for the period during two prediction points by using the accumulated times. For instance, a 75% rate is obtained on P0 in the example.

**Local Prediction:** The next step is to locally predict the new configuration for the upcoming period of a process based on the latest real-time  $CR$  rate and the threshold obtained offline. A rate higher than the threshold means that the time the process makes progress in the MPI stack should be sufficiently long to potentially cause operation overaggregation if redirection is enabled (i.e.,  $rT_{hd} > E$ ). Conversely, a rate lower than the threshold indicates a large proportion of computation (i.e.,  $rT_{hd} < E$ ) on this process; enabling asynchronous progress in this case becomes more beneficial. We further use a two-level threshold  $HIGH\_CR$  and  $LOW\_CR$  to avoid frequent fluctuations among large varieties of communication patterns and data characteristics. To ensure a sufficient base of profiling time for every prediction, we also define the threshold  $PREDICT\_INT$  in order to control the interval between two predictions.

### 5.3.2 Window-Wide Synchronization

After the local prediction on every target process, we need to coordinate with the origin side. Thus, the origin process can decide whether to redirect to a ghost process when issuing operation to that target. Similar to the restriction in the user-guided approach, the window-wide synchronization must be done with either a window allocation or special synchronization calls such as `MPI_WIN_FENCE` or `MPI_WIN_SET_INFO` with a symmetric hint.

This component is implemented in a straightforward way such that every process in a window collectively exchanges the last predicted configuration by using `MPI_ALLGATHER` and stores the exchanged data in a local array, as demonstrated in Figure 2(b). Therefore, when the next communication starts, any two operations issued to the same target on the window must both be redirected to the ghost process or issued to the original target process.



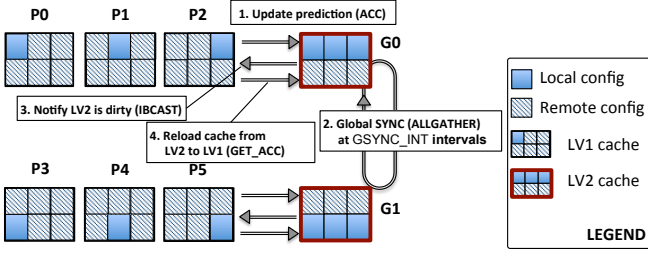


Fig. 3. Ghost-offloaded synchronization.

### 5.3.3 Ghost-Offloaded Synchronization for PUT/GET

Although the window-wide synchronization guarantees semantics correctness, it also limits the adaptation to be valid only at several MPI calls. Unlike `ACCUMULATEs`, `PUT` and `GET` do not require ordering or atomicity. Thus they should be able to be adapted in a more flexible way that does not rely on the existence of special MPI calls in the application code. Therefore, we further investigate a ghost-offloading approach for `PUT` and `GET` where the background ghost processes periodically perform an *asynchronous global synchronization* to exchange the predicted results for all processes. To minimize the overhead, we carefully decouple the local synchronization between the application processes and local ghost processes and the global synchronization among ghost processes, by utilizing a *two-level cache* mechanism as shown in Figure 3. We describe the detailed implementation.

**Two-Level Caches:** Every application process allocates the level-1 cache (denoted by LV1) on its local memory to ensure lightweight querying at frequent `PUT/GET` calls; a shared window is then allocated among the application processes and the first ghost process on every node called “sync ghost” as the level-2 cache (denoted by LV2). Each cache is an array that stores the latest configuration of all application processes in the system and is created only once at MPI initialization. The offset of the configuration for a particular application process is consistent in all processes’ caches. The example shown in Figure 3 demonstrates the cache construction of six application processes distributed on two computing nodes. Every cache is a six-integer-elements array, where the elements from offset 0 to 5 are responsible for the cached configuration from process P0 to P5, respectively.

**Local Updating:** Whenever an application process performs the prediction (see Section 5.3.1) on its local stack, it immediately updates the result to the corresponding element in the LV1 cache. If the value is different from the previous one (e.g., changed from `ON` to `OFF`), the element in the LV2 cache is also updated by issuing an `ACCUMULATE` to ensure atomicity when accessing the shared window.

**Ghost-Offloaded Global Synchronization:** Regardless of the execution on application processes, the sync ghosts perform a global exchange of the LV2 cache at specific intervals defined by the threshold `GSYNC_INT`. Each of the ghost processes sends out the elements corresponding to its local application processes (shown as the solid blue blocks in the LV2 cache in Figure 3) and receives remote values from others through an `MPI_IALLGATHER` collective call.

**Dirty Notification and Reloading:** After the completion of a global synchronization, each sync ghost broadcasts a

dirty notification to its local application processes in an `MPI_IBCAST` call. Each process can then reload its LV1 data from the LV2 cache by using an atomic `GET_ACCUMULATE`.

**Per-Operation Query:** At the issuing of every `PUT` or `GET`, the origin process queries the latest configuration of its target through the LV1 cache, in order to decide whether to redirect that operation.

**Performance Optimization:** The additional synchronization can result in extra overhead in both global synchronization and access to the LV2 caches. Avoiding any unnecessary synchronization is nontrivial. For example, after a window-wide synchronization on most processes in the system such as that with a window allocation call, the synchronized data can be directly updated into the nodes’ LV2 cache and the sync ghosts can skip the upcoming synchronization.

## 6 EXPERIMENTAL ENVIRONMENT

We performed our experiments on two platforms. The first platform is the NERSC Edison Cray XC30 supercomputer<sup>2</sup>, and the second is the Argonne Bebop cluster<sup>3</sup>. Table 2 summarizes the hardware and software configuration of the two platforms. We highlight two important features: (1) Each node of Edison comprises two sockets of the 12-core Intel® Xeon® E5-2695 v2 processor (Ivy Bridge) with Hyper-Threading (HT) enabled, whereas the Bebop node uses two sockets of the 18-core Intel® Xeon® E5-2695 v4 processor (Broadwell) without Hyper Threading; (2) the Cray MPI 7.6.0 on Edison offloads contiguous `PUT` and `GET` operations to hardware and handles other operations in software, whereas the Cray MPI 7.2.1 and the Intel MPI on Bebop handles all operations in software.

For the application case study, we use the large-scale computational chemistry application NWChem version 6.3, with MKL (version 11.2.1 and 2017.3.196 on Edison and Bebop, respectively) as the external math library.

Table 2 summarizes the available asynchronous progress options on the two platforms. We compare the proposed adaptable Casper with original MPI and several static asynchronous progress approaches as listed in Table 3. We note that the Cray MPI 7.2.1 supports a *DMAPP mode* which executes contiguous `PUT` and `GET` in hardware and provides the interrupt-based asynchronous progress for other operations. It is omitted in the evaluate because of its known overhead due to frequent interrupts, and in fact, this mode is deprecated in Cray MPI 7.6.0.

## 7 MICROBENCHMARKS

In this section, we analyze the performance of the adaptive approaches on five microbenchmarks. We use the Cray MPI 7.2.1 as the primary MPI version on Edison.

### 7.1 Overhead Analysis

We first evaluate the overheads caused by the proposed adaptation in comparison with static Casper on Edison. We ran the experiments on a single node with one ghost process, and we vary the total number of application processes.

2. <https://www.nersc.gov/users/computational-systems/edison>

3. <https://www.lcr.anl.gov/systems/resources/bebop>

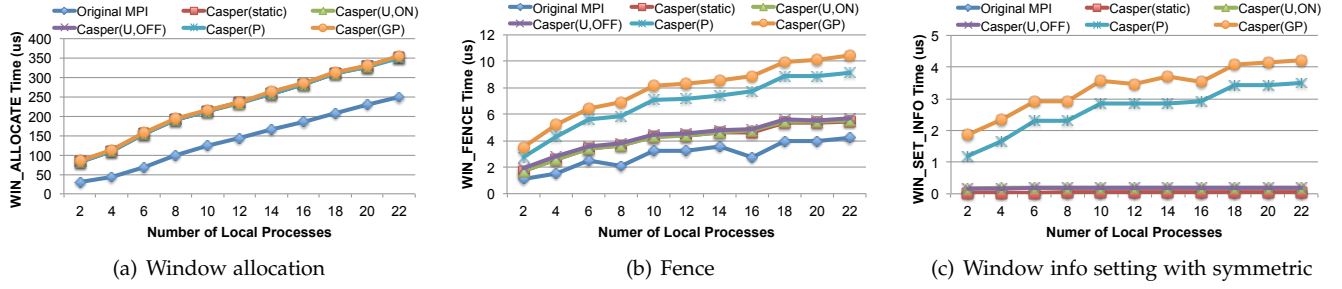


Fig. 4. Adaptation overhead at window collective synchronization on Edison (average of ten runs; error is less than 3%).

TABLE 2

Hardware and software configuration on two experimental platforms.

	CPU	Memory	Interconnect	HT enabled
Edison	2×12-core Ivy Bridge	64 GB DDR3	Cray Aries	Yes
Bebop	2×18-core Broadwell	128 GB DDR4	Omni-Path	No

	MPI	RMA implementation	Default Async	Casper
Edison	Cray MPI 7.2.1	All SW	Thread	static/adaptive
	Cray MPI 7.6.0	HW config PUT/GET	Thread	- *
Bebop	Intel MPI 17.0.4	All SW	Thread	static/adaptive

\*Casper can not execute with the Cray MPI 7.6.0 because of an issue in the MPI implementation which has been reported and being fixed. Thus we use the Cray MPI 7.2.1 as the primary MPI version on Edison.

TABLE 3

Summary of evaluation approaches.

Original MPI	SW RMA without async support.
Casper(static)	SW RMA with Casper static redirection.
Casper(U)	SW RMA with Casper user-guided adaptation.
Casper(GP)	SW RMA with Casper profiling-based adaptation.
Casper(P)	SW RMA with Casper profiling-based adaptation (only window-wide sync).
Thread(D)	SW RMA with thread async, dedicated 50% cores.
Thread(O)*	SW RMA with thread async, oversubscribe core.
Original MPI/HW	Partial HW RMA without async.
Thread(O)/HW	Partial HW RMA with thread async, oversubscribe core.

\*Because the Cray MPI 7.2.1 is not available when we measured the thread(O) approach, we used the *FALLBACK* mode of 7.6.0 which behaves the same as 7.2.1, handling all operations in software.

**Window Allocation:** Figure 4(a) shows the overhead of `MPI_WIN_ALLOCATE` on an application process. Since we focus only on the fence or lockall synchronization in our evaluation, we set the `epoch_type=fence,lockall` info at the allocation call for all the Casper approaches. We also pass the `async_config` hint with either `ON` or `OFF` in the user-guided approaches (denoted by `Casper(U,ON/OFF)`). Both the two adaptive approaches show performance similar to that of static Casper, delivering about 40% to 100% overhead compared with the original MPI implementation. We note that this overhead is because of the internal window creation in Casper, which is unrelated to the proposed adaptation. Although the `Casper(U,OFF)` approach disables the communication redirection, it still suffers from this overhead because we always need to initialize the internal windows in case the user re-enables in the future phase.

**Fence:** Figure 4(b) compares the overhead at `MPI_WIN_FENCE`. The overhead of static Casper compared with the original MPI is due to the passive mode translation in Casper, as discussed in our previous work. The user-guided approaches show performance similar to that of the static version, because they do not involve any additional communication. `Casper(P)` and `Casper(GP)` result in extra overhead because of the additional `MPI_ALLGATHER` that exchanges the

value of predicted new configurations among all processes. Moreover, we see a consistent gap between `Casper(P)` and `Casper(GP)` at close to  $1 \mu s$ ; this is because in `Casper(GP)` the first application process on the node also updates the synchronized data into the LV2 cache (see Section 5.3.3).

**Symmetric Info Setting:** When the user passes the `symmetric=true` hint into the `MPI_WIN_SET_INFO` call, we can also perform adaptation. Figure 4(c) compares the associated overhead. The profiling-based approaches involve additional `MPI_ALLGATHER` communication, thus showing increasing overhead with increasing numbers of processes. The additional  $1 \mu s$  overhead of `Casper(GP)` compared with `Casper(P)` is the same as in the fence experiment.

## 7.2 Offline Estimation for Predictive Threshold

We estimate the thresholds of  $CR$  rate in the profiling-based adaptation through an offline preprocessing step. We use a benchmark set to demonstrate a common RMA communication pattern where every process performs *RMA-computation-RMA* in multiple iterations following with a barrier. The computation part is simulated as busy waiting allowing us to flexibly set different computation costs ( $E$ ). The RMA portion is dynamically generated to cover all the combinations of the factor values as listed in Table 1. For each test, the program automatically adjusts the communication cost by increasing the number of operations until the average execution time with asynchronous progress redirection in static Casper is more expensive than the time with redirection disabled. We record all measured  $CR$  rates that indicate an execution time difference in the range of  $\pm 5\%$ . The benchmark set is available online<sup>4</sup>.

We compare the trend of estimated rates with different sets of factors on both the Edison (with Cray MPI 7.2.1) and the Bebop platforms. We then conclude an approach that calculates the thresholds for later experiments.

**Varying Operation Types, Datatypes, and Blocking Patterns:** We first summarize the trends of estimated  $CR$  rates when we change only one of the following factors: operation types, datatypes, and blocking patterns. Figures 5 (a-c) and Figures 5 (f-h) show the trends on Edison over 192 cores with 1 ghost process per node ( $n=192, g=1$ ), and the trend on Bebop with ( $n=288, g=1$ ), respectively. Roughly speaking, the estimated rates on Edison do not show significant differences for different settings. However, the trends on Bebop show significant diversity. For instance, the strided `ACCUMULATES` delivered much higher  $CR$  rate than the

4. <https://github.com/pmodels/casper-dev/tree/dev-dynamic-sched/preprocess>

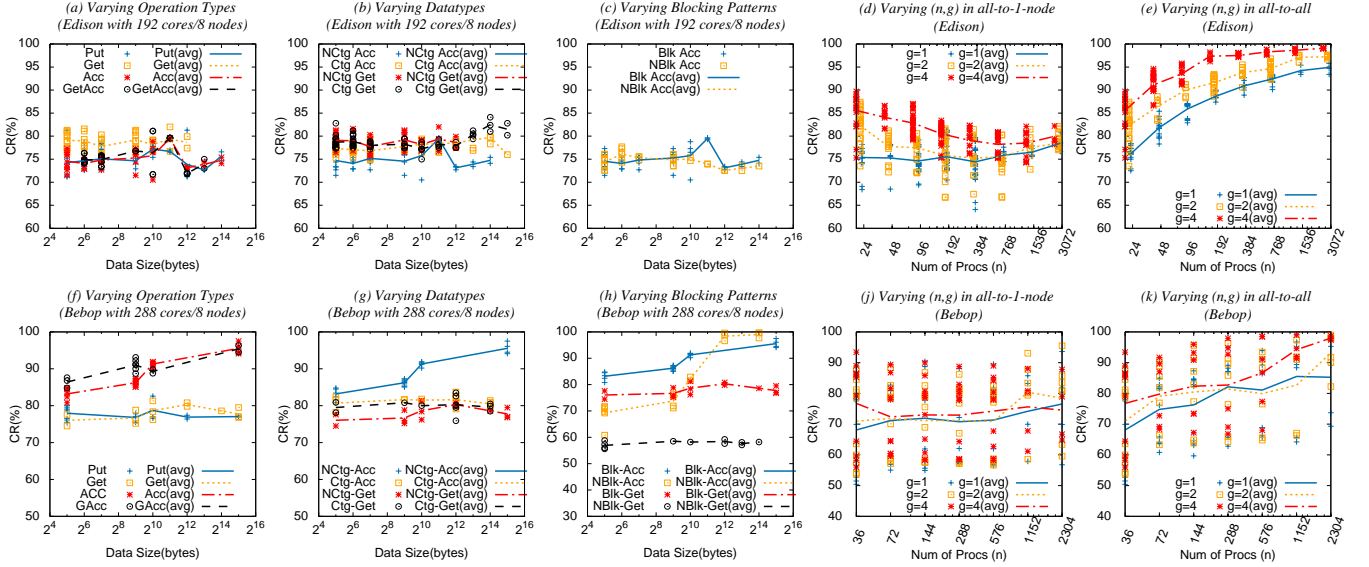


Fig. 5. Analyzing  $CR$  rate. Figures (a)–(e) show the rates measured on Edison, and (f – k) show the rates on Bebop. List of acronyms: Ctg is Contiguous; Std is Strided; Blk is Blocking; and NBk is nonblocking. [update c–e.]

other types, and the change of the blocking patterns with GET also resulted in highly different rates.

Obviously, the diversity of  $CR$  rates can be highly platform dependent. To make effective predictions for applications that often involve a mixture of multiple different patterns, we processed the obtained benchmark results following a simple statistical approach. We calculate the mean value of the results for each combination (e.g., a nonblocking strided ACCUMULATE) called a *basic-mean*. Then we vary the process deployment setting ( $n, g$ ) with different target patterns, and analyze the distribution of these basic-means.

**Varying Process Deployment:** Figures 5 (d) and (j) show the *all-to-1-node* pattern with increasing number of  $n$  and varying number of  $g$  on the test platforms. We notice that the basic-means on Bebop are clearly distributed into two ranges. This is the large diversity we have observed in the previous comparison. Figure 5 (e) and (k) then show the same measurement but with the *all-to-all* target pattern. We make two observations from the figures. The first is that the greater the number of ghost processes, the higher the  $CR$  rate that is required in order to reach a performance bottleneck. This is because the ratio of target processes to a ghost process (abstracted as  $r$  in Section 4) is reduced and thus the bottleneck becomes harder to reach. The second observation is that the larger the number of target processes, the higher the estimated rate. This is because of the reduced proportion of  $T_{hd}$  in the overall communication time.

We then define two strategies that calculate the thresholds. In the first strategy, we calculate the overall boundaries of basic-means for every set of  $(n, t, g)$  on the platform. When the deployment of  $(n, t, g)$  is already known, we directly use the corresponding boundaries. Since  $t$  might change in applications, however, we also define the second strategy: taking the overall boundaries of basic-means for every set of  $(n, g)$ .

### 7.3 Single-Phase Benchmark

Our third set of experiments focuses on the usage of static and adaptive asynchronous progress approaches in

two single-phase microbenchmarks. Specifically, the first one demonstrates a typical computation-intensive pattern (denoted by COMP), and the second demonstrates a communication-intensive pattern (denoted by COMM). Every process performs two times the *100-times[RMA-computation-RMA]-barrier* pattern in an *all-to-all* fashion. In the COMP benchmark, we compute DGEMM in every iteration with a total problem size  $M=N=K=192000$ , and we issue a single *GET-flush* and *ACCUMULATE-flush* in the first and second RMA steps, respectively. In the COMM benchmark, we reduce the total size of DGEMM to  $M=N=K=9600$  and increase the number of operations to 100 at the RMA steps. Every RMA operation carries data with a  $2^3$  3D subarray on the  $8^3$  window region as the target datatype and 8 contiguous double elements as the origin data structure.

We measure each experiment on 192 cores (8 nodes) on Edison. Because the strided operations are handled in software in all the MPI versions, we omit the measurements with original MPI/HW and thread(O)/HW. The original MPI approach uses 24 processes on every node, and we vary the number of ghost processes from 1 to 4 in static Casper. Each serves 23, 22, and 20 application processes per node, denoted by C(1), C(2), and C(4), respectively. We also compare the thread-based approaches (denoted by TH(D) and TH(O)) and the adaptive approaches (denoted by C(U), C(P), and C(GP)) as defined in Table 3. We specify two ghost processes in each adaptive approach. To enable adaptation also in the passive-target communication, we add the `win_set_info` with symmetric hint after the barrier call (see Section 5.2). In C(U), we specify the global `CSP_ASYNC_CONFIG=ON` in the COMP benchmark and set to OFF in the other one. In C(P) and C(GP), we set the  $CR$  thresholds to  $\{89\%, 95\%\}$  according to the offline estimation for  $(n, t, g) = (192, all-to-all, 2)$ , and we empirically set `PREDICT_INT` to 1 second and `GSYNC_INT` to 2 seconds.

Figure 6 compares the performance results. Static Casper always reduces the communication cost significantly in the COMP benchmark (Figure 6(a)) because of asynchronous progress, but it also degrades the computation performance when using more ghost processes (e.g., C(1) achieves the



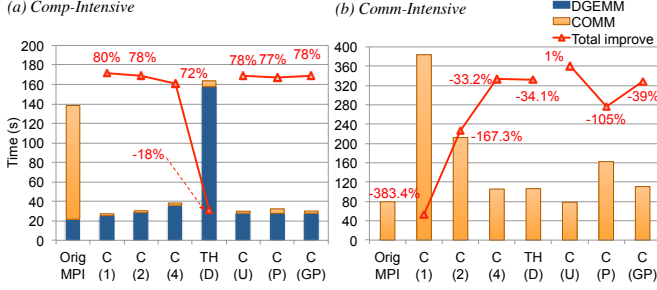


Fig. 6. Performance trend of static asynchronous progress approaches in different communication patterns over 192 cores on Edison (average of three runs). [Add TH(O) and stdev.]

TABLE 4

Expected adaptation of Casper configurations (C), (P), and (GP) in multiphase benchmark involving up to 8 phases.

Approach	1	2	3	4	5	6	7	8
	COMP	COMP	COMM	COMM	COMP	COMP	COMM	COMM
(U)	ON	ON	OFF	OFF	ON	ON	OFF	OFF
(P)	ON	ON	ON	OFF	OFF	ON	ON	OFF
(GP)	ON	ON	ON/OFF	OFF	OFF/ON	ON	ON/OFF	OFF

best performance at 80% improvement compared with the original MPI) because of losing computing power. On the other hand, a small number of ghost processes can result in severe degradation in the COMM benchmark (Figure 6(b)) because of operation overaggregation. For instance, C(1) reports a 383.4% degradation. Such overhead can be reduced by using more ghost processes, and the issue can be completely resolved by disabling the redirection, shown as C(U). However, the profiling-based adaptation, shown as C(P) and C(GP), deliver significant overhead in the COMM benchmark. The reason is that they can adapt only after the first barrier, although C(GP) can partially help the GET operations at an earlier time.

## 7.4 Multiphase Benchmark

Although the user can adjust the setting of static Casper for the execution with a different communication pattern, achieving optimal performance is impossible if a single execution contains both patterns. Our fourth set of experiments focuses on such a multiphase benchmark. The benchmark contains two sequential windows, each consisting of both a heavy-computing period and a heavy-communicating period (combination of the two single benchmarks in the third set of experiments). For convenience, let us call every 100-times[RMA-computation-RMA]-barrier a phase.

We use two ghost processes in all Casper approaches. In Casper(U), we set the user hint `async_config=ON` at each window allocation call for the upcoming COMP phases, and we set `async_config=OFF` through `win_set_info` in front of the third and the seventh phases for the next COMM phases. The configuration of Casper(P) and Casper(GP) remains the same as that in the preceding set of experiments.

As listed in Table 4, the three adaptive approaches can result in different reconfigurations in every internal phase. To be specific, Casper(U) can deliver the most precise adaptation that enables redirection (ON) in every COMP phase and disables it (OFF) in every COMM phase. Casper(P), however, cannot promptly adapt to the third, fifth, and seventh phases because it cannot apply the new predicted

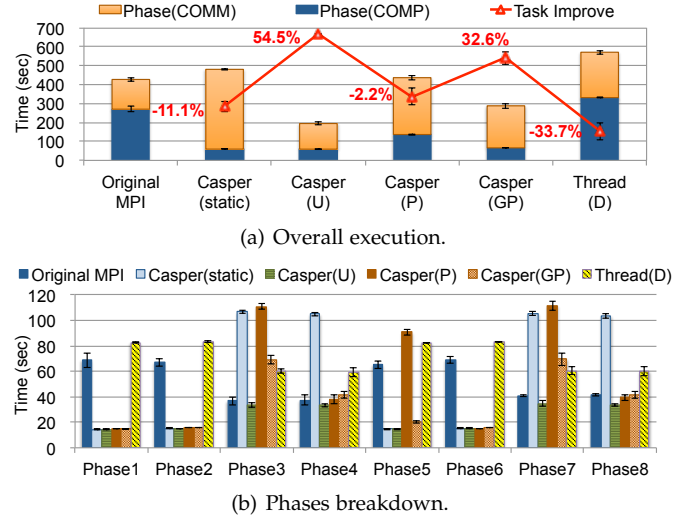


Fig. 7. Comparison of asynchronous progress approaches in the multiphase benchmark over 192 cores on Edison (average of three runs). [add TH(O)]

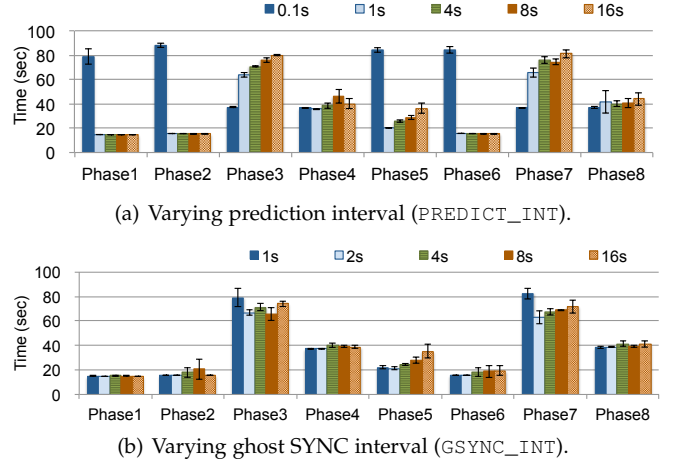


Fig. 8. Comparison of varying adapting intervals in the multiphase benchmark over 192 cores on Edison (average of three runs).

results until it reaches `win_set_info`. Casper(GP) partially addresses this issue; for example, the ghost-offloaded synchronization disables redirection at the third and seventh phases only for GETs.

Figure 7 shows the results. Although static Casper can reduce the cost of the COMP phases, it also degrades the other phases that perform intensive communication, resulting in an 11.1% degradation compared with the original MPI. As expected, Casper(U) achieves the greatest improvement at 54.5%; Casper(P) cannot provide appropriate adaptation at Phases 3, 5, and 7, as shown in Figure 7(b), resulting in a 2.2% degradation; Casper(GP) reduces the overhead at those three phases by adapting GETs.

## 7.5 Varying Adaptation Intervals

For our fifth set of experiments, we use the same multiphase benchmark to observe the impact of two interval thresholds in the profiling-based approaches: `PREDICT_INT` (see Section 5.3.1) and `GSYNC_INT` (see Section 5.3.3).

Figure 8(a) compares the per-phase costs with varying `PREDICT_INT` in Casper(GP) with a fixed `GSYNC_INT` at 2 seconds. We omit the graph of Casper(P) because of space limitations. We notice that the smallest interval, 0.1 seconds,

can result in imprecise adaptation especially in the COMP phases (i.e., Phases 1, 2, 5, and 6 in Casper(GP)). The reason is that the fragment of the executed period is so short that includes only the last few MPI calls, thus the profiling data cannot represent the heavy computing characteristics. With increasing internal time, Casper(GP) shows increasing overhead in Phases 3, 5, and 7, where processes reconfigure for heavy communication, because of the delay in prediction.

Figure 8(b) compares the per-phase overhead with varying `GSYNC_INT` and a fixed `PREDICT_INT` at 1 second. It indicates a visible overhead in the COMM phases that are contiguous to the preceding COMP phases (i.e., Phases 3 and 7) when a 1-second interval is set, because of the frequent reloading executed on every application process. Increasing the interval can lead to delays in adaptation especially in the heavy computing phase (i.e., Phase 5).

## 8 CASE STUDY: CHEMISTRY APPLICATION

In previous work we evaluated the NWChem application with static asynchronous progress by focusing on particular internal phases of the CCSD(T) method [13], [16]. Here we focus on the overall multiphase execution.

**NWChem Background:** NWChem [1] is a widely used computational chemistry application suite [17], [18]. NWChem is developed on top of the Global Arrays [3] toolkit over MPI RMA [19]. A typical *get-compute-update* mode is widely used in all the internal phases of NWChem, which every process essentially performs by varying the size of matrix-matrix multiplication for multidimensional tensor contraction by coordinating with others through RMA GET/ACCUMULATE operations. Furthermore, *NXTASK* is the generic task-scheduling component that assigns the “owner” for subdomain computing tasks. It is implemented as a single *FETCH\_AND\_OP* operation (denoted by FOP).

We note that most of the RMA operations in NWChem exchange the subblocks of the global matrix. The subblock data is represented as a strided subarray in MPI. Thus, the hardware-offloaded PUT/GET can not help performance.

**Experimental Setup:** We insert a `win_set_info` call with symmetric info at the Global Arrays `GA_SYNC` call<sup>5</sup>, since its semantics guarantee the completion of all outstanding operations on all processes. We note that Casper(GP) requires three kinds of predefined thresholds: `LOW|HIGH_CR`, `PREDICT_INT`, and `GSYNC_INT`. We use the estimated *CR* thresholds from the results in Section 7.2 following the second strategy. With regard to the interval thresholds, we decide the value according to each task’s execution time within the original MPI.

We choose two widely used modules of NWChem in our case study, the single-phase density functional theory (DFT) and the multiphase CCSD(T).

### 8.1 Single-Phase DFT

Density functional theory is one of the most broadly used methods in NWChem. It provides a good mix of efficiency and accuracy to investigate the structural and electronic properties of atoms and molecules. It contains only a single

5. `GA_SYNC` internally calls `MPI flush_all` on all processes followed by a barrier.

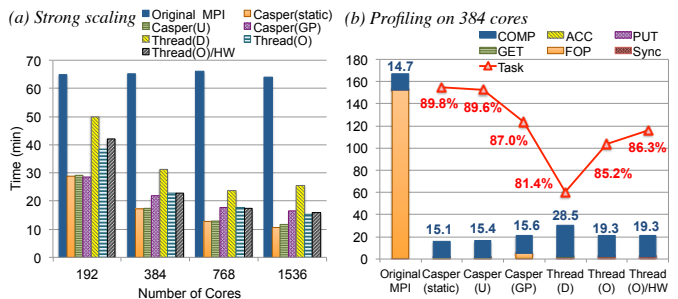


Fig. 9. Single-phase DFT task for C240 with asynchronous progress on Edison. Original MPI/HW cannot complete in 5 hours, thus it is omitted. All Casper approaches use 1 ghost process per node. COMP values are shown in (b).

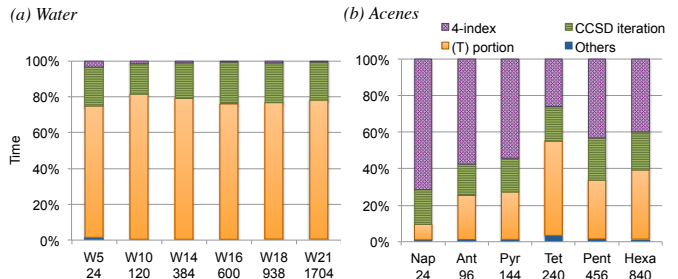


Fig. 10. Internal phases in multiphase CCSD(T) with pVDZ on Edison.

internal phase in the implementation, which follows the *get-compute-update* mode and utilizes *NXTASK* task scheduling.

We evaluate the DFT calculation for Carbon 240 (denoted by C240) with the 6-31G\* basis set. We use one ghost process in all Casper approaches; set `CSP_ASYNC_CONFIG=ON` in Casper(U); and in Casper(GP) set `PREDICT_INT` to 2 seconds, `GSYNC_INT` to 120 seconds, and a *CR* rate range {75%,90%}, {75%,90%}, {80%,90%}, and {85%,90%} for 192, 384, 768, and 1,536 cores, respectively.

Figure 9(a) compares the strong scaling of both static and adaptive asynchronous progress approaches over a varying number of system cores. The original MPI does not scale because of the significant delay in the blocking FOP operations in *NXTASK* as shown in Figure 9(b). All the asynchronous progress approaches can eliminate such overhead; however, the thread-based approaches are not as efficient as the Casper approaches because of increased overhead in computation. We compare the static and adaptive approaches in Casper. The static approach is clearly the best solution for the single-phase DFT. Casper(U) gives similar performance, but Casper(GP) shows visible communication overhead primarily because of the extra synchronization and prediction error.

### 8.2 Multiphase CCSD(T)

The coupled cluster theory is one of the most popular approaches in quantum chemistry for solving electron correlation in atoms and molecules with arbitrary accuracy requirements. The “gold standard” *coupled cluster with singles and doubles and perturbative triples* method, known as CCSD(T), is one of the most accurate CC methods applicable to large molecules to date.

The CCSD(T) method comprises four internal phases: self-consistent field (SCF), four-index transformation (4-index), CCSD iteration, and the noniterative (T) portion [20].

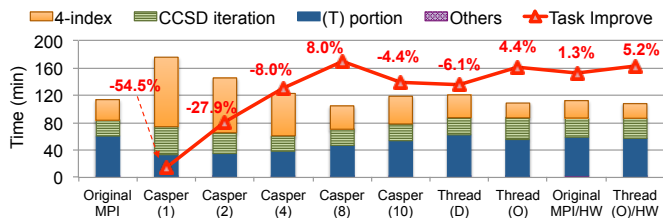


Fig. 11. Trade-off in NWChem CCSD(T) for Tet-aug-cc-pVDZ with static asynchronous progress over 240 cores on Edison. Casper( $g$ ) denotes Casper with  $g$  ghost processes ( $g=1;2;4;8;10$ ). The “Task Improve” rate is calculated based on Original MPI.

The overhead proportion among these phases can vary in particular molecular problems. Figure 10 compares the overhead breakdown of two sets of molecular problems with the original MPI: the water molecule ( $H_2O$ ) $_n$  problems ( $n = 5; 10; 14; 16; 18; 21$ , denoted as  $W_n$ ), with the cc-pVDZ basis set, and the Acenes problems, including naphthalene, anthracene, tetracene, pentacene, and hexacene (denoted Nap, Ant, Tet, Pent, and Hexa, respectively) with the aug-cc-pVDZ basis set. Each problem is measured over the appropriate number of cores fitting its memory requirement, as listed below the x-axis.

In all the water problems, the (T) portion consistently dominates the cost of the entire task by close to 80%, and the CCSD iteration takes the other 20%; the remaining phases represent less than 2% of the time. The Acenes series shows a different trend in each problem, where the (T) portion indicates only a 52% cost in Tet and an even lower proportion in others. Instead, the 4-index contributes more overhead, representing 40–71% of the time. We note that the SCF always takes less than 1% of the cost; thus it is merged into the “Others” portion for simplicity.

### 8.2.1 Analysis with Static Asynchronous Progress

Next we looked into the performance issue of static asynchronous progress. We chose two problem types: large *W21* molecule over 1,704 cores and *Tet* over 240 cores. We compared the performance impact on each internal phase by utilizing the static Casper and thread-based approaches. We used the same total number of cores on every computing node in all approaches, some of which are dedicated to asynchronous ghost processes/threads.

**Trade-Off in Overall Execution:** Figure 11 shows the task execution time of the *Tet* problem. Casper(1) delivers the maximum improvement in the (T) portion by close to 50%, but it also leads to more expensive CCSD iteration and 4-index. With increasing numbers of ghost processes such overhead can be gradually decreased, but the overhead of the (T) portion increases. Thread(D) follows the same trend, because it occupies half of the computing cores. The Thread(O) approaches do not perform better because of core oversubscription. As a result, only an 8% improvement is achievable with 8 ghost processes.

The internal phases of *W21* indicate trends similar to those observed in *Tet*. Static Casper delivers the best improvement for the overall execution at 28% by using 2 ghost processes, because the deduction of the degradations in other phases can be reimbursed by the improvement in the (T) portion, which dominates the entire cost by 80%.

Having studied the overall performance trend, we then analyzed each specific internal phase. Since we observed similar trends in each phase in both the *W21* and *Tet* problems, we have omitted the results of *W21*.

**Four-Index Phase:** Figure 12(a) shows the overhead of the computation and RMA operations in the 4-index phase. The degradation with small numbers of ghost processes is caused mostly by ACCUMULATES, which degrade performance by 40x with one ghost process; but the GET portion, which dominates the cost of the 4-index in the original MPI, can benefit from the redirection in Casper. After careful code reading and profiling, we confirmed that this difference is due to the different target patterns executed in these operations. To be specific, all ACCUMULATES are issued as the *all-to-1-node* pattern described in Section 7.2. GETs, on the other hand, are issued following the *all-to-all* pattern.

**CCSD Iteration Phase:** Figure 12(b) shows the profiling of the CCSD iteration phase. Different from the overhead construction in the 4-index, the numerous *all-to-all* GETs dominate the execution time by close to 80%, and the DGEMM computation (shown as COMP) takes less than 10%. Such intensive communication can rarely benefit from asynchronous progress if the operations are aggregated to only a few ghost processes. Thus, 4 ghost processes are required in order to balance the overaggregation overhead.

**(T) Portion Phase:** Figure 12(c) shows the overhead profiles of the noniterative phase: (T). With the original MPI, the heavy computation takes 30 minutes, and GETs dominate the other half of the cost. The overhead of GETs clearly indicates the delay caused by lack of asynchronous progress. All the static approaches can asynchronously complete GET operations; thus such overhead can be eliminated. With more cores dedicated to ghost processes or threads, however, the computation resources are also reduced, resulting in significant degradation in the computation. The TH(O) approaches show similar degradation because of core oversubscription.

### 8.2.2 Dynamic Adaptation

We next evaluate the dynamic adaptive strategies on both the Edison and Bebop platforms.

**Weak and Strong Scaling:** We evaluate Casper(U) and Casper(GP) in both weak and strong scaling of the Acenes problems, by comparing them with the original MPI and the static approaches studied in the preceding section. In both the static and adaptable Casper approaches, we specify two ghost processes per node. In Casper(U), we specify {ON, OFF, OFF, ON} as the value of global `CSP_ASYNC_CONFIG` and the `async_config` infos passed to three internal phases: 4-index, CCSD iteration, and (T) portion, respectively. In Casper(GP), we specify the thresholds as listed in Table 5.

Figure 13 shows the results on Edison. In the weak scaling graph, we increase the problem sizes and numbers of cores. Static Casper shows significant degradation in the execution of all multi-node problems. The thread-based approaches still cannot achieve consistent improvement. The adaptable Casper(U) and Casper(GP), on the other hand, consistently improve the execution for each problem type by up to 23.2% at Hexa and 16.3% at Tet, respectively. In the strong scaling, the static approaches show consistent degradation in performance, while Casper(U) and Casper(GP)



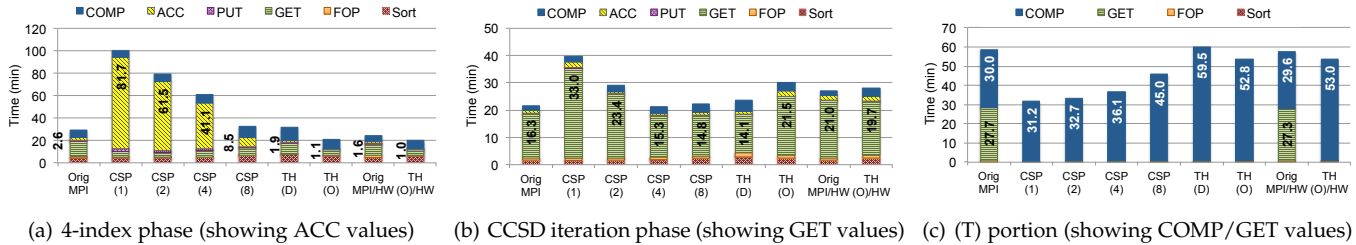


Fig. 12. Profiling multiphase CCSD(T) for Tet-aug-cc-pVDZ with static asynchronous progress over 240 cores on Edison.

TABLE 5  
Environment variable setting for NWChem CCSD(T) with profiling-based adaptation.

Edison: {CR(%)}, PREDICT_INT(s), GSYNC_INT(s)	Bebop
Nap/24: {81,88}, 60, 2	Nap/36: {50,95}, 60, 2
Ant/96: {78,85}, 120, 2	Ant/72: {50,90}, 120, 2
Pyr/144: {75,90}, 120, 2	Pyr/108: {50,90}, 120, 2
Pent/456: {80,90}, 240, 2	Tet/288: {60,95}, 240, 2
Hexa/840: {85,90}, 240, 2	Pent/576: {60,95}, 240, 2
-	Hexa/1008: {65,95}, 240, 2

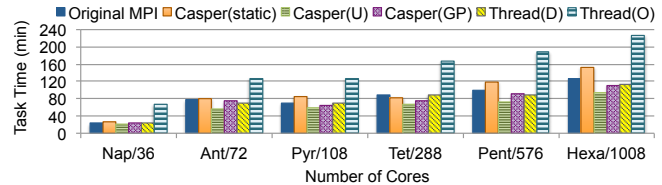
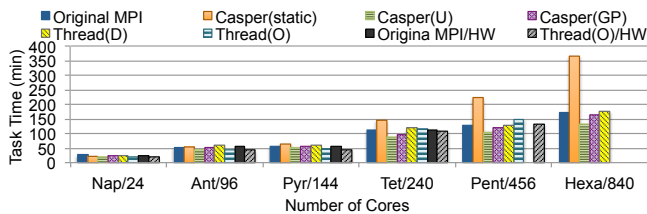
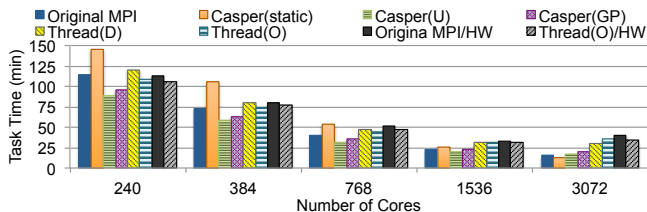


Fig. 14. Weak scaling for dynamic adaptation with NWChem CCSD(T) on Bebop (using 2 ghost processes in all Casper approaches).



(a) Weak scaling on Acenes series



(b) Strong scaling on Tet-aug-cc-pVDZ

Fig. 13. Dynamic adaptation with NWChem CCSD(T) on Edison (using 2 ghost processes in all Casper approaches). The unfinished bars in (a) are because of out of memory error.

can resolve such inefficiency. Casper(U) delivers the best performance by utilizing user hints, achieving up to 21.8% speedup; Casper(GP) provides a fully automatic solution based on self-profiling and prediction; and it improves performance up to 16.3%. When scaling to 3,072 cores, static Casper becomes the best option because the 4-index becomes dominated by numerous *all-to-all* GETs that benefit from asynchronous progress with only two ghost processes.

Figure 14 shows the weak scaling on Bebop. Thread(O) delivers significant overhead because it oversubscribes without HT. Casper(GP) shows higher overhead than the results on Edison because of the overestimated range of CR thresholds. For instance, the 50% `LOW_CR` used in Ant was generated by the nonblocking Get patterns in preprocessing (see Figure 5 (h)), which is never used in the application. This caused the delay of adaptation in (T).

**Internal Phase Overhead:** We choose the Tet problem with 240 cores of Edison as the base of our profiling. We first compare the overhead of each phase with different approaches. We also add the profiling-based adaptation without ghost-

offloaded synchronization, denoted by CSP(P) in this experiment. As shown in Figure 15(a), both CSP(U) and CSP(GP) can correctly resolve the overhead in 4-index and improve the performance for the (T) portion, but Casper(P) cannot improve the overall performance because of the expensive overhead in the (T) portion.

We then compare the overhead distributed in the 4-index phases. Figure 15(b) clearly indicates that all the adaptive approaches can resolve the overhead caused by overaggregated ACCUMULATES. With regard to the (T) portion, as shown in Figure 15(c), Casper(U) behaves the same as static Casper because it immediately re-enables the redirection at the beginning of (T). Casper(P), on the other hand, cannot reduce the GET overhead because no synchronization call exists in the application code. Casper(GP) eliminates the overhead of GET by re-enabling asynchronous progress through ghost-offloaded synchronization. In addition, we notice an overhead of close to 3 minutes in GET and FOP portions in Casper(GP) compared with Casper(U), because of the interval set for ghost synchronization.

Another observation we made is that, although Casper(GP) predicts on each process separately, the majority of the processes always make the same decision (e.g., 99% of the processes disabled redirection in 4-index, and all of them enabled in (T)).

## 9 RELATED WORK

Casper is a portable, process-based asynchronous progress model for MPI one-sided communication. In this paper, we focus on the use of this model in multiphase applications, and we propose several adaptive methods to dynamically reconfigure asynchronous progress to resolve operation aggregation imbalance. We divide the related work into two broad topics: communication asynchronous progress and dynamic adaptation for load balancing.

**Communication Asynchronous Progress:** *Thread-based asynchronous progress* is considered the most common approach for supporting software progress and is found in many MPI implementations such as MPICH and its derivatives [6] [7] [8]. This model allows every MPI process to utilize a background thread to asynchronously handle the



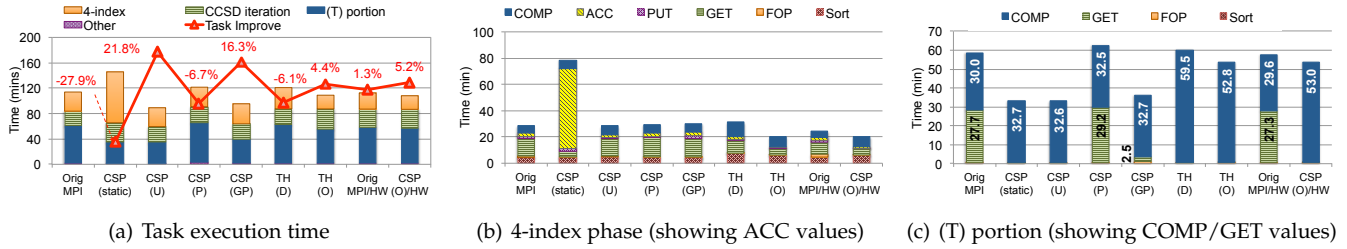


Fig. 15. Profiling multiphase CCSD(T) for Tet-aug-cc-pVDZ with dynamic adaptation over 240 cores on Edison.

incoming messages from other processes. While being a generic approach for various MPI communication models, this approach also raises the restriction that a background thread can make progress only for the process that spawned it. Thus it has to deploy at least as many background threads as MPI processes on every computing node. Consequently, the user must choose either to dedicate half of the computing resources or to involve expensive core oversubscription. Furthermore, this model forces MPI runtimes to support multithreaded safety, which may result in further overheads because of thread synchronization [10].

PIOMan [21] is a multithreaded communication engine supporting thread-based asynchronous progress. It divides rendezvous handshakes into multiple tasks and offloads them to background threads running only on idle cores. This approach, however, also suffers from a non-negligible overhead derived from the necessary multithreaded safety [22].

Vaidyanathan et al. [23] contributed an approach for asynchronous progress in the “MPI+X” model by utilizing a dedicated thread together with a lock-free command queue. The “MPI+X” model often utilizes multiple threads over multi- or many-core systems to parallelize computation and employs only a single MPI process per node for internode communication. Thus, only a single asynchronous thread is required per node.

The other well-known approach in the MPI community is the *interrupt-based asynchronous progress*, which has been supported on both Cray [24] and IBM systems [25] [26]. The core concept of this approach assumes that all processes are busy in external computation, thus utilizing a system interrupt to awaken the kernel thread to asynchronously complete incoming messages. The design is straightforward; however, the implementation often relies on a platform-specific lightweight interrupt engine; otherwise, severe performance degradation might occur because of frequently issued interrupts [20].

Supporting asynchronous progress is an essential task for using the portable MPI in other runtime systems. Daily et al. [27] proposed the approach to build the PGAS ComEX runtime on top of MPI two-sided model, and designed a progress rank engine in ComEX that splits the MPI world communicator and uses a subset of processes to help communication progress.

**Dynamic adaptation and load balancing:** Dynamic adaptation is a popular approach to dynamically balance irregular workloads or adapt heterogeneous execution environment and communication methods in both application and runtime systems. Flaherty et al. [28] and Biswas et al. [29] introduced their dynamic load balancer approaches for irregular workloads in mesh applications by repartitioning domains.

As examples of runtime level adaptation, Bhandarkar et al. [30] proposed an MPI implementation on top of the Charm++ environment that provides support for processor visualization and balances the workloads by dynamically measuring idle time or through user hints. Some researchers [31] [32] concentrated on generic autonomic runtime management for workloads on distributed memory systems by implementing dedicated system modules.

Different from these works, we propose adaptive strategies in a portable MPI asynchronous progress library to resolve the operation overaggregation imbalance when providing asynchronous progress.

## 10 CONCLUSION AND FUTURE WORK

Casper is a portable process-based asynchronous progress model for MPI RMA on multi- and many-core architectures. Our previous work presented the core framework of Casper that sets aside a small number of cores as background ghost processes and redirects the user RMA operations targeting an application process to the bound ghost process, thus enabling asynchronous completion of RMA communication. This redirection-based design, however, might also result in operation overaggregation bottlenecks because of the limited progress resources, especially when communication becomes dominant. Therefore, a performance trade-off has to be made in multiphase applications.

In this paper, we proposed the adaptive mechanism for Casper that resolved the overaggregation issue by disabling operation redirection in communication-intensive phases without affecting the benefit of asynchronous progress in other computation-heavy phases.

We chose an approximate prediction model in the adaptation to detect performance in order to maintain the portability of Casper. This model relies on the offline prepossessing to sample the system performance matrices from a large set of benchmarks. However, it might be imprecise if the pattern of an application phase is not covered or large performance diversity exists among different patterns such as the trends observed on the Bebop system. Moreover, the real-time prediction can be further affected by several factors such as system noise or temporary network delay. Although we usually expect such noises to be small on HPC supercomputers, we should give careful consideration. Therefore, we plan to optimize the prediction model based on dynamic heuristic in future work.

## ACKNOWLEDGMENTS

This material was based upon work supported by the U.S. Dept. of Energy, Office of Science, Advanced Scientific Computing Research (SC-21), under contract DE-AC02-06CH11357. The experimental resources for this paper were

provided by the National Energy Research Scientific Computing Center (NERSC) on the Edison Cray XC30 supercomputer and by the Laboratory Computing Resource Center on the Bebop cluster at Argonne National Laboratory. Antonio J. Peña is co-financed by the Spanish Ministry of Economy and Competitiveness under Juan de la Cierva fellowship number IJCI-2015-23266.

## REFERENCES

- [1] E. J. Bylaska et al., “NWChem, A Computational Chemistry Package for Parallel Computers, Version 6.3,” 2013.
- [2] T. Shiozaki, “BAGEL: Brilliantly Advanced General Electronic-structure Library,” *Wiley Interdisciplinary Reviews: Computational Molecular Science*. [Online]. Available: <http://dx.doi.org/10.1002/wcms.1331>
- [3] J. Nieplocha, R. J. Harrison, and R. J. Littlefield, “Global Arrays: A Portable “Shared-Memory” Programming Model for Distributed Memory Computers,” in *ACM/IEEE conference on Supercomputing*, 1994.
- [4] H. E. Trease, “Code Development for NWGrid/NWPhys,” *Laboratory Directed Research and Development Annual Report-Fiscal Year 2000*, p. 103, 2001.
- [5] W. Jiang, J. Liu, H.-W. Jin, D. Panda, D. Buntinas, R. Thakur, and W. Gropp, “Efficient Implementation of MPI-2 Passive One-Sided Communication on InfiniBand Clusters,” in *Euro PVM/MPI*, ser. Lecture Notes in Computer Science, 2004, vol. 3241, pp. 68–76.
- [6] Argonne National Laboratory, “MPICH — High-Performance Portable MPI,” <http://www.mpich.org>, 2014.
- [7] The Ohio State University, “MVAPICH: MPI over InfiniBand, 10GigE/iWARP and RoCE,” <http://mvapich.cse.ohio-state.edu>, 2014.
- [8] Intel Corporation, “Intel MPI Library,” <http://software.intel.com/en-us/intel-mpi-library>, 2014.
- [9] Cray Inc., “Cray XC Series Network,” <https://www.cray.com/sites/default/files/resources/CrayXCNetwork.pdf>, Cray Inc., Tech. Rep., 2004.
- [10] W. Gropp and R. Thakur, “Thread-Safety in an MPI Implementation: Requirements and Analysis,” *Parallel Comput.*, vol. 33, no. 9, pp. 595–604, 2007.
- [11] K. Kandalla, P. Mendygral, N. Radcliffe, B. Cernohous, D. Knaak, K. McMahon, and M. Pagel, “Optimizing Cray MPI and SHMEM Software Stacks for Cray-XC Supercomputers based on Intel KNL Processors,” *Proceedings of 2016 Cray User Group (CUG)*, 2016.
- [12] M. Gilge, *IBM System Blue Gene Solution: Blue Gene/P Application Development*. IBM, Jun. 2013.
- [13] M. Si, A. J. Peña, J. Hammond, P. Balaji, M. Takagi, and Y. Ishikawa, “Casper: An Asynchronous Progress Model for MPI RMA on Many-Core Architectures,” in *2015 IEEE International Parallel and Distributed Processing Symposium*, May 2015, pp. 665–676.
- [14] “MPI: A Message-Passing Interface Standard,” <http://mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf>, Sep. 2012.
- [15] T. Hoefler, J. Dinan, R. Thakur, B. Barrett, P. Balaji, W. Gropp, and K. Underwood, “Remote Memory Access Programming in MPI-3,” Argonne National Laboratory, Tech. Rep., 2013.
- [16] M. Si, A. J. Peña, J. Hammond, P. Balaji, and Y. Ishikawa, “Scaling NWChem with Efficient and Portable Asynchronous Communication in MPI RMA,” in *The 15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, May 2015, pp. 811–816.
- [17] S. Hirata, “Tensor Contraction Engine: Abstraction and Automated Parallel Implementation of Configuration-Interaction, Coupled-Cluster, and Many-Body Perturbation Theories,” *J. Phys. Chem. A*, vol. 107, pp. 9887–9897, 2003.
- [18] E. Aprà et al., “Liquid Water: Obtaining the Right Answer for the Right Reasons,” in *SC*, 2009.
- [19] J. S. Dinan, P. Balaji, J. R. Hammond, S. Krishnamoorthy, and V. Tipparaju, “Supporting the Global Arrays PGAS Model Using MPI One-Sided Communication,” in *2012 IEEE International Parallel and Distributed Processing Symposium*, May 2012.
- [20] J. R. Hammond, S. Krishnamoorthy, S. Shende, N. A. Romero, and A. D. Malony, “Performance Characterization of Global Address Space Applications: A Case Study with NWChem,” *Concurrency and Computation: Practice and Experience*, vol. 24, no. 2, pp. 135–154, 2012.
- [21] F. Trahay and A. Denis, “A Scalable and Generic Task Scheduling System for Communication Libraries,” in *IEEE Cluster*, Sep. 2009.
- [22] F. Trahay, É. Brunet, and A. Denis, “An Analysis of the Impact of Multi-Threading on Communication Performance,” in *9th Workshop on Communication Architecture for Clusters (CAC)*, May 2009.
- [23] K. Vaidyanathan, D. D. Kalamkar, K. Pamnany, J. R. Hammond, P. Balaji, D. Das, J. Park, and B. Jo, “Improving Concurrency and Asynchrony in Multithreaded MPI Applications Using Software Offloading,” in *The International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, Nov. 2015, pp. 1–12.
- [24] H. Pritchard, D. Roweth, D. Henseler, and P. Cassella, “Leveraging the Cray Linux Environment Core Specialization Feature to Realize MPI Asynchronous Progress on Cray XE Systems,” in *Proceedings of the Cray User’s Group Meeting (CUG)*, 2012.
- [25] S. Kumar, G. Doza, G. Almasi, P. Heidelberger, D. Chen, M. E. Giampapa, M. Blocksome, A. Faraj, J. Parker, J. Ratterman, B. Smith, and C. J. Archer, “The Deep Computing Messaging Framework: Generalized Scalable Message Passing on the Blue Gene/P Supercomputer,” in *Proceedings of the 22nd Annual International Conference on Supercomputing*, 2008, pp. 94–103.
- [26] S. Kumar, Y. Sun, and L. V. Kale, “Acceleration of an Asynchronous Message Driven Programming Paradigm on IBM Blue Gene/Q,” in *Proceedings of the IEEE 27th International Symposium on Parallel and Distributed Processing*, 2013, pp. 689–699.
- [27] J. Daily, A. Vishnu, B. Palmer, H. van Dam, and D. Kerbyson, “On the Suitability of MPI as a PGAS Runtime,” in *21st International Conference on High Performance Computing (HiPC)*, Dec. 2014.
- [28] J. E. Flaherty, R. M. Loy, C. zturan, M. S. Shephard, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz, “Parallel Structures and Dynamic Load Balancing for Adaptive Finite Element Computation,” *Applied Numerical Mathematics*, vol. 26, no. 12, pp. 241–263, 1998.
- [29] R. Biswas, S. K. Das, D. J. Harvey, and L. Oliker, “Parallel Dynamic Load Balancing Strategies for Adaptive Irregular Applications,” *Applied Mathematical Modelling*, vol. 25, no. 2, pp. 109–122, 2000.
- [30] M. A. Bhandarkar, L. V. Kalé, E. de Sturler, and J. Hoeflinger, “Adaptive Load Balancing for MPI Programs,” in *Proceedings of the International Conference on Computational Science-Part II*, 2001.
- [31] J. Yang, H. Chen, S. Hariri, and M. Parashar, “Autonomic Runtime Manager for Adaptive Distributed Applications,” in *Proceedings of the 14th IEEE International Symposium on High Performance Distributed Computing (HPDC)*, 2005.
- [32] T. Estrada and M. Taufer, “On the Effectiveness of Application-aware Self-Management for Scientific Discovery in Volunteer Computing Systems,” in *The International Conference on High Performance Computing, Networking, Storage and Analysis (SC)*, 2012.