## Quantifying the Performance Benefits of Partitioned Communication in MPI

Thomas Gillis

Argonne National Laboratory Lemont, Illinois, USA Ken Raffenetti

Argonne National Laboratory Lemont, Illinois, USA Hui Zhou

Argonne National Laboratory Lemont, Illinois, USA

## Yanfei Guo

Argonne National Laboratory

Lemont, Illinois, USA

## Lemont, Illinois, USA

**Rajeev** Thakur

Argonne National Laboratory

## ABSTRACT

Partitioned communication was introduced in MPI 4.0 as a userfriendly interface to support pipelined communication patterns, particularly common in the context of MPI+threads. It provides the user with the ability to divide a global buffer into smaller independent chunks, called *partitions*, which can then be communicated independently. In this work we first model the performance gain that can be expected when using partitioned communication. Next, we describe the improvements we made to MPICH to enable those gains and provide a high-quality implementation of MPI partitioned communication. We then evaluate partitioned communication in various common use cases and assess the performance in comparison with other MPI point-to-point and one-sided approaches. Specifically, we first investigate two scenarios commonly encountered for small partition sizes in a multithreaded environment: thread contention and overhead of using many partitions. We propose two solutions to alleviate the measured penalty and demonstrate their use. We then focus on large messages and the gain obtained when exploiting the delay resulting from computations or load imbalance. We conclude with our perspectives on the benefits of partitioned communication and the various results obtained.

#### **KEYWORDS**

distributed systems, MPI, partitioned communication

#### **ACM Reference Format:**

Thomas Gillis, Ken Raffenetti, Hui Zhou, Yanfei Guo, and Rajeev Thakur. 2023. Quantifying the Performance Benefits of Partitioned Communication in MPI. In *52nd International Conference on Parallel Processing (ICPP 2023), August 7–10, 2023, Salt Lake City, UT, USA.* ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3605573.3605599

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0843-5/23/08...\$15.00

https://doi.org/10.1145/3605573.3605599

1

## **1 INTRODUCTION**

A hybrid MPI+threads model is commonly used nowadays to program parallel systems comprising nodes with multiple cores or accelerators such as GPUs. A common scenario in such a model is that multiple threads perform operations at different locations on the same buffer. In this situation, a bulk synchronization of the threads followed by a single communication is usually the chosen approach to avoid heavy congestion on the MPI resources [12], as illustrated in Figure 1. However, the load imbalance between threads or the computational load can lead to some threads idling before the communication. It delays the start of the send operation and misses the opportunity to overlap communication with computation. An alternative approach is the pipelined communication model, as illustrated in Figure 2. Instead of sending the entire buffer using one thread, each thread now sends its own section of the buffer: each thread performs computations independently and communicates the results immediately. Therefore, the first thread to end its computation gets a head start in the communication (also called the *early-bird* effect). While it enables the send operation to be started as soon as one thread completes the computation, it also brings another challenge of coordinating the communication from multiple threads. This multithreaded MPI communication pattern usually scales poorly because of the contention for shared resources such as message queues and communication contexts. Several approaches have been proposed to tackle this issue and provide performance for the pipelined communication pattern. Some of the most well-known ones are scalable endpoints [13], finepoint communication [4], thread-based MPI implementation [8], and more recently MPIX\_Stream [15]. Inspired by all these works, the MPI 4.0 standard [9] introduced point-to-point partitioned communication [2] to provide better support for the pipelined communication model and improve its adoption among users. This new feature divides the communication buffer into non-overlapping partitions where threads can operate individually. When one partition is ready, the thread marks it as "ready to be sent". The send operation is completed once the main thread completes the communication, after all the partitions have been marked as ready.

This design allows the early threads to start the communication of partitions when they become ready. It also offers two advantages compared with other MPI functionalities:

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of the United States government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

ICPP 2023, August 7-10, 2023, Salt Lake City, UT, USA

- Easy-to-use multithreaded MPI communication. The semantics of partitioned communication provides a simple interface to the user, hiding the complexity of multithreaded performance. It is now easy to benefit from the *early-bird* effect and achieve performance gain.
- (2) Flexibility in message transmission. The new API gives the implementation the opportunity to perform optimizations in order to reduce the latency, otherwise tedious to implement for users. A commonly considered approach is to aggregate partitions together into a single message in order to avoid the overhead for small partitions.

Partitioned communication has been supported in MPICH since the release of MPI 4.0. However, the initial implementation is focused primarily on correctness instead of performance. In this work we present improvements to the implementation, enabling the user to achieve the expected performance gains. Specifically, we are now able to (1) aggregate small partitions together; (2) if required by the user, reduce thread congestion when performing communication; and (3) reduce the time-to-solution using the early-bird effect. In Section 2 we present some background and assess the expected gain of using the pipelined communication pattern for very small and very large messages. In Section 3 we present the improvements to the existing implementation in MPICH. This work puts a particular focus on the user experience and measurable gain. To validate and assess the obtained performance, in Section 4 we compare the improved implementation in MPICH with other MPI 3.1 approaches, relying on both the point-to-point and one-sided semantics. In Section 5 we present our conclusions and discuss future directions.

#### **Related work and novelty**

Prior contributions have focused on assessing the partitioned communication benefits. In [2], the authors present initial performance metrics, with an emphasis on the perceived bandwidth metric for large messages. More recently, in [5] the authors use four metrics to measure the performance: overhead, perceived bandwidth, application availability, and *early-bird* communication. Still with a focus on large messages, they describe the behavior with different noise models and detail the usage for sweep- and halo-based algorithms. Partitioned collective communication has also been proposed in [6] as an extension to the MPI 4.0 semantics. Despite those efforts, a comprehensive evaluation of the newly proposed semantics against existing ones is still missing, especially for small message sizes. With this work we aim to bridge that gap and therefore guide MPI users in making informed and evidence-driven choices for their own applications.

## 2 PIPELINED COMMUNICATION: PERFORMANCE MODEL AND IMPLEMENTATION APPROACHES

A detailed view of pipelined communication is presented in Figure 3, where we highlight the different operations performed by each thread. Our presentation is intentionally general: we will later detail how to implement the pattern using different strategies. To initiate the pattern, the master thread performs a start operation, which implies a thread barrier afterward. Then each thread performs computations and marks the partition as ready. The master thread can finalize the communication using wait, which usually entails a thread barrier beforehand depending on the MPI API used.

#### 2.1 Performance measurements

Measuring the performance of the pipelined communication pattern can be done in different ways [1, 5]. In our case we focus on the user experience, and therefore the time-to-solution is the most relevant metric. As illustrated in Figure 3 in red, the latter runs from the start operation up to the completion of the communication on the receiver side. Since we benchmark the MPI-related operations and not the computation, we remove the time of each thread spent in the computation. By doing so we have a measure relevant to the user, namely, the overhead coming from the communications only. We note that this metric is close to the perceived bandwidth proposed by [2]. However, we go a step further and include the start operation and the following thread barrier into our metric.

#### 2.2 Performance prediction

The performance of the pipelined communication pattern compared with the bulk thread-synchronization can be expressed as

$$\eta = \frac{T_b}{T_p} \quad , \tag{1}$$

where  $T_b$  is the communication time with the bulk thread synchronization and  $T_p$  is the communication time with the pipelined communication pattern. For large message sizes, the value of  $\eta$  is driven by the delay coming from the computations and the load imbalance between threads. For small message sizes, however, the latency of the communication will prevail over the delay time and dictate the gain. In this section we further describe these two factors and the performance gain that a user might expect in both situations.

2.2.1 Large messages and delay time. For very large messages, the communication time will be given by  $S_{\text{part}}/\beta$ , where  $\beta$  is the bandwidth of the network and  $S_{\text{part}}$  is the size of one partition. The communication time associated with bulk thread synchronization is given by

$$T_b \approx N_{\rm part} \; \frac{S_{\rm part}}{\beta} \; ,$$
 (2)

where a total of  $N_{\text{part}}$  very large partitions will be used. Introducing  $\theta$  as the number of partitions per thread and N the number of threads, we obtain that  $N_{\text{part}} = N\theta$ . When using pipelined communication, the communication time,  $T_p$ , is given by

$$T_p \approx \max\left\{ (N_{\text{part}} - 1) \frac{S_{\text{part}}}{\beta} - D, 0 \right\} + \frac{S_{\text{part}}}{\beta} , \quad (3)$$

where *D* is the delay time between the first and the last partition to be ready and the max ensures that we overlap at most the communication time of the  $N_{\text{part}} - 1$  first partitions with the delay. The latter is assumed to depend linearly on the partition size, leading to the definition of the delay rate,  $\gamma$ , such that  $D = \gamma S_{\text{part}}$ . The



Figure 1: Bulk thread synchronization followed by the send operation. The time idle due to imbalance and computation delays is wasted.

expression of  $\gamma$  is itself a function of  $\theta$  and other parameters; see Appendix A.

Combining the equations, we obtain the theoretical gain associated with the pipelined communication as

$$\eta = \frac{T_b}{T_p} = \frac{N\theta}{\max\left\{N\theta - \gamma_\theta \ \beta \ , \ 1\right\}} \quad . \tag{4}$$

In practice, the bandwidth  $\beta$ , the algorithm, and the number of threads *N* are usually fixed. The gain is then a function of the number of partitions per thread,  $\theta$ , a user-controlled parameter. For example, with  $\theta = 1$ ,  $\beta = 25 GB/s$ , and N = 8 threads, typical values are  $\gamma \approx [1 ; 10] \ \mu s/MB$ , which lead to  $\eta = 1.003$  and  $\eta = 1.032$ , respectively. However, increasing the number of partitions per thread would enable the communication to be started earlier and therefore increase the delay rate. With  $\theta = 8$ , the value of  $\gamma$  goes up to  $\approx 1000 \ \mu s/MB$  and the gain to  $\eta = 1.641$ , leading to a more significant benefit.

We conclude that using multiple partitions per thread is therefore crucial for performance. However, it is hard to achieve in practice at large message sizes because  $\theta$  is inversely proportional to the size of the message, and our assumption of nonsignificant latency quickly becomes invalid.

2.2.2 Small messages and latency. In some situations the latency overwhelms the cost of the communication, due either to a small buffer or to a very large number of partitions per thread. In practice, this is usually the case for messages  $\ll 16kB$ . Furthermore, the delay generated by the computations is irrelevant for those small messages. Assuming a delay rate of  $\gamma = 100\mu s/MB$  (see the preceding section) and a latency of 1  $\mu s$ , a buffer of 1kB would generate enough delay to offset 10% of the latency of a single message.

Therefore, assuming that latency is the only relevant metric and that the delay is negligible, we obtain

$$\eta = \frac{1}{N\theta} \quad . \tag{5}$$

In this situation, issuing multiple messages in the pipeline communication scheme will increase the overhead and decrease the performance. To avoid this issue, the user must either aggregate



Figure 2: Pipelined send operations, initiated from each thread. The imbalance and computation delays provide a gain through the *early-bird* effect.

messages or decrease the number of partitions (and therefore increase the partition size). Furthermore, this prediction does not take into account the thread contention that will also impact the performance for small messages.



Figure 3: Benchmark for the pipelined communication pattern, illustrated with two threads. The application to each of the proposed implementation is summarized in Tables 1 and 2. Red boxes represent operations added for benchmarking purposes.

#### 2.3 Possible user approaches

To implement the pipelined communication pattern, the user might consider different approaches, divided into three categories: partitioned communications (MPI 4.0), point-to-point (MPI 3.1), and RMA-based (MPI 3.1) APIs. In Figure 3 we present the general template for the different implementations. For each of the steps, the list of MPI API calls actually used is summarized in Table 1 for the sender side and in Table 2 for the receiver side. The actual implementation for each of them is available at [3].

2.3.1 Partitioned communication. The partitioned communication API provides the user with a simple way of exploiting the pipeline communication pattern. The communication is initialized during the call to MPI\_Psend\_init. Then, the main thread calls MPI\_Start on the partitioned requests, which is followed by a thread barrier. Once the computation on a partition is completed, the thread can call MPI\_Pready to signal to MPI that the partition is ready to be sent. After a barrier, the master thread calls the MPI\_Wait function to complete the communication. The receiver side is similar. MPI\_Precv\_init is used to initialize the communication, along with MPI\_Start to start an iteration. A thread can then query the status of a partition using MPI\_Parrived, and the master thread completes the communication using MPI\_Wait.

2.3.2 Point-to-point MPI 3.1. A first approach is be to use a single message to communicate once the threads have completed their work on the different partitions. This strategy implements a bulk thread synchronization instead of the pipelined communication and is referred to as Pt2Pt single. After a thread barrier, the master thread issues the persistent communication with MPI\_Start. The receiver uses a single persistent request to receive the message. Another approach is to send one message from every thread as soon as the computation is over. This approach is denoted here as Pt2Pt many. To avoid competition on the same resource, we first duplicate the communicator per thread. Different communicators will be mapped to different communication contexts, hence removing the contention between the threads [14]. Then each thread can send and receive its own partition independently of the status of other threads. This approach, more complicated for the user than the partitioned communication, is the traditional way of taking advantage of a pipelined communication pattern.

2.3.3 One-sided communication. Our third category of implementation uses the MPI one-sided (RMA) semantics. Similar to the point-to-point variations, it can be implemented either on a single window shared by all threads or by using one window per thread, each over the entire buffer. The approaches also can be distinguished by their synchronization API: active or passive. By design, a send-receive operation is an active RMA communication pattern, which means that the target of the RMA call is involved in the communication. This pattern can be naturally implemented with the active synchronization API. However, an enhanced use of the passive synchronization API can also be used to implement an active communication pattern, at the cost of added synchronization.

In the active synchronization API, the origin opens and closes the *access epochs* through MPI\_Start and MPI\_Complete, respectively.

The target controls when its memory is exposed (*exposure epochs*) using MPI\_Post and MPI\_Wait. The main advantage of the active API is to offer explicit control to the user on the target readiness to handle the data. The active synchronization approach on a single window and on many windows is denoted as RMA single - active and RMA many - active, respectively, in the rest of this work.

In contrast, the passive synchronization API is based only on managing the *access epochs*. One must still control the *exposure epochs*, which can be done by using 0*B* send and receive messages. We note that ensuring progress with passive synchronization can be a challenge, especially when no global progress is done in MPICH. Different approaches exist to address this issue; see [11] for a thoughtful discussion of them. For our specific case, we have chosen to use MPI\_MODE\_NOCHECK when locking the window, to avoid requiring the receiver to be involved in the synchronization at that stage. In the rest of this work, we denote the passive synchronization on a single window and on multiple windows, respectively, as RMA single – passive and RMA many – passive.

## 3 PARTITIONED COMMUNICATION IMPLEMENTATION IN MPICH

In this section we briefly describe the existing implementation of the partitioned communications in MPICH and the improvements made as part of this work. This section focuses on the underlying mechanisms, which are useful for understanding the performance obtained.

## 3.1 Existing implementation

The current implementation of partitioned communication in MPICH uses a single-message approach, done through the active messaging (AM) code path. When the user calls MPI\_Psend\_init, an atomic counter is associated with the partition request. A "ready-to-send" (RTS) message is sent to the receiver with some of the basic information about the size of the data and the number of partitions. During MPI\_Start, the counter is set to the number of partitions given by the user, plus one. The "plus one" takes into account that for each iteration the sender has to wait for a "clear-to-send" (CTS) message from the receiver. Because of the AM nature, this mandatory CTS avoids early sends from the sender to a receiver still in the previous iteration. Upon receiving the CTS, the sender will decrement the counter by one. Then, when a partition is ready and the user calls MPI\_Pready, the counter is decremented by one. Once all the partitions are ready and the CTS has been received, the value of the counter is zero, and the whole buffer is sent to the receiver.

The use of AM, together with the CTS needed at each iteration, delivers a semantically correct implementation, yet not the expected performance for the user. Specifically we would see no benefit of the *early-bird* effect coming with the pipelined communication.

#### 3.2 Improvements

We have improved the implementation in MPICH to use multiple internal tag-matched messages, instead of a single AM communication. Another option would have been to rely on an RMA-supported

#### Table 1: MPI operations for the sender side.

	init	start	ready	wait
Pt2Pt part	MPI_Psend_init	MPI_Start	MPI_Pready	MPI_Wait
Pt2Pt single	MPI_Send_init			MPI_Start MPI_Wait
Pt2Pt many	MPI_Comm_dup MPI_Send_init		MPI_Start MPI_Wait	
RMA single - passive	MPI_Comm_dup MPI_Win_create MPI_Win_lock	MPI_Recv	MPI_Put	MPI_Win_flush MPI_Send
RMA many - passive	MPI_Win_create MPI_Win_lock	MPI_Recv	MPI_Put MPI_Win_flush	MPI_Send
RMA single - active	MPI_Comm_dup MPI_Win_create	MPI_Start	MPI_Put	MPI_Complete
RMA many - active	MPI_Win_create		MPI_Start MPI_Put MPI_Complete	

#### Table 2: MPI operations for the receiver side.

	init	start	ready	wait
Pt2Pt part	MPI_Precv_init	MPI_Start	MPI_Parrived	MPI_Wait
Pt2Pt single	MPI_Recv_init	MPI_Start		MPI_Wait
Pt2Pt many	MPI_Comm_dup MPI_Recv_init		MPI_Start MPI_Wait	
RMA single – passive RMA many – passive	MPI_Win_create	MPI_Send		MPI_Recv
RMA single – active RMA many – active	MPI_Win_create	MPI_Post		MPI_Wait

implementation, as suggested in [2]. We decided not to follow this strategy for two reasons. First, the difference between the two approaches matters only for small messages. Second, an RMA-based approach requires exposure control (see Section 2.3.3), which increases the overhead. To alleviate this overhead, we rely on the repetitive use of a put operation, faster than a tag-matching send. For small messages, however, optimal performance is obtained with a few messages (see Section 2.2.2). In this configuration an RMA-based implementation is then slower, as illustrated in Section 4.1.

3.2.1 Initialization. During the initialization, the sender and the receiver will agree on using the tag-matching code path and on a fixed number of messages to be sent. The tag matching can be used only if there is enough tag space to isolate the traffic of partitioned communication from other communications coming from the user. The sender keeps a count of the number of partitioned requests created for each of the receiver ranks. If that number exceeds the tag space reserved for the partitioned communications, the AM code path is used instead.

The sender and the receiver have to agree on the number of messages to be actually sent. The most general protocol is to let the receiver decide this number once the RTS has been received. Then, that information is sent to the sender with the CTS. However, this general approach incurs a performance overhead for two reasons. First, the sender has to wait for the CTS during the first iteration. Second, a CTS does not naturally fit within a tag-matching communication protocol. Another approach would be to let the sender decide on the number of messages to be sent. However, this strategy adds complexity when the sender and/or the receiver uses noncontiguous datatypes. In this case the receiver might receive a partial datatype, and dealing with this scenario efficiently is challenging.

In this work we have chosen to implement the first, yet suboptimal, approach. The receiver decides the number of messages to be sent using gcd  $\left(N_{\text{part}}^{\text{send}}, N_{\text{part}}^{\text{recv}}\right)$ , which guarantees that a partition will contribute only to a single message. In our implementation, the receiver is also in charge of message aggregation, based on the user-defined threshold MPIR\_CVAR\_PART\_AGGR\_SIZE. This value is

used as an upper bound for aggregation: if the size of multiple messages fit within the prescribed threshold, then they are aggregated together. The number of messages obtained from this logic is then sent to the sender as part of the CTS.

We note that the use of MPI\_Parrived is in contradiction with message aggregation. Indeed, the former is used to reduce the overhead by exploiting a coarser-grained communication strategy, whereas the latter suggests that the user could exploit a finegrained communication pattern. In our implementation we have chosen to optimize MPICH toward achieving low latency, and therefore we have not spent much effort in optimizing the usage of MPI\_Parrived.

3.2.2 Sending and receiving partitions. On the sender side, at each iteration, each message to be sent is associated with an atomic counter whose value is set to the number of partitions contributing to the message. When a partition is marked as ready by the user, the associated counter is decremented. If the value reaches zero, the message is then sent using tag-matching or an AM send/receive MPICH internal API. On the receiver side, each message is associated with a receive request. The user can then query the reception of a given partition, hence reading the status of the request.

We also allow the user to use different communication resources (known as VCI in MPICH) to send different partitions. This is done by encoding the source and destination VCI id into the tag, using an experimental feature in MPICH. With this, we reduce the thread congestion that occurs when sending from multiple threads using the same resource. However, despite the ubiquitous multithreaded context when using partitioned communication semantics, the user has no standard way of conveying the thread granularity to the MPI implementation. Therefore, our implementation assumes a round-robin attribution of the threads to the partitions. This assumption is inflexible and likely to break when used in practice with  $\theta > 1$ . We note that info hints provided during communication initialization or the usage of MPIX\_Stream [15] with partitioned communication could be used to express the thread granularity to MPI. Such improvements are left for future work.

In summary, the improvements offer better performance than the existing implementation, as detailed in Section 4.1. Moreover, the user now has the opportunity to take advantage of three possible gains:

- (1) Thread Congestion (experimental): When multiple threads send different partitions simultaneously, they will compete for the same resources. The congestion is overwhelming especially for small partitions. To alleviate this issue, we use an experimental MPICH capability to use different resources for each partition. Further details can be found in Section 4.2.1.
- (2) Message Aggregation: When dealing with small partitions sizes, different partitions can be aggregated together under a single message to reduce the overhead. Results on this performance gain are detailed in Section 4.2.2.
- (3) *Early-Bird Effect*: When sending large partitions, the user now benefits from the gain offered by the pipelined communication model as shown in Section 4.3.

#### 4 PERFORMANCE RESULTS

In this section we compare the performance for each of the possible pipelined communication implementations (see Section 2.3), using the benchmark template described in Figure 3. We first assess the benefits of the improved implementation and compare it with other MPI-3.1 approaches. To streamline our analysis and avoid implementation artifacts, we assume that the number of partitions is the same on both the sender and the receiver side. Next, we investigate the performance for small messages, together with the impact of message aggregation and thread congestion. We then compare the different approaches to achieve pipelined communication with the expected gain when using large messages, as detailed in Section 2.2.1.

All the results in this work have been obtained using MPICH and ucx-1.13.1 between two nodes on MeluXina.<sup>1</sup> The openmp threads are closely bound to the cores,<sup>2</sup> and the MPI processes are bound to as many cores as there are threads.<sup>3</sup> The benchmark [3] has been run for 150 iterations and 1 warm-up iteration to get rid of the overhead, explained in Section 3.2. For each of the data results, we present the time as the average on the iterations (excluding the warm-up), and we obtain a 90% confidence interval assuming a Student's t-distribution. To avoid network noise, we rerun the measure if the half-width of the confidence interval is larger than 5% of the average time, with a maximum retry at 50. Confidence intervals are displayed as a shaded area around the results on figures displaying time.

#### 4.1 Improvement over existing implementation

To demonstrate the gain of not using the AM path, we measure the time in the case of N = 1 threads,  $\theta = 1$  partition per thread with no delay ( $\gamma = 0$ ). Although not representative of the usual operation space of pipelined communications, the configuration is well suited to highlight the performance gain made possible by our improvements. In Figure 4 we show the time needed by each of the approaches to complete the communication. For reference, we also give the time corresponding to the theoretical bandwidth of the system (25*GB/s*). The difference between the existing AM-based implementation (Pt2Pt part - old) and our improved version (Pt2Pt part) is noticeable for all message sizes. The latency associated with the copy needed in the AM code path implies a significant overhead and degrades the performance. With the new implementation we match the performance of the Pt2Pt single approach, as expected.

For point-to-point-based approaches, we note that the time jumps when switching protocols over the different messages sizes. In particular, we note the change from the short protocol to the bcopy one between 1,024 and 2,048 *B* and to the *rendez-vous* (zcopy) protocol from 8,192 to 16,384 *B* [10]. The difference between the two families of approaches (point-to-point and RMA) is also clearly

<sup>2</sup>with OMP\_PROC\_BIND=CLOSE and OMP\_PLACES=cores

<sup>3</sup>using -bind-to cores:\${OMP\_NUM\_THREADS}

 $<sup>^1</sup>$ System in Luxembourg (#379 top500 06/2023), with 73,344 cores (AMD EPYC 7H12 64C) on its CPU partition connected through a Mellanox 200Gb/s and 1.22 $\mu s$  latency HDR200-IB network

observed at small message sizes. The RMA-based approaches require two additional synchronizations to be performed, resulting in a larger overhead. We also note that the gap vanishes when considering message sizes above the *rendezvous* threshold. The reason is that the bandwidth is dominant for large message sizes and all the approaches use the same communication protocol. The zcopy protocol used in point-to-point is actually based on the RMA network capability.



Figure 4: Time across message sizes with 1 thread and 1 partition: comparison of the existing and improved partitioned communication implementation with other MPI-3.1 approaches.

#### 4.2 Small messages

4.2.1 Thread congestion. In practice, partitioned communication is used in a multithreaded environment, which will lead to thread congestion. To highlight the issue, we present in Figure 5 the time needed to communicate when using 32 threads and one partition per thread ( $\theta = 1$ ).

As expected for small messages, the Pt2Pt single approach performs the best. The single message does not suffer from any of the downsides of the multithreading since the communication happens on one thread only. However, we still note a higher latency compared with Figure 4 due to the needed synchronization barrier. The Pt2Pt part and Pt2Pt many communication strategies both suffer from thread contention, with little difference between the achieved overheads. With the RMA-based passive synchronization approaches, the results are more sparse. We observe that the RMA approaches using many windows (one per partition) suffer from an additional overhead compared with the single RMA window. Regarding the MPI\_Put operations, there is no significant difference because in both cases the threads will compete for the same resource. However, the RMA many - passive approach adds an overhead in the progress engine compared with the RMA single - passive since it has to operate on multiple windows simultaneously. This causes the upward shift observed in Figure 5.



Figure 5: Thread congestion: communication time across message sizes for 32 partitions with 32 threads.

In order to reduce the overhead, MPICH can be configured to use multiple *virtual communication interfaces* (VCIs) [14]. This is achieved by using MPIR\_CVAR\_NUM\_VCIS to control the number of VCIs used by the implementation. Different communicators/windows will then be mapped onto different VCIs, which allow multiple threads to access different resources. As detailed in Section 3.2, in the improved partitioned communication implementation, we map each partition on a different VCIs using a round-robin strategy.<sup>4</sup> The results obtained when using one VCI per thread are presented in Figure 6. In this setting, the Pt2Pt many strat-



Figure 6: Thread congestion: communication time across message sizes for 32 partitions with 32 threads using 32 VCIs.

egy reaches the same performance as the Pt2Pt single method.

<sup>&</sup>lt;sup>4</sup> if the user has used -enable-vci-method=tag during the configuration

As we expect, different communicators are assigned to different VCIs, which then leads to no contention for the Pt2Pt many approach. The Pt2Pt single method is further slightly penalized by the needed thread barrier before starting the send operation. The Pt2Pt part code path exploits the different VCIs as well; however, it still suffers from an overhead compared with Pt2Pt single. Compared with the non-VCI usage, we have decreased the cost of thread contention by a factor of  $\approx$  10. Regarding the RMA-based implementations, the RMA many – passive is now faster than the RMA single – passive. The former approach relies on different VCIs (one per window) and therefore avoids the cost of contention.

As pointed out in Section 3.2.2, the partitioned communication implementation may lack necessary information to avoid VCI contention in a multithreaded environment. If this is the usage model of one's application, we would favor the use of the Pt2Pt many approach to get better performance. In the rest of this work, we will consider a single VCI to illustrate the expected application context of the partitioned communication API.

4.2.2 Partition aggregation. To reduce the latency, one can also gather multiple partitions as a single message to avoid multiplication of the individual overheads. As explained in Section 3.2.1, the user can use MPICH's environment variable MPIR\_CVAR\_PART\_AGGR\_SIZE (in bytes) to request an upper bound on the aggregation size. We note that this technique is compatible with other approaches such as Pt2Pt many, but it would require significant code changes from the user. To avoid interference with other delays and to focus our analysis on the message aggregation, we consider that all the partitions are ready immediately and processed in order by each thread. The results of this approach are shown in Figure 7, where the messages are aggregated from 512 up to 16,384 *B*.

We observe that the Pt2Pt part reduces the overhead for small messages significantly compared with the Pt2Pt many approach, which has the same performance as the no-aggregation based Pt2Pt part. For a given aggregation size, the number of messages actually sent increases with the size of the global buffer. Therefore, message aggregation is beneficial for global message buffer size only below  $N_{\text{part}}$  times the aggregation size. As illustrated in Figure 7, larger aggregation sizes. Regardless of the size, however, we do not match the latency of the Pt2Pt single approach. The reason is the added overhead of partitioned communications such as the atomic update on the message counter performed by every partition when ready. The latter becomes more significant for an increased number of partitions. For an infinite aggregation size and neglecting the overhead, Pt2Pt single is then the upper bound of performance.

4.2.3 Take-away. For a user focused on a usage based on many small partitions, the partitioned communication offers an easy interface with little to no overhead compared with most advanced point-to-point APIs. Achieving the same performance using the standard MPI 3.1 API would complicate significantly the user's code, especially when message aggregation is desired. For a user focused on performance, however, moving to a more explicit yet more complex API will take full advantage of the state-of-the-art



Figure 7: Message aggregation: time for  $\theta$  = 32 partitions per thread and 4 threads.

features in MPICH. For cases with many threads, we recommend the use of the Pt2Pt many API with multiple VCIs. For cases with many partitions per threads and a few threads, we recommend instead the use of the Pt2Pt single approach to reduce the latency.

# 4.3 Large messages: benefit of the *early-bird* effect

When using large messages, the delay generated by computations and load imbalance is significant, and it will drive the performance, as detailed in Section 2.2. One could measure the gain obtained for different values of  $\theta$ , different algorithm parameters, and so forth. As demonstrated in Appendix A, however, the value of  $\gamma$ accurately models the delay obtained in those different situations. Therefore, we rely instead on the value of  $\gamma$  to characterize the delay obtained in various practical situations, which allows us to measure only cases with  $\theta = 1$ . Specifically, the last partition is delayed compared with the other  $N_{\text{part}}$  – 1 partitions, where the delay time is given by  $\gamma S_{part}$ , where  $S_{part}$  is the partition size. Then, we measure the obtained bandwidth and compare it with the theoretical gain predicted in Section 2.2. In Figure 8 we present the results with  $\gamma = 100 \ \mu s/MB$ , which represents a practical delay; see Appendix A.2. With a total of 4 partitions and 4 threads, we measure a gain of  $\approx$  2.54, close to the theoretical value of 2.67. The difference comes from the latency involved in the actual communications and the thread congestion, both left out of the model; see Section 2.2.1.

We note that, as expected, the gain obtained from the *early-bird* effect is independent of the approach used by the user. Since the messages are dominated by bandwidth and have (almost) negligible latency, every possible variation of the MPI API will provide the same gain. As highlighted earlier, the partitioned communication API provides a simple interface to the user in order to achieve that gain. In real-life cases, however, the actual gain from an application



Figure 8: Gain obtained with the *early-bird* effect ( $\gamma = 100\mu s/MB$ , which stands for a value of  $\theta > 1$ ) with 4 threads and 4 partitions

perspective is tightly coupled to the size of the partitions and the delay achieved by the application.

In summary, the results of Figure 8 represent perfectly the usage of partitioned communication for pipelined communication. As expected by our performance modeling in Section 2.2, we observe that with small messages it adds an overhead, due to the thread congestion and multiple latency costs. Therefore, for a fixed number of threads, using  $\theta = 1$  will lead to the best performances. For larger messages, however, the gain is significant as one hides communications behind computations. A larger number of partitions per thread ( $\theta \gg 1$ ) will lead to a larger delay rate ( $\gamma$ ) and therefore to a larger measured gain. With this example, we observe the trade-off to be around 100 *kB*, a value driven by thread congestion.

#### **5** CONCLUSIONS

In this work we investigate the pipelined communication pattern and the expected gain from the early-bird effect. First, we introduce a theoretical model to quantify the expected gain and to identify in which cases it will be beneficial. Then, we present the improvements made to the MPICH implementation in order to deliver the expected performance to the user. Specifically, we provide three features: (1) thread congestion alleviation, (2) message aggregation, and (3) early-bird effect gain obtained by starting the communication as soon as the data is ready. Further, we explore various other approaches that rely on MPI 3.1 features and could also be used to implement pipelined communication. We study the use of point-to-point-based approaches, such as using a single message or one message per partition. We also consider various one-sided strategies relying on a single window, multiple windows, and active and passive synchronization. We use a pipelined communication benchmark to compare them with the partitioned communication semantics, including the existing and improved implementations

in MPICH. Then, we investigate three specific cases across the spectrum of typical usage of the pipelined pattern. First, we consider the case of small messages where multiple threads contend for the same MPI resources. Relying on existing thread contention alleviation strategies in MPICH, we are able to reduce dramatically the associated overhead: Compared with a single-message approach, we reduce the penalty from a factor of  $\approx 30$  to  $\approx 4$ . Second, we demonstrate message aggregation and how it reduces the overhead associated with multiple messages to a single-message latency, at the cost of a few atomic updates. With that strategy we are able to reduce the penalty factor from  $\approx 10$  to  $\approx 3$  compared with a single-message approach. Third, we demonstrate how the user can benefit from a significant bandwidth improvement when using pipelined communication with large messages, even with thread contention. In the context of the presented results, we measure a benefit for messages larger than  $\approx 100 \text{ kB}$ . We also demonstrate that this benefit is agnostic to the type of method used (point-topoint or one-sided). We conclude that the best configuration to use partitioned communication depends on the partition size. To avoid significant overhead with small messages, the user should use message aggregation, or other existing MPI-3.1 semantics, in order to send as few messages as possible. With large partition size, however, a higher number of partitions will lead to greater performance benefit, as latency and thread contention become negligible. Partitioned communication then delivers the expected benefits of pipelined communication, similar to other existing MPI methods.

From our work, we estimate that the strength and the weakness of the partitioned communication semantics are in the ease of use of its interface. The latter leads to suboptimal performance when used with small messages, because of thread contention. The reason is not new to MPI: to provide a well-performing implementation for both many partitions per thread and many threads, the implementation needs to be able to exploit a user-provided thread context identifier. Doing so would guarantee no conflict when accessing the resources and the optimal performance in every scenario. A user worried about performance should therefore use other existing strategies such as MPI\_Comm\_dup or the more lightweight MPIX\_Stream, to isolate communications issued from different threads.

Extending partitioned communication to GPU accelerators is an active topic, which is at the center of our future work. We also plan to further improve our implementation and to remove the need of synchronization between the sender and the receiver during the first iteration.

#### ACKNOWLEDGEMENTS

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration. Computational resources have been provided by EuroHPC for the access to MeluXina (*EHPC-BEN-2022B01*).

#### REFERENCES

 M. G. F. Dosanjh, T. Groves, R. E. Grant, R. Brightwell, and P. G. Bridges. 2016. RMA-MT: A Benchmark Suite for Assessing MPI Multi-threaded RMA Performance, In 2016 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid). 2016 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid), 550–559. https://doi.org/10.1109/CCGrid.2016.84

- [2] Matthew G. F. Dosanjh, Andrew Worley, Derek Schafer, Prema Soundararajan, Sheikh Ghafoor, Anthony Skjellum, Purushotham V. Bangalore, and Ryan E. Grant. 2021. Implementation and evaluation of MPI 4.0 partitioned communication libraries. *Parallel Comput.* 108 (2021), 102827. https://doi.org/10.1016/j. parco.2021.102827
- [3] Thomas Gillis. 2023. bench-pcomm. https://github.com/pmodels/bench-pcomm.
- [4] Ryan E. Grant, Matthew G. F. Dosanjh, Michael J. Levenhagen, Ron Brightwell, and Anthony Skjellum. 2019. Finepoints: Partitioned Multithreaded MPI Communication. In *High Performance Computing* (2019//). Springer International Publishing, 330–350. https://doi.org/10.1007/978-3-030-20656-7\_17
- [5] Yiltan Hassan Temucin, Ryan E. Grant, and Ahmad Afsahi. 2023. Micro-Benchmarking MPI Partitioned Point-to-Point Communication. In Proceedings of the 51st International Conference on Parallel Processing (Bordeaux, France) (ICPP '22). Association for Computing Machinery, New York, NY, USA, Article 64, 12 pages. https://doi.org/10.1145/3545008.3545088
- [6] D. J. Holmes, A. Skjellum, J. Jaeger, R. E. Grant, P. V. Bangalore, M. G. F. Dosanjh, A. Bienz, and D. Schafer. 2021. Partitioned Collective Communication, In 2021 Workshop on Exascale MPI (ExaMPI). 2021 Workshop on Exascale MPI (ExaMPI), 9–17. https://doi.org/10.1109/ExaMPI54564.2021.00007
- [7] Huda Ibeid, Luke Olson, and William Gropp. 2020. FFT, FMM, and multigrid on the road to exascale: Performance challenges and opportunities. *J. Parallel and Distrib. Comput.* 136 (2020), 63–74. https://doi.org/10.1016/j.jpdc.2019.09.014
- [8] H. Kamal and A. Wagner. 2010. FG-MPI: Fine-grain MPI for multicore and clusters, In 2010 IEEE International Symposium on Parallel & Distributed Processing, Workshops and Phd Forum (IPDPSW). 2010 IEEE International Symposium on Parallel & Distributed Processing, Workshops and Phd Forum (IPDPSW), 1–8. https: //doi.org/10.1109/IPDPSW.2010.5470773
- [9] Message Passing Interface Forum. 2021. MPI: A Message-Passing Interface Standard Version 4.0. https://www.mpi-forum.org/docs/mpi-4.0/mpi40-report.pdf
  [10] P. Shamis, M. G. Venkata, M. G. Lopez, M. B. Baker, O. Hernandez, Y. Itigin, M.
- [10] P. Shamis, M. G. Venkata, M. G. Lopez, M. B. Baker, O. Hernandez, Y. Itigin, M. Dubman, G. Shainer, R. L. Graham, L. Liss, Y. Shahar, S. Potluri, D. Rossetti, D. Becker, D. Poole, C. Lamb, S. Kumar, C. Stunkel, G. Bosilca, and A. Bouteiller. 2015. UCX: An Open Source Framework for HPC Network APIs and Beyond, In 2015 IEEE 23rd Annual Symposium on High-Performance Interconnects. 2015 IEEE 23rd Annual Symposium on High-Performance Interconnects, 40–43. https://doi.org/10.1109/HOTI.2015.13
- [11] M. Si, A. J. Peña, J. Hammond, P. Balaji, M. Takagi, and Y. Ishikawa. 2015. Casper: An Asynchronous Progress Model for MPI RMA on Many-Core Architectures, In 2015 IEEE International Parallel and Distributed Processing Symposium. 2015 IEEE International Parallel and Distributed Processing Symposium, 665–676. https: //doi.org/10.1109/IPDPS.2015.35
- [12] Rohit Zambre and Aparna Chandramowlishwaran. 2022. Lessons Learned on MPI+threads Communication. In Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (Dallas, Texas) (SC '22). Article 77, 16 pages.
- [13] R. Zambre, A. Chandramowlishwaran, and P. Balaji. 2018. Scalable Communication Endpoints for MPI+Threads Applications, In 2018 IEEE 24th International Conference on Parallel and Distributed Systems (ICPADS). 2018 IEEE 24th International Conference on Parallel and Distributed Systems (ICPADS), 803–812. https://doi.org/10.1109/PADSW.2018.8645059
- [14] Rohit Zambre, Aparna Chandramowliswharan, and Pavan Balaji. 2020. How I Learned to Stop Worrying about User-Visible Endpoints and Love MPI. In Proceedings of the 34th ACM International Conference on Supercomputing (Barcelona, Spain) (ICS '20). Association for Computing Machinery, New York, NY, USA, Article 35, 13 pages. https://doi.org/10.1145/3392717.3392773
- [15] Hui Zhou, Ken Raffenetti, Yanfei Guo, and Rajeev Thakur. 2022. MPIX Stream: An Explicit Solution to Hybrid MPI+X Programming. In Proceedings of the 29th European MPI Users' Group Meeting (Chattanooga, TN, USA) (EuroMPI/USA'22). Association for Computing Machinery, New York, NY, USA, 1–10. https://doi. org/10.1145/3555819.3555820

## A DELAY RATE

#### A.1 Definition

The delay in the pipelined communication pattern comes from the computation time assumed to be proportional to the partition size ( $S_{part}$ ) and to the average computation rate  $\mu$ , and the computation noise to follow a normal distribution,  $\mathcal{N}(0; \sigma\mu)$ , whose standard deviation is proportional to  $\mu$ .

The average computation rate,  $\mu$ , depends on several factors:

- the CPU, represented by its frequency, *F*, and the number of *flops* per cycles; and
- the algorithm used, described by the arithmetic intensity (*AI*, given in *flop/B*), and the communication intensity (*CI*), that is, the number of bytes actually sent/received compared with the memory used by the algorithm.

We then obtain that

$$\mu = \frac{AI}{CI} \frac{1}{(8\ F)} \quad . \tag{6}$$

On the other hand, the noise accumulated during the computations,  $\sigma$ , depends on two other factors: the algorithmic imbalance in the computations (different branches lead to different computations),  $\delta$ ; and the noise in the system execution[2],  $\epsilon$ . Hence we get that  $\sigma = (\epsilon + \delta)/2$ . Finally, we obtain that the computation time of a given partition is given by

$$T_{cmpt} = \mu \mathcal{N}\left(1 \ ; \ \frac{\epsilon + \delta}{2}\right) \quad .$$
 (7)

The delay time in the pipelined communication pattern is the time elapsed between the first partition to be ready and the last one. Assuming the Gaussian model described earlier, the first partition will be ready after  $\mu S_{\text{part}}(1-\sigma)$ . The last partition will be ready once the  $\theta$  partitions on the thread have been processed with some delay,  $\mu S_{\text{part}}\left(\theta + \sqrt{\theta}\sigma\right)$ . The delay time is then obtained as the difference between the two:

$$D = \gamma_{\theta} S_{\text{part}} = \mu \left( \theta + \frac{\epsilon + \delta}{2} (\sqrt{\theta} + 1) - 1 \right) S_{\text{part}} \quad , \tag{8}$$

leading to the definition of the delay rate as being

$$\gamma_{\theta} = \mu \left( \theta + \frac{\epsilon + \delta}{2} (\sqrt{\theta} + 1) - 1 \right) \quad . \tag{9}$$

#### A.2 Numerical examples

*A.2.1* Fourier transform. A distributed FFT has an  $AI \approx 5$ , CI = 1, and  $\delta = 0$  (no algorithmic delay) [7]. Assuming a reasonable level of noise ( $\epsilon = 0.04$ ) and 8 threads, we get a delay rate of  $\gamma_1 = 7.1428$  for  $\theta = 1$ ,  $\gamma_2 = 187.1936$  for  $\theta = 2$  and  $\gamma_8 = 1263.67$  for  $\theta = 8$ . The associated gains would then be  $\eta = 1.0228$ ,  $\eta = 1.4134$ , and  $\eta = 1.9748$ .

*A.2.2* Finite difference stencil. For a distributed 3D finite difference stencil, considering one cubic block of data per rank of size  $64^3$  and two ghost points, the *CI* is  $(66/64)^3 - 1 \approx 0.1$ . The *AI* values of a finite difference stencils are usually around  $\approx 1/13$  (4th order) and the  $\delta$  can be high in some applications,  $\delta = 0.5$  indicating that some algorithmic branches can lead to 50% more computations. With N = 8, we obtain  $\gamma_1 = 15.3398$ ,  $\gamma_2 = 46.92385411$ , and  $\gamma_8 = 228.21310932$ . The associated gains are then given by  $\eta = 1.1060$  for  $\theta = 1$ ,  $\eta = 1.1718$  for  $\theta = 1$ , and  $\eta = 1.2169$  for  $\theta = 8$ .