

## Assessing Contention Effects on MPI\_Alltoall Communications

Luiz Angelo Steffemel, Maxime Martinasso, Denis Trystram

► **To cite this version:**

Luiz Angelo Steffemel, Maxime Martinasso, Denis Trystram. Assessing Contention Effects on MPI\_Alltoall Communications. 2nd International Conference on Grid and Pervasive Computing - GPC 2007, May 2007, France. pp.00. hal-00136204

**HAL Id: hal-00136204**

**<https://hal.archives-ouvertes.fr/hal-00136204>**

Submitted on 12 Mar 2007

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Assessing Contention Effects on MPI\_Alltoall Communications

Luiz Angelo Steffene<sup>1</sup>, Maxime Martinasso<sup>2</sup> and Denis Trystram<sup>2</sup>

<sup>1</sup>Université Nancy-2, LORIA, AlGorille Team, Nancy, France  
*Luiz-Angelo.Steffene@univ-nancy2.fr*

<sup>2</sup>LIG - Laboratoire d'Informatique de Grenoble, Grenoble, France  
*{Maxime.Martinasso,Denis.Trystram}@imag.fr*

**Abstract.** One of the most important collective communication patterns used in scientific applications is the *complete exchange*, also called *All-to-All*. Although efficient algorithms have been studied for specific networks, general solutions like those available in well-known MPI distributions (e.g. the MPI\_Alltoall operation) are strongly influenced by the congestion of network resources. In this paper we present an integrated approach to model the performance of the All-to-All collective operation, which consists in identifying a contention signature that characterizes a given network environment, using it to augment a contention-free communication model. This approach, assessed by experimental results, allows an accurate prediction of the performance of the All-to-All operation over different network architectures with a small overhead.

**Keywords.** Network Contention, MPI, Collective Communications, Performance Modeling

## 1 Introduction

One of the most important collective communication patterns for scientific applications is the *total exchange* [1] (also called *All-to-All*), in which each process holds  $n$  different data items that should be distributed among the  $n$  processes, including itself. An important example of this communication pattern is the All-to-All operation, where all messages have the same size  $m$ .

Although efficient All-to-All algorithms have been studied for specific networks structures like meshes, hypercubes, tori and circuit-switched butterflies, general solutions like those found in well-known MPI distributions rely on direct point-to-point communications among the processes. Because all communications are started simultaneously, architecture independent algorithms are strongly influenced by the saturation of network resources and subsequent loss of packets - the network contention.

In this paper we present a new approach to model the performance of the All-to-All collective operation. Our strategy consists in identifying a *contention signature* that characterizes a given network environment. Using such *contention signature*, we are able to accurately predict the performance of the All-to-All operation, with an arbitrary number of processes and message sizes. To demonstrate

our approach, we present experimental results obtained with different network architectures (Fast Ethernet, Gigabit Ethernet and Myrinet). We believe that this model can be extremely helpful on the development of application performance prediction frameworks such as PEMPIs [2], but also in the optimization of grid-aware collective communications (e.g.: LaPIe [3] and MagPIe [4]).

This paper is organized as follows: Section 2 presents a survey of performance modeling under communication contention. Section 3 presents the network models used in this paper, and in section 4 we formalize the *total exchange* problem, as well as some performance lower bounds. In Section 5 we propose a strategy to characterize the *contention signature* of a given network and for instance, to predict the performance of the All-to-All operation. Section 6 validates our model against experimental data obtained on different network architectures (Fast Ethernet, Gigabit Ethernet and Myrinet). In Section 7 we provide a study case for predicting the performance of a grid-aware All-to-All algorithm. Finally, Section 8 presents some conclusions and the future directions of our work.

## 2 Related Works

In the *All-to-All* operation, every process holds  $m \times n$  data items that should be equally distributed among the  $n$  processes, including itself. Because general implementations of the All-to-All collective communication rely on direct point-to-point communications among the processes the network can easily become saturated, and by consequence, degrade the communication performance. Indeed, Chun and Wang [5][6] demonstrated that the overall execution time of intensive exchange collective communications are strongly dominated by the network contention and congestive packet loss, two aspects that are not easy to quantify. As a result, a major challenge on modeling the communication performance of the All-to-All operation is to represent the impact of network contention.

Unfortunately, most communication models like those presented by Christara *et al.* [1] and Pjesivac-Grbovic *et al.* [7] do not take into account the potential impacts of network contention. Indeed, these works usually represent the All-to-All operation as parallel executions of the *personalized one-to-many* pattern [8], as presented by the linear point-to-point model below, where where  $\alpha$  is the start-up time (the latency between the processes),  $\frac{1}{\beta}$  is the bandwidth of the link,  $m$  represents the message size in bytes and  $n$  corresponds to the number of processes involved in the operation:

$$T = (n - 1) \times (\alpha + \beta m) \quad (1)$$

The development of contention-aware communication models is relatively recent, as shown by Grove [9]. For instance, Adve [10] presented one of the first models to take into account the effects of resource contention. This model considers that the total execution time of parallel programs is the sum of four components, namely:

$$T = t_{computation} + t_{communication} + t_{resource-contention} + t_{synchronization} \quad (2)$$

While conceptually simple, this model was non-trivial in practice because of the non-deterministic nature of resource contention, and because of the difficulty to estimate average synchronization delays.

While the non-deterministic behavior of the network contention is a major obstacle to modeling communication performance, some authors suggested a few techniques to adapt the existing models. As consequence, Bruck [11] suggested the use of a *slowdown factor* to correct the performance predictions. Similarly, Clement *et al.* [12] introduced a technique that suggested a way to account contention in shared networks such as non-switched Ethernet, consisting in a contention factor  $\gamma$  that extends the linear communication model T:

$$T = \alpha + \beta \times m \times \gamma \quad (3)$$

where  $\gamma$  is equal to the number of processes. A restriction on this model is that it assumes that all processes communicate simultaneously, which is only true for a few collective communication patterns. Anyway, in the cases where this assumption holds, they found that this simple contention model enhanced the accuracy of their predictions for essentially zero extra effort.

The use of a contention factor was supported by the work of Labarta *et al.* [13], that intent to approximate the behavior of the network contention by considering that if there are  $m$  messages ready to be transmitted, and only  $b$  available buses, then the messages are serialized in  $\lceil \frac{m}{b} \rceil$  communication waves. Also, König *et al.* [14] have shown indeed that some All-to-All algorithms that are optimal with unlimited buffers become less efficient when communications depend on restricted buffers size.

A similar approach was followed by Jeannot *et al.* [15], who designed scheduling algorithms for data redistribution through a backbone. In their work, they suppose that at most  $k$  communications can be performed at the same time without causing network contention (the value of  $k$  depending on the characteristics of the platform). Using the knowledge of the application transfer pattern, they proposed two algorithms to schedule the messages transfer, performing an application-level congestion control that in most cases outperforms the TCP contention control mechanism.

Most recently, some works aimed to design contention-aware performance models. For instance, LoGPC [16] presents an extension of the LogP model that tries to determine the impact of network contention through the analysis of  $k$ -ary  $n$ -cubes. Unfortunately, the complexity of this analysis makes too hard the application of such model in practical situations.

Another approach to include contention-specific parameters in the performance models was introduced by Chun [6]. In his work, the contention is considered as a component of the communication latency, and by consequence, the model uses different latency values depending on the message size. Although easier to use than LoGPC, this model does not take into account the number of messages passing in the network nor the link capacity, which are clearly related to the occurrence of network contention.

### 3 Network Models Definition

In this work we assume that the network is fully connected, which corresponds to most current parallel machines with distributed memory.

**Communication Model:** The links between pairs of processes are bidirectional, and each process can transmit data on at most one link and receive data on at most one link at any given time.

**Transmission Model:** We use Hockney's notation [17] to describe our transmission model. Therefore, the time to send a message of size  $w_{i,j}$  from a process  $p_i$  to another process  $p_j$ , is  $\alpha + w_{i,j}\beta$ , where  $\alpha$  is the start-up time (the communication latency between the processes) and  $\frac{1}{\beta}$  is the bandwidth of the link. As in this paper we assume that all links have the same latency and bandwidth, and because we only investigate the regular version of the MPI\_Alltoall operation where all messages have the same size  $m$ ,  $\forall i, \forall j, w_{i,j} = m$ , and therefore the time to send a message from a process  $p_i$  to a process  $p_j$  is  $\alpha + m\beta$ .

**Synchronization Model:** We assume an asynchronous communication model, where transmissions from different processes do not have to start at the same time. However, all processes start the algorithm simultaneously. This synchronization model corresponds to the execution of the MPI\_Alltoall operation, used as reference in this work.

### 4 Problem Definition

In the *total exchange* problem,  $n$  different processes hold each one  $n$  data items that should be evenly distributed among the  $n$  processes, including itself. Because each data item has potentially different contents and sizes according to their destinations, all processes engage a total exchange communication pattern. Therefore, a total exchange operation will be complete only after all processes have sent their messages to their counterparts, and received their respective messages.

Formally, the *total exchange problem* can be described using a weighted digraph  $dG(V, E)$  of order  $n$  with  $V = \{p_0, \dots, p_{n-1}\}$ . This digraph is called a message exchange digraph or MED for short. In a MED, the vertices represent the process nodes, and the arcs represent the messages to be transmitted. An integer  $w(e)$  is associated with each arc  $e = (p_i, p_j)$ , representing the size of the message to be sent from process  $p_i$  to process  $p_j$ . Note that there is not necessarily any relationship between a MED and the topology of the interconnection network.

The port capacity of a process for transmission is the number of other processes to which it can transmit simultaneously. Similarly, the port capacity for reception is the number of other processes from which it can receive simultaneously. We will concentrate on the performance modeling problem with all port capacities restricted to one for both transmitting and receiving. This restriction is well-known in the literature as *1-port full-duplex*.

## 4.1 Notation and lower bounds

In this section, we present theoretical bounds on the minimum number of communications and on the bandwidth for the general message exchange problem. The number of communications determines the number of start-ups, and the bandwidth depends on the message weights.

Given a MED  $dG(V; E)$ , we denote the in-degree of each vertex  $p_i \in V$  by  $\Delta_r(p_i)$ , and the out-degree by  $\Delta_s(p_i)$ . Let  $\Delta_r = \max_{p_i \in V} \{\Delta_r(p_i)\}$  and  $\Delta_s = \max_{p_i \in V} \{\Delta_s(p_i)\}$ . Therefore, we obtain the following straightforward bound on the number of start-ups.

**Claim 1.** *The number of start-ups needed to solve a message exchange problem on a digraph  $dG(V; E)$  without message forwarding is at least  $\max(\Delta_s, \Delta_r)$ .*

Given a MED  $dG(V, E)$ , the bandwidth bounds are determined by two obvious bottlenecks for each vertex - the time for it to send its messages and the time for it to receive its messages. Each vertex  $p_i$  has to send messages with sizes  $\{w_{i,j} \mid j = 0 \dots n-1\}$ . The time for all vertices to send their messages is at least  $t_s = \max_i \sum_{j=0}^{n-1} w_{i,j} \beta$ . Similarly, the time for all vertices to receive their messages is at least  $t_r = \max_j \sum_{i=0}^{n-1} w_{i,j} \beta$ .

**Claim 2.** *The time to complete a personalized exchange is at least  $\max\{t_s, t_r\}$ .*

We can combine the claims about the number of start-ups and the bandwidth when message forwarding is not allowed.

**Claim 3.** *If message forwarding is not allowed, and either the model is synchronous or both maxima are due to the same process, the time to complete a personalized exchange is at least  $\max(\Delta_s, \Delta_r) \times \alpha + \max\{t_s, t_r\}$ .*

Because in this paper we do not assume messages forwarding, the fan-in and fan-out of a process must be  $(n-1)$ . Further, as we consider messages to be the same size and the network to be homogeneous, we can simplify Claim 3 so that the following bound holds.

**Proposition 1.** *If message forwarding is not allowed, and all messages have size  $m$ , and both bandwidth and latency are identical to any connection between two different processes  $p_i$  and  $p_j$ , the time to complete a total exchange is at least  $(n-1) \times \alpha + (n-1) \times \beta m$ .*

*Proof.* The proof is trivial, as the time to complete a total exchange is at least the time a single process needs to send one message to each other process.

■

## 5 Contention Signature Approach

To cope with this problem and to model the contention impact on the performance of the All-to-All operation, we adopt an approach similar to Clement *et al.* [12], which considers the contention sufficiently linear to be modeled. Our approach, however, tries to identify the behavior of the All-to-All operation with regard to the theoretical lower bound (Proposition 1) on the *1-port* communication model. In our hypothesis, the network contention depends mostly on the physical characteristics of the network (network cards, links, switches),

and consequently, the ratio between the theoretical lower bound and the real performance represents a “*contention signature*” of the network. Once identified the *signature* of a network, it can be used in further experiments to predict the communication performance, provided that the network infrastructure does not change.

Initially, we consider communication in a contention-free environment. In this case, a process that sends messages of size  $m$  to  $n - 1$  processes needs at least  $(n - 1) \times \alpha + (n - 1) \times m\beta$  time units. Further, by the properties of the *1-port* communication model, the total communication time of the All-to-All operation must be at least  $(n - 1) \times \alpha + (n - 1) \times m\beta$  time units if all processes start communicating simultaneously, as stated by Proposition 1.

In the case of the All-to-All operation, however, the intensive communication pattern tends to saturate the network, causing message delays and packet loss that strongly impact on the communication performance of this collective communication. In this network congestion situation, traditional models such as those presented by Christara [1] do not hold anymore, even if the communication pattern has not changed.

Therefore, our approach to model the performance of the MPI\_Alltoall operation despite network contention consists on determining a *contention ratio*  $\gamma$  that express the relationship between the theoretical performance (lower bound) and the real completion time. For simplicity, we consider that this *contention ratio*  $\gamma$  is constant and depends exclusively on the network characteristics. Therefore, the simplest way to integrate this *contention ratio*  $\gamma$  in our performance model would be as follows:

$$T = (n - 1) \times (\alpha + m\beta \times \gamma) \quad (4)$$

### 5.1 Non-linear aspects

Although the performance model augmented by use of the *contention ratio*  $\gamma$  improves the accuracy of the predictions, we observe nonetheless that some network architectures are still subject to performance variations according to the message size. To illustrate this problem, we present in Fig. 1, a detailed mapping of the communication time of the MPI\_Alltoall operation in a Gigabit Ethernet network. We observe that the communication time does not increase linearly with the message size, but instead, present a non-linear behavior that prevents our model to accurately predict the performance when dealing with small messages.

To cope with this non-linearity, we propose an extension of the *contention ratio* model to better represent this phenomenon when messages are sufficiently large. Hence, we augment the model with a new parameter  $\delta$ , which depends on the number of processes but also on a given message size  $M$ . As a consequence, the association of different equations helps to define a more realistic performance model for the MPI\_Alltoall operation, as follows:

$$T = \begin{cases} (n - 1) \times (\alpha + m\beta \times \gamma) & \text{if } m < M \\ (n - 1) \times (\alpha + m\beta \times \gamma + \delta) & \text{if } m \geq M \end{cases} \quad (5)$$

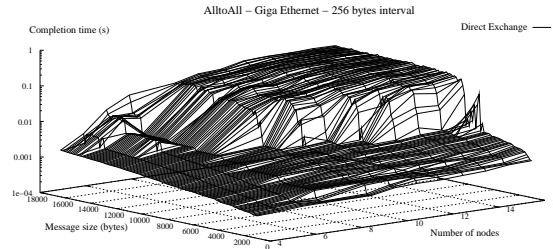


Fig. 1. Non-linearity of communication cost with small messages

## 6 Validation

To validate the approach proposed in this paper, this section presents our experiments to model the performance of MPI\_Alltoall operation using three network architectures, Fast Ethernet, Gigabit Ethernet and Myrinet. As previously explained, our approach consists on comparing the expected and real performance of the MPI\_Alltoall operation using a sample experiment with  $n'$  nodes; the relationship between these two measures allows us to define the  $\gamma$  and  $\delta$  parameters that characterize the "network contention signature".

To obtain these parameters, we compare the sample data obtained from both theoretical lower bound and experimental measure, when varying the message size. Indeed, the lower bound comes from Proposition 1, with parameters  $\alpha$  and  $\beta$  obtained from a simple point-to-point measure. The parameters  $\gamma$  and  $\delta$  are obtained through a linear regression with the Generalized Least Squares method, comparing at least four measurement points in order to better fit the performance curve.

The different experiments presented in this paper represent the average of 100 measures for each set of parameters (message size, number of processes), and were conducted over two clusters of the Grid'5000<sup>1</sup>:

**The *icluster2* cluster**, located at INRIA-Rhone-Alpes, composed of 104 dual Itanium2 nodes at 900 MHz, used for the experiments with the Fast Ethernet network (5 Fast Ethernet switches - 20 nodes per switch - interconnected by 1 Gigabit Ethernet switch) and the Myrinet 2000 network (one 128 ports M3-E128 Myrinet switch). All machines run Red Hat Enterprise Linux AS release 3, with kernel version 2.4.21.

**The *GdX* (Grid'eXplorer) cluster**, operated by INRIA-Futurs. This cluster includes 216 nodes with dual AMD Opteron processors at 2 GHz running Debian Linux kernel 2.6.8 and a Broadcom Gigabit Ethernet network.

### 6.1 Fast Ethernet

Taking as basis the measured performance for a 24 machines network, we were able to approximate the performance of the Fast Ethernet network with a *con-*

<sup>1</sup> <http://www.grid5000.org/>



contention ratio  $\gamma = 1.0195$ . Indeed, this relatively small difference must be considered in the light of the retransmission policy: although the communication latency (and therefore the timeouts) is relatively small (around  $60 \mu s$ ), the reduced bandwidth of the links minimizes the impact of the retransmission of a lost packet. More important, we observe that the experimental measure behave like an affine equation, showing a start-up cost usually not considered by the traditional performance model which corresponds to the  $\delta$  parameter proposed in our model. Therefore, we determined  $\delta = 8.23 ms$  for messages larger than  $M = 2 kB$ , which means that each simultaneous communication induces an overload of  $8.23 ms$  to the completion time of the All-to-All operation. Applying both  $\gamma$  and  $\delta$  parameters we were able to approximate our predictions from the performance of the MPI Alltoall operation with an arbitrary number of processes, as demonstrate in Fig. 2a. We observe indeed that our error rate is usually smaller than 10% when there are enough processes to saturate the network, as presented in Fig. 2b.

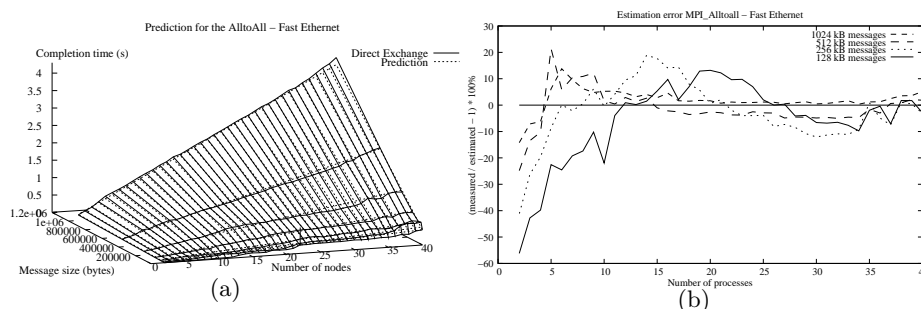


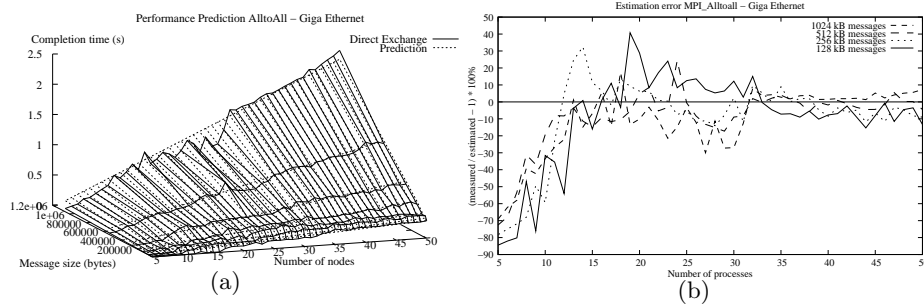
Fig. 2. Performance prediction on a Fast Ethernet network

## 6.2 Gigabit Ethernet

To compute the *contention ratio*  $\gamma$  and a *start-up cost*  $\delta$ , we use sample data for an arbitrary number of processes. Indeed, we chose in this example the results for an execution of the All-to-All operation with 40 processes (one by machine). Using linear regression on these data we obtain  $\gamma = 4.3628$  and  $\delta = 4.93 ms$  (to be used only for messages larger than  $M = 8 kB$ ). As a result, the performance predictions from our model correspond to the curve presented on Fig. 3a. As in the case of the Fast Ethernet network, the error rate is quite small when the network becomes saturate, even when we consider different message sizes (Fig. 3b).

## 6.3 Myrinet

Although the two previous experiments give important proofs on the validity of our modeling method, they share many similarities on both network architecture and transport protocol (TCP/IP). To ensure that our method is not bounded to

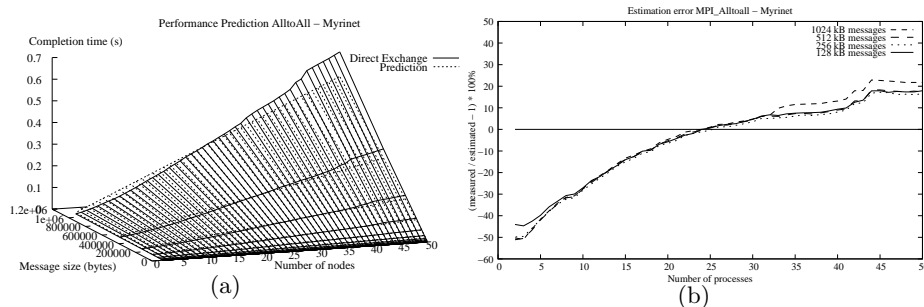


**Fig. 3.** Performance prediction on a Gigabit Ethernet network

a specific infrastructure, we chose to validate our performance model also in a Myrinet network, using the *gm* transport protocol. Because of the *Myrinet+gm* stack differs considerably from the Ethernet+TCP/IP stack, any systematic behavior introduced into our sampling data by these architectures should be exposed.

Indeed, the Myrinet network differs from Ethernet-based architectures due to an *start-up* cost almost inexistent (one of the main characteristics of the *Myrinet+gm* stack). Indeed, we were able to fit the performance of a 24-processes All-to-All operation using only the *contention ratio*  $\gamma = 2,49754$  (as the linear regression pointed a *start-up cost*  $\delta$  smaller than 1 microsecond).

Nevertheless, when applying this factor to an arbitrary number of machines, as presented in Fig. 4a, we observe that our predictions do not follow the experimental data as observed before with Fast Ethernet and Gigabit Ethernet. Actually, a close look at the error rate (Fig. 4b) indicates that network saturation occurs only when there are more than 40 communicating processes (evidenced by the constant error rate from that point). These results demonstrate the limitations of our approach: while a *contention ratio* may provide precise performance predictions, it depends on the data used to define the network signature. By using reference data from a partially saturated network we are subjected to inaccurate approximations (even if they are better than the contention unaware predictions).



**Fig. 4.** Performance prediction on a Myrinet network

## 7 Applications to Grid-Aware Communications

Actually, most of the complexity of the All-to-All problem in grid environments resides on the need to exchange different messages through different networks (local and distant). The traditional implementation of the MPI\_Alltoall operation cannot differentiate these networks, leading to poor performances. However, if we assume that communications between clusters are slower than intra-clusters ones, it might be useful to collect data in the local level before sending it through the backbone, in a single transmission. Indeed, in [18] we propose a grid-aware solution which performs on two phases. In the first phase only local communications are performed. During this phase the total exchange is performed on local nodes on both cluster and extra buffers are prepared for the second (inter-cluster) phase. During the second phase data are exchanged between the clusters. Buffers that have been prepared during the first phase are sent directly to the corresponding nodes in order to complete the total exchange.

More precisely, our algorithm works as follow. Without loss of generality, let us assume that cluster  $\mathcal{C}_1$  has less nodes than  $\mathcal{C}_2$  ( $n_1 \leq n_2$ ). Nodes are numbered from 0 to  $n_1 + n_2 - 1$ , with nodes from 0 to  $n_1 - 1$  being on  $\mathcal{C}_1$  and nodes from  $n_1$  to  $n_1 + n_2 - 1$  being on cluster  $\mathcal{C}_2$ . We call  $\mathcal{M}_{i,j}$  the message (data) that has to be send form node  $i$  to node  $j$ . For instance, the algorithm proceeds in two phases:

**First phase** During the first phase, we perform the local exchange: Process  $i$  sends  $\mathcal{M}_{i,j}$  to process  $j$ , if  $i$  and  $j$  are on the same cluster. Then it prepares the buffers for the remote communications. On  $\mathcal{C}_1$  data that have to be send to node  $j$  on  $\mathcal{C}_2$  is first stored to node  $j \bmod n_1$ . Data to be sent from node  $i$  on  $\mathcal{C}_2$  to node  $j$  on  $\mathcal{C}_1$  is stored on node  $\lfloor i/n_1 \rfloor \times n_1 + j$ .

**Second phase** During the second phase only  $n_2$  inter-cluster communications occurs. This phase is decomposed in  $\lceil n_2/n_1 \rceil$  steps with at most  $n_1$  communications each. Steps are numbered from 1 to  $\lceil n_2/n_1 \rceil$ . During step  $s$  node  $i$  of  $\mathcal{C}_1$  exchange data stored in its local buffer with node  $j = i + n_1 \times s$  on  $\mathcal{C}_2$  (if  $j < n_1 + n_2$ ). More precisely  $i$  sends  $\mathcal{M}_{k,j}$  to  $j$  where  $k \in [0, n_1]$  and  $j$  sends  $\mathcal{M}_{k,i}$  to  $i$  where  $k \in [n_1 \times s, n_1 \times s + n_1 - 1]$ .

As our algorithm minimizes the number of inter-cluster communications between the clusters, we need only  $2 \times \max(n_1, n_2)$  messages in both directions (against  $2 \times n_1 \times n_2$  messages in the traditional algorithm). For instance, the exchange of data between two clusters with the same number of process will proceed in one single communication step of the second phase. Our algorithm is also wide-area optimal since it ensures that a data segment is transferred only once between two clusters separate by a wide-area link.

### 7.1 Performance prediction in a grid environment

As shown above, the algorithm we propose to optimize All-to-All communications in a grid environment rely on the relative performances of both local and

remote networks. Indeed, we extend the total exchange among nodes in the same cluster in order to reduce transmissions through the backbone.

This approach has therefore two consequences for performance prediction: First, it prevents contention in the wide-area links, which are hard to model. Second, transmission of messages packed together is easy to be predicted in a wide-area network (large messages are less subjected to network interferences). For instance, we can design a wide-area performance model by composing local-area predictions obtained with our performance model and wide-area predictions that can be easily obtained from traditional methods. Hence, an approximate model for the communication between two clusters would be similar to:

$$T = \max(T_{C_1}, T_{C_2}) + \lceil n_2/n_1 \rceil \times (\alpha_w + \beta_w \times m \times n) \quad (6)$$

Although not in the scope of this work, preliminary experiments indicate that this model holds. We expect to develop this subject in a future work.

## 8 Conclusions and Future Works

In this paper we address the problem of modeling the performance of *Total Exchange* communication operations, usually subject to important variations caused by network contention. Because traditional performance models are unable to predict the real completion time of an All-to-All operation, we try to cope with this problem by identifying the *contention signature* of a given network. In our approach, two parameters  $\gamma$  and  $\delta$  are used to augment a linear performance model in order to fit the performance of the MPI\_Alltoall operation. Because these parameters characterize the network contention and are independent of the number of communicating processes, they can be used to accurately predict the communication performance when communications tend to saturate the network. Indeed, we demonstrate our approach through experiments conducted on popular network architectures, Fast Ethernet, Gigabit Ethernet and Myrinet.

We intend to pursue this research by validating our model under other network architectures like Infiniband. Indeed, we expect to extend our models to other collective communication operations, which are especially affected by contention when scaling up to a grid level. We are also investigating different strategies to model collective communications in grid environments.

## Acknowledgments

We are grateful to the anonymous referees for many comments and helpful suggestions which helped us improve the focus of the paper.

## References

1. Christara, C., Ding, X., Jackson, K.: An efficient transposition algorithm for distributed memory computers. In: Proceedings of the High Performance Computing Systems and Applications. (1999) 349–368

2. Midorikawa, E.T., Oliveira, H.M., Laine, J.M.: PEMPIs: A new methodology for modeling and prediction of MPI programs performance. In: Proceedings of the SBAC-PAD 2004, IEEE Computer Society/Brazilian Computer Society (2004) 254–261
3. Barchet-Steffenel, L.A., Mounie, G.: Scheduling heuristics for efficient broadcast operations on grid environments. In: Proceedings of the Performance Modeling, Evaluation and Optimization of Parallel and Distributed Systems Workshop - PME0'06 (associated to IPDPS'06), Rhodes Island, Greece, IEEE Computer Society (April 2006)
4. Kielmann, T., Bal, H., Gortlatch, S., Verstoep, K., Hofman, R.: Network performance-aware collective communication for clustered wide area systems. *Parallel Computing* **27**(11) (2001) 1431–1456
5. Chun, A.T.T., Wang, C.L.: Realistic communication model for parallel computing on cluster. In: Proceedings of the International Workshop on Cluster Computing. (1999) 92–101
6. Chun, A.T.T.: Performance Studies of High-Speed Communication on Commodity Cluster. PhD thesis, University of Hong Kong (2001)
7. Pjesivac-Grbovic, J., Angskun, T., Bosilca, G., Fagg, G.E., Gabriel, E., Dongarra, J.J.: Performance analysis of MPI collective operations. In: Proceedings of the Workshop on Performance Modeling, Evaluation and Optimisation for Parallel and Distributed Systems (PME0), in IPDPS 2005. (2005)
8. Johnsson, S.L., Ho, C.T.: Optimum broadcasting and personalized communication in hypercubes. *IEEE Transactions on Computers* **38**(9) (September 1989) 1249–1268
9. Grove, D.: Performance Modelling of Message-Passing Parallel Programs. PhD thesis, University of Adelaide (2003)
10. Adve, V.: Analysing the Behavior and Performance of Parallel Programs. PhD thesis, University of Wisconsin, Computer Sciences Department (1993)
11. Bruck, J., Ho, C.T., Kipnis, S., Upfal, E., Weathersby, D.: Efficient algorithms for all-to-all communications in multiport message-passing systems. *IEEE Transactions on Parallel and Distributed Systems* **8**(11) (November 1997) 1143–1156
12. Clement, M., Steed, M., Crandall, P.: Network performance modelling for PM clusters. In: Proceedings of Supercomputing. (1996)
13. Labarta, J., Girona, S., Pillet, V., Cortes, T., Gregoris, L.: DiP: A parallel program development environment. In: Proceedings of the 2nd Euro-Par Conference. Volume 2. (1996) 665–674
14. König, J.C., Rao, P.S., Trystram, D.: Analysis of gossiping algorithms with restricted buffers. *Parallel Algorithms and Applications* **13**(2) (feb 1998) 117–133
15. Jeannot, E., Wagner, F.: Two fast and efficient message scheduling algorithms for data redistribution through a backbone. In: Proceedings of the IPDPS. (2004)
16. Moritz, C.A., Frank, M.I.: LoGPC: Modeling network contention in message-passing programs. *IEEE Transactions on Parallel and Distributed Systems* **12**(4) (2001) 404–415
17. Hockney, R.: The communication challenge for MPP: Intel paragon and meiko cs-2. *Parallel Computing* **20** (1994) 389–398
18. Jeannot, E., Steffenel, L.A.: Fast and efficient total exchange on two clusters. Submitted to EuroPar'07 - 13th International Euro-Par Conference European Conference on Parallel and Distributed Computing