

Fast and Scalable Startup of MPI Programs in InfiniBand Clusters^{*}

Weikuan Yu, Jiasheng Wu, Dhabaleswar K. Panda

Network-Based Computing Lab
Dept. of Computer Science and Engineering
The Ohio State University
{yuw,wuj,panda}@cse.ohio-state.edu

Abstract. One of the major challenges in parallel computing over large scale clusters is fast and scalable process startup, which typically can be divided into two phases: process initiation and connection setup. In this paper, we characterize the startup of MPI programs in InfiniBand clusters and identify two startup scalability issues: serialized process initiation in the initiation phase and high communication overhead in the connection setup phase. To reduce the connection setup time, we have developed one approach with data reassembly to reduce data volume, and another with a bootstrap channel to parallelize the communication. Furthermore, a process management framework, Multi-Purpose Daemons (MPD) system is exploited to speed up process initiation. Our experimental results show that job startup time has been improved by more than 4 times for 128-process jobs, and the improvement can be more than two orders of magnitude for 2048-process jobs as suggested by our analytical models.

1 Introduction

The MPI (Message Passing Interface) Standard [12] has evolved as a *de facto* parallel programming model for distributed memory systems. Traditional research over MPI has been largely focusing on the high performance communication between processes. As cluster computing becomes a prominent platform of high performance computing, scalable process management of MPI applications becomes an active research topic [3, 1]. One of the major challenges in process management is the fast and scalable startup of large-scale applications [2, 6, 10, 4, 9]. This issue becomes even more pronounced in the large scale systems with thousands of nodes. A parallel job is usually launched by a process manager, which is often referred to as the *process initiation phase*. These initiated processes usually require assistance from the process manager to set up peer-to-peer connections before starting communication and computation. This is referred to as the *connection setup phase*.

InfiniBand Architecture (IBA) [8] has been recently standardized in industry to design next generation high-end clusters for both data-center and high performance computing. Large cluster systems with InfiniBand are being deployed. For example, in the Top500 list released in November 2003 [15], the 3rd, 111th, and 116th most powerful supercomputers use InfiniBand as their parallel application communication interconnect. These three systems have 2200, 256, and 512 processors, respectively. The startup

^{*} This research is supported in part by a DOE grant #DE-FC02-01ER25506, NSF Grants #CCR-0204429 and #CCR-0311542, and a grant from Los Alamos National Laboratory.

of MPI applications in InfiniBand clusters at such a large scale is a challenging issue. It may take more than ten minutes to go through the above mentioned process initiation and connection setup phases for an application with 1000 processes without scalable and high performance startup support.

In this paper, we have taken on the challenge to support a scalable and high performance startup of MPI programs over InfiniBand clusters. With MVAPICH [13] as the platform of study, we have analyzed the startup bottlenecks. Accordingly, different approaches have been developed to speed up the connection setup phase, one with data reassembly at the process manager and another using pipelined all-to-all broadcast over a ring of InfiniBand queue pairs (referred to as a bootstrap channel). In addition, we have exploited a process management framework, Multi-Purpose Daemons (MPD) system to further speed up the startup. The bootstrap channel is also utilized to reduce the impact of communication bottlenecks in MPD, including multiple process context switches and quadratically increasing data volume over the MPD management ring. Over 128 processes, our work improves the startup time by more than 4 times. Scalability Models derived from these results suggest that the improvement can be more than two orders of magnitude for the startup of 2048-process jobs.

The rest of the paper is structured as follows. Section 2 gives an overview of InfiniBand. Section 3 describes the challenge of scalable startup faced by parallel programs over InfiniBand and related work on process management. Section 4 describes the design of startup with different approaches to improve the connection setup time and the process initiation phase. Experiments results are provided in 5. Finally, we conclude the paper in Section 6.

2 Overview of InfiniBand Architecture

The InfiniBand Architecture (IBA) [8] defines a System Area Network (SAN) for interconnecting computing nodes and I/O nodes. In an InfiniBand network, a switched communication fabric is defined to allow many devices to communicate concurrently at high bandwidth and low latency. Processing nodes are connected as end-nodes to the fabric with Host Channel Adapters (HCAs).

InfiniBand provides four types of transport services: Reliable Connection (RC), Reliable Datagram (RD), Unreliable Connection (UC), and Unreliable Datagram (UD). The often used service is RC in the current InfiniBand products and software. It is also our focus in this paper. To support RC, a connection must be set up between two QPs before any communication. In the current InfiniBand SDK, each QP has a unique identifier, called *QP-ID*. This is usually an integer. For network identification, each HCA also has a unique 16-bit local identifier (*LID*). To make a connection, a pair of QPs must exchange their QP IDs and LIDs.

3 Problem Statement and Related Work

This section first characterizes the scalability constraints of the startup of MPI programs in InfiniBand clusters. It then provides a brief discussion of related work and motivates the study for a scalable startup scheme.

3.1 Startup of MPI Applications using MVAPICH

MVAPICH [13] is a high performance implementation of MPI over InfiniBand. Its design is based on MPICH [5] and MVICH [11]. The current implementation of MVAPICH utilizes the Reliable Connection (RC) service for the communication between processes. The connection-oriented nature of IBA RC-based QPs requires each process to create at least one QP for every peer process. To form a fully connected network of N processes, a parallel application needs to create and connect at least $N \times (N - 1)$ QPs during the initialization time. Note that it is possible to have these QPs be allocated and connected in an on-demand manner [16], however this requires that the connection management subsystem of IBA can handle either peer-to-peer or client-server model connection establishment, which is not mature yet in the current IBA software. Another reason for the fully-connected connection model is simplicity and robustness. Therefore, this connection model has been used in many MPI implementations, including MVAPICH.

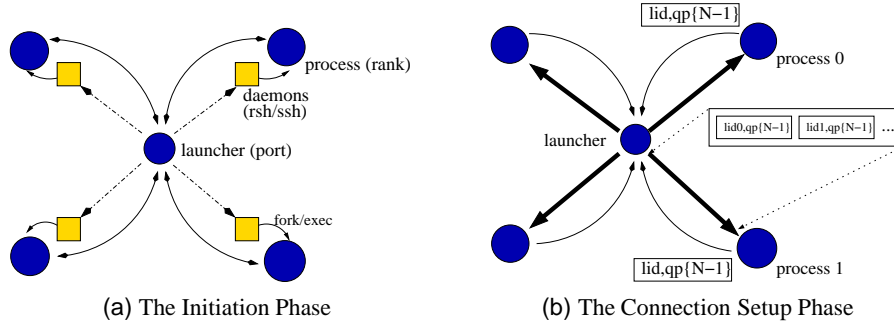


Fig. 1: The Startup of MPI Applications in Current MVAPICH

The startup of an MPI application using MVAPICH can also be divided into two phases. As shown in Fig. 1(a), an MPI application using MVAPICH is launched with a simple process launcher iterating over UNIX remote shell (rsh) or secure shell (ssh) to start individual processes. Each process connects back to the launcher via a port exposed by the launcher. Except the rank of the process, each process has no global knowledge about the parallel program. In the second phase of connection setup, as shown in Fig. 1(b), each process creates $N - 1$ QPs, one for each peer process, for an N -process application. Then, these processes exchange their local identifiers (LIDs) and corresponding QP identifiers (QP-IDs), as mentioned in Section 2 for connection setup. Since each process is not connected to its peer processes, the data exchange has to rely on the connections that are created to the launcher in the first phase. The launcher collects data about LIDs and QP-IDs from each process, and then sends the combined data back to each process. Each process in turn sets up connections over InfiniBand with the received data. A parallel application with fully connected processes is then created.

3.2 The Scalability Problem

The startup paradigm described above is able to handle the startup of small scale parallel applications. However, as the size of an InfiniBand cluster goes to 100s–1000s,

the limitation of this paradigm becomes pronounced. For example, launching a parallel application with 2000 processes may take tens of minutes. There are two main scalability bottlenecks, one in each phase. The first bottleneck is *rsh/ssh-based startup* in the process initiation phase. This process startup mechanism is simple and straightforward, but its performance is very poor on large systems. The second bottleneck is the communication overhead for exchanging LIDs and QP-IDs in the connection setup phase. To launch an N-process MPI application, the launcher has to receive data containing $(N - 1)$ QP-IDs from each process. Then it returns the combined data with $N \times (N - 1)$ QP-IDs to each process. In total, the launcher has to communicate data in the amount of $O(N^3)$ for an N-process application. Each QP-ID is usually a four-byte integer, for a 1024-process application the launcher will receive almost 4 MegaBytes data and sends almost 4 Gigabytes of data. This communication typically goes through the management network which is normally Fast Ethernet or Gigabit Ethernet. This incurs significant communication overhead and slowdown to the application startup.

3.3 Related Work

Numerous work have been done to provide resource management framework for collections of parallel processes, ranging from basic iterative rsh/ssh-based process launch in MVICH [11] to more sophisticated packages like MPD [3], Cplant [2], PBS [14], LoadLeveler/POE [7], to name a few. Compared to the rsh/ssh-based iterative launch of processes, all these packages can provide more scalable startup and retain better monitoring and control of parallel programs. However, they typically lack efficient support for complete exchange of LIDs and QP-IDs as required by parallel programs over InfiniBand clusters. In this paper, we focus on providing an efficient support for the complete exchange of LIDs and QP-IDs, and applying such a scheme to one of these package, MPD, in order to obtain efficient process initiation support. We choose to study MPD [3] because it is one of the systems widely distributed along with MPICH [5] releases and has a large user base.

4 Designing Scalable Startup Schemes

This section describes the design of scalable startup schemes in InfiniBand clusters. We first describe different approaches used to enhance the connection setup phase while the processes are still launched via rsh/ssh daemons. Then we exploit the advantages of *MPD* [3], to replace the rsh/ssh based scheme and achieve efficient process initiation. We also characterize some MPD features and their limitations to the scalable startup of MPI applications in InfiniBand clusters. We also introduce the concept of a bootstrap channel which can be used to overcome these limitations.

4.1 Efficient Connection Setup

As mentioned in the previous section, because the launcher has to collect, combine and broadcast QP IDs, the volume of these data scales up in the order of $O(N^3)$, which leads to prolonged connection setup time. One needs to consider two directions in order to reduce the connection setup time. The first direction is to reduce the volume of data that needs to be communicated. The other direction is to parallelize communication for the exchange of QP IDs.

Approach 1: Reducing the Data Volume with Data Reassembly (DR) To have processes fully connected over InfiniBand, each process needs to connect with another peer process via one QP. This means that each process needs to obtain $N - 1$ QP IDs, one for each peer. That is to say, out of the combined data of $N \times (N - 1)$ QP IDs in the launcher, each process only needs to receive $N - 1$ QP IDs that is specifically targeted for itself. This requires a centralized component, i.e., the launcher, to collect and re-assemble QP IDs. The biggest advantage of this data reassembly (DR) scheme is that the data volume exchanged can be reduced down to an order of $O(N^2)$. But there are several disadvantages associated with this scheme. First, the entire set of QP IDs need to be reassembled before sending them to each client processes. This constitutes another performance/scalability bottleneck at the launcher. Second, the whole procedure of receive-reassembly-send is also serialized at the launcher.

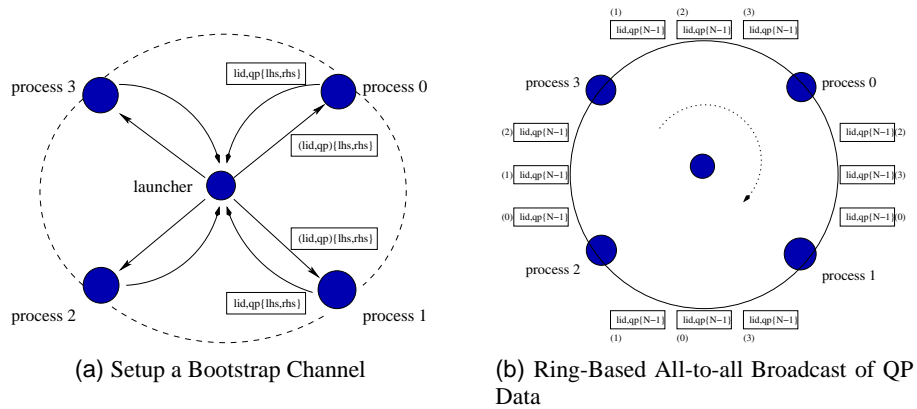


Fig. 2: Parallelizing the Total Exchange of InfiniBand Queue Pair Data

Approach 2: Parallelizing Communication with a Bootstrap Channel (BC) More insights can be gained on the possible parallelism with further examination of the startup. Essentially, what needs to be achieved at the startup time is an all-to-all personalized exchange of QP IDs, i.e., each process receives the specific QP IDs from other processes. In the original startup scheme as shown in Fig. 1, the launcher performs a gather/broadcast to help the all-to-all broadcast of their QP data. On top of that, the DR scheme in Section 4.1 reassembles and “personalizes” QP data to reduce the data volume. Both do not exploit the parallelism of all-to-all personalized exchange. Algorithms that parallelize an all-to-all personalized exchange can be used here. These algorithms are usually based on a ring-, hypercube- or torus-based topology, which requires more connections to be provided among processes. With the initial star topology in the original startup scheme, providing these connections has to be done through the launcher. However, since a parallel algorithm can potentially overlap both sending and receiving QP data, it promises better scalability over clusters with larger sizes.

Among the three possible parallel topologies, the ring-based topology requires the least number of additional connections, i.e., 2 per process. This would minimize the impact of the ring setup time. Another design option to be considered is that which type

of connections should be provided. Either TCP/IP- or InfiniBand-based connections can be used. Since the communication over InfiniBand is much faster than that over TCP/IP (see [17] for detail latency comparison between them), we choose to use a ring of InfiniBand QPs as a further boost to the parallelized data exchange.

The second approach works as follows. First, each process creates two QPs for its left hand side (lhs) and right hand side (rhs) processes, respectively. We call these QPs *bootstrap QPs*. Second, the DR scheme mentioned in Section 4.1 is used to set up connections between these bootstrap QPs as shown in Figure 2(a). Thus, a ring of connections over InfiniBand is created, as shown by the dotted line in Figure 2(a). We refer to this ring as a *bootstrap channel (BC)*. After this channel is set up, each process initiates a broadcast of its own QP IDs through the channel in the clockwise direction as shown in Fig. 2(b) with four processes. Each process also forwards what it receives to its next process. In this scheme, we take advantage of both communication parallelism and high performance of InfiniBand QPs to reduce the communication overhead.

4.2 Fast Process Initiation with MPD

MPD [3] is designed to be a general process manager interface that provides the needed support for MPICH, from which MVAPICH is developed. It mainly provides fast startup of parallel applications and process control to the parallel jobs. MPD achieves its scalable startup by instantly spreading a job launch request across its ring of daemons, then launches one ring of manager and another ring of application processes in a parallel fashion (see [17] for detailed description of MPD systems). For processes to exchange individual information MPD system also exposes a BNR interface with a put/fence/get model. A process stores (puts) a (key,value) pair at its manager process, a part of the MPD database, then another process retrieves (gets) that value by providing the same key after a synchronization phase (fence).

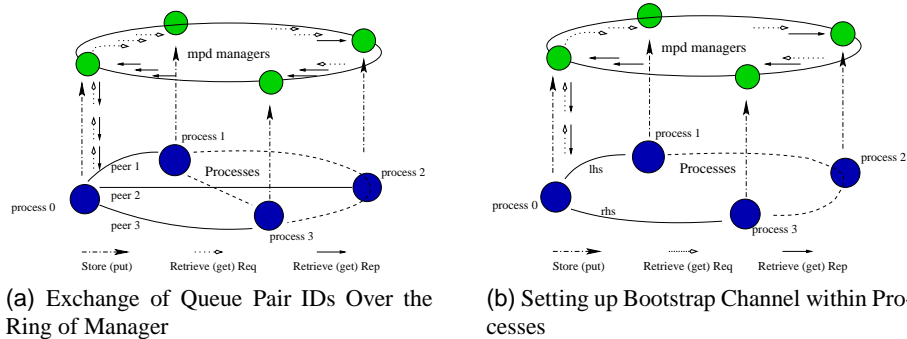


Fig. 3: Improving the Scalability of MPD-Based Startup

Although this fast and parallelized process startup from MPD solves the process initiation problem, the significant volume of QP data still poses a great challenge to the MPD model. As shown in Fig. 3(a), the database is distributed over the ring of manager processes when each process stores (puts) their process-specific data to its manager. To collect the data from every peer process, one process has to send a request and get the reply back for the target process. At the completion of these data exchanges,

each process then sets up connections with all the peers, as shown with process 0 in Fig. 3(a). Together, messages for the request and the reply make a complete round over the manager ring. For a parallel job with N processes, there are $N \times (N - 1)$ message exchanges in total. Each of these messages is in the order of $O(N)$ bytes and has to go through the ring of manager processes. In addition, since application processes store and retrieve data through their corresponding manager processes at each node, process context switches are very frequent and they further degrade the performance of ring-based communication. Furthermore, the message passing is over TCP/IP sockets, which delivers lower performance than InfiniBand-based connections (see [17] for latency comparisons).

There are different alternatives to overcome these limitations. One way of doing that is to replace the connections for the MPD manager ring with VAPI connections to provide fast communications. In addition, copies of QP data can be saved at each manager process as the first copy of QP data passes through the ring. Then further retrieve (get) requests can get the data from the local manager directly instead of the MPD manager ring. This approach will improve the communication time, however, the process context switches still exist between the application processes and manager processes. In addition, retrieve requests made before QP data reaches the local manager process still has to go through the manager ring. Last but not least, this approach necessitates a significant amount of instrumentation of MPD code and has only limited portability to InfiniBand-ready clusters.

Instead of exchanging all the QP data over the ring of MPD manager processes, we propose to exchange QP IDs over the bootstrap channel described in Section 4.1. Though setting up the bootstrap channel still needs help from the ring of manager processes. As shown in Fig. 3(b), each process first creates and stores QP IDs for its left side (lhs) and right hand side (rhs) processes to the local manager. Then, from the database, they retrieve QP IDs for its left hand side and right hand side processes, and then set up InfiniBand connections. Eventually a ring of such connections are constructed and together form a bootstrap channel. This bootstrap channel is then utilized to perform a complete exchange of QP IDs as described in Section 4.1. Since this bootstrap channel is provided within the application processes and over InfiniBand, this approach will not only provide fast communication and eliminate the process context switches, but also reduce the number of communications through each manager process.

5 Performance Evaluation

Our experiments were conducted on a 256-node cluster of 4GB DRAM dual-SMP 2.4GHz Xeon at the Ohio Supercomputing Center. For fast network discovery with data reassembly (DR) or the bootstrap channel (BC), we used ssh to launch the parallel processes. Performance comparisons were provided against MVAPICH 0.9.1 (Original). Since Networked File System (NFS) performance could be a big bottleneck in a large cluster and mask out the performance improvement of startup, all binary executable files were duplicated at local disks to eliminate its impact.

5.1 Experimental Results

Table 1 shows the startup time for parallel jobs of different number processes using different approaches. SSH-DR represents ssh-based startup with QP data assembly (DR)

Table 1: Comparisons of Parallel Job Startup Time over MVAICH with Different Approaches

Number of Processes	4	8	16	32	64	128
Original (sec)	0.59	0.92	1.74	3.41	7.3	13.7
SSH-DR (sec)	0.58	0.94	1.69	3.37	6.77	13.45
SSH-BC (sec)	0.61	0.95	1.70	3.38	6.76	13.3
MPD-BC (sec)	0.61	0.63	0.64	0.84	1.58	3.10

at the process launcher. SSH-BC represents ssh-based startup using the bootstrap channel (BC) to exchange QP IDs. MPD-BC represents MPD-based startup with a bootstrap channel for the exchange of QP IDs.

As the number of processes increases, both SSH-DR and SSH-BC reduce the startup time, compared to the original approach. This is because data reassembly can reduce the data volume by an order of $O(N)$ and the bootstrap channel can parallelize the communication time. Note that the BC-based approach performs slightly worse than the original and DR-based approach for small number of processes. This is due to the overhead from setting up the additional ring over InfiniBand. As the number of processes increases, the benefits become greater. Both SSH-BC and SSH-DR will be able to provide more scalable startup for a job with thousands of processes since they remove the major communication bottleneck imposed by potentially large volume of QP data. In contrast, the MPD-based approach with a bootstrap channel provides the most scalable startup. On one hand, MPD-BC provides efficient parallelized process initialization, compared to the ssh-based schemes. On the other hand, it also pipelines the QP data exchange over a ring of VAPI connections, hence this approach speeds up the connection setup phase. Compared to the original approach, the MPD-BC approach reduces the startup time for a 128-process job by more than 4 times.

5.2 Analytical Models and Evaluations for Large Clusters

As indicated by the results from Section 5.1, the benefits of the designed schemes will be more pronounced for parallel jobs with larger number of processes. In this section, we further analyze the performance of different startup schemes and provide parameterized models to gain insights about their scalability over large clusters. The total startup time $T_{startup}$ can be divided into the process initiation time and the connection setup time, denoted as T_{init} and T_{conn} respectively. Based on the scalability analysis, we use the following model to describe the startup time of the original scheme (Original), ssh-based scheme with data reassembly (SSH-DR) and the MPD-based scheme with the bootstrap channel (MPD-BC). Each of the models shows the time for the startup of N processes, and the last component describes the time for other overheads that are not quantified in the models, for example, process switching overhead.

Original: $T_{startup} = (O_0 * N) + (O_1 * N * (W_N + W_{N^2})) + O_2$

The process initiation phase time T_{init} scales linearly as the number of processes increases with ssh/rsh-based approaches, while during the connection setup there are $2N$ messages communicated over TCP/IP. Half of them are gathered by the launcher, each being in the order of $O(N)$ bytes; the other half are scattered by the launcher, each of $O(N^2)$ bytes.

SSH-DR: $T_{startup} = (D_0 * N) + (D_{comp} * N^3 + D_1 * 2N * W_N) + D_2$

The process initiation time T_{init} scales linearly with ssh/rsh. During the

connection setup phase, the amount of computation scales in the order of $O(N^3)$ (the constant D_{comp} can be very small, being the time for extracting one QP Id), and there are $2*N$ message communicated over TCP/IP. Half of them are gathered by the launcher, each being in the order of $O(N)$ bytes; The other half are scattered by the launcher, each of them is only $O(N)$ bytes due to reassembly.

MPD-BC: $T_{startup} = (M_0 + N * W_{req}) + (M_{ch_setup} * N + M_1 * N * W_N) + M_2$

The process initiation time T_{init} scales constantly using MPD, however there is a small fractional increase of communication time for the request message W_{req} . During the connection setup phase, the time to setup a bootstrap channel increases in the order of $O(N)$. Each process also handles N message in the pipeline, each in the order of $O(N)$ bytes.

Original: $T_{startup} \text{ (sec)} = (0.100 * N) + (10.5 * N * (W_N + W_{N^2})) + 0.12$

SSH-DR: $T_{startup} \text{ (sec)} = (0.100 * N) + (8.5e^{-9} * N^3 + 10.5 * N * W_N) + 0.12$

MPD-BC: $T_{startup} \text{ (sec)} = (0.20 + 0.0010 * N) + (0.0180 * N + 2.5 * N * W_N) + 0.30$

The above scalability models are parameterized based on our analytical modeling. As shown in Fig. 4, the experiment results confirm the validity of these models for jobs with 4 to 128 processes. Fig. 5 shows the scalability of different startup schemes when applying the same models to larger jobs from 4 to 2048 processes. Both SSH-DR and MPD-BC improves the scalability of job startup significantly. Note that MPD-BC scheme improves the startup time by about two orders of magnitudes for 2048-process jobs.

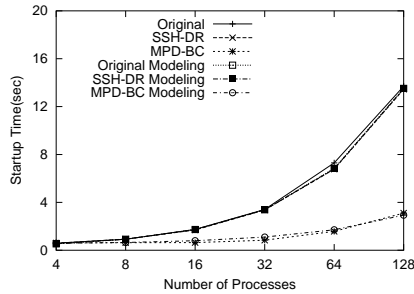


Fig. 4: Performance Modeling of Different Startup Schemes

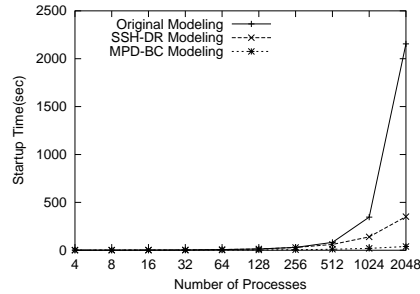


Fig. 5: Scalability Comparisons of Different Startup Schemes

6 Conclusions and Future Work

In this paper, we have presented schemes to support scalable startup of MPI programs in InfiniBand clusters. With MVAPICH as the platform of study, we have characterized the startup of MPI jobs into two phases: process initiation and connection setup. To speed up connection setup phase, we have developed two approaches, one with queue pair data reassembly at the launcher and the other with a bootstrap channel. In addition, we have exploited a process management framework, Multi-Purpose Daemons (MPD) system,

to improve the process initiation phase. The performance limitations in the MPD's ring-based data exchange model, such as exponentially increased communication time and numerous process context switches, are eliminated by using the proposed bootstrap channel. We have implemented these schemes in MVAPICH [13]. Our experimental results show that, for 128-process jobs, the startup time has been reduced by more than 4 times. We have also developed an analytical model to project the scalability of the startup schemes. The derived models suggest that the improvement can be more than two orders of magnitudes for the startup of 2048-process jobs with the MPD-BC startup scheme.

In future, we want to provide a file broadcast mechanism to MPD system to achieve efficient loading of jobs [10]. Furthermore, we intend to provide a hypercube-based scalable startup over really large systems, e.g., future Peta-scale clusters with tens of thousands of processors.

References

- [1] M. Baker, G. Fox, and H. Yau. Cluster Computing Review, November 1995.
- [2] R. Brightwell and L. A. Fisk. Scalable parallel application launch on Cplant. In *Proceedings of Supercomputing, 2001*, Denver, Colorado, November 2001.
- [3] R. Butler, W. Gropp, and E. Lusk. Components and interfaces of a process management system for parallel programs. *Parallel Computing*, 27(11):1417–1429, 2001.
- [4] E. Frachtenberg, F. Petrini, J. Fernandez, S. Pakin, and S. Coll. STORM: Lightning-Fast Resource Management. In *Proceedings of the Supercomputing '02*, Baltimore, MD, November 2002.
- [5] W. Gropp, E. Lusk, N. Doss, and A. Skjellum. A High-Performance, Portable Implementation of the MPI Message Passing Interface Standard. *Parallel Computing*, 22(6):789–828, 1996.
- [6] E. Hendriks. Bproc: The beowulf distributed process space. In *Proceedings of the International Conference on Supercomputing*, New York, New York, June 2002.
- [7] IBM. Using the Parallel Operating Environment, Version 4, Release 1, 2004.
- [8] Infiniband Trade Association. <http://www.infinibandta.org>, 2000.
- [9] M. Jette and M. Grondona. SLURM: Simple Linux Utility for Resource Management. In *Proceedings of the International Conference on Linux Clusters*, San Jose, CA, June 2003.
- [10] A. Kavas, D. Er-El, and D. G. Feitelson. Using Multicast to Pre-Load Jobs on the ParPar Cluster. *Parallel Computing*, 27(3):315–327, 2001.
- [11] Lawrence Berkeley National Laboratory. MVICH: MPI for Virtual Interface Architecture. <http://www.nersc.gov/research/FTG/mvich/index.html>, August 2001.
- [12] Message Passing Interface Forum. MPI: A message-passing interface standard. *The International Journal of Supercomputer Applications*, 8(3–4):159–416, 1994.
- [13] Network-Based Computing Laboratory. MVAPICH: MPI for InfiniBand on VAPI Layer. <http://nowlab.cis.ohio-state.edu/projects/mpi-iba/index.html>.
- [14] OpenPBS Documentation. <http://www.openpbs.org/docs.html>, 2004.
- [15] TOP 500 Supercomputers. <http://www.top500.org/>, 2003.
- [16] J. Wu, J. Liu, P. Wyckoff, and D. K. Panda. Impact of On-Demand Connection Management in MPI over VIA. In *Proceedings of the International Conference on Cluster Computing*, 2002.
- [17] W. Yu, J. Wu, and D. K. Panda. Fast and Scalable Startup of MPI Programs in InfiniBand Clusters. Number OSU-CISRC-5/04-TR33, Columbus, OH 43210, May 2004.