

Hybrid Preemptive Scheduling of MPI Applications on the Grids

Aurelien Bouteiller, Hinde-Lilia Bouziane, Thomas Herault, Pierre Lemarinier, Franck Cappello
INRIA/LRI, Université Paris-Sud, Orsay, France

E-mail: {bouteill, bouziane, herault, lemarini, fci}@lri.fr
phone: (+33) 1 69 15 4222
fax: (+33) 1 69 15 4213

Abstract—Time sharing between all the users of a Grid is a major issue in cluster and Grid integration. Classical Grid architecture involves a higher level scheduler which submits non overlapping jobs to the independant batch schedulers of each cluster. The sequentiality induced by classical batch scheduling is not fitted for the number of users and the heterogeneity of the jobs in Grids. Time sharing techniques address this issue by allowing the simultaneous execution of many applications on the same resources.

Classical sharing techniques are co-scheduling and gang-scheduling. Co-scheduling rely on the independant operating system of each node to schedule every process of every application. Gang-scheduling ensures that the same application is scheduled on all nodes simultaneously. Previous work has proven that co-scheduling techniques outperforms gang-scheduling while physical memory is not exhausted.

In this paper, we introduce a new hybrid sharing technique providing checkpoint based explicit memory management. It consists in co-scheduling some of the simultaneous applications up to the memory capacity, and use gang-scheduling related techniques to share fairly the execution time between all the applications. We compare experimentally the merits of the three solutions. The experiments show that the hybrid solution is as efficient as the co-scheduling technique when the physical memory is not exhausted, and is more efficient than gang-scheduling and co-scheduling when physical memory is exhausted.

I. INTRODUCTION

Two of the most fundamental principles of Grid are 1) the capacity to establish virtual organizations spanning over the several administration domains in order to extend the number of resources accessible by users and 2) the coordination of these resources in order to cooperate in the resolution of user problems. From these two principles eventually arises the need for efficient and fair resource sharing mechanisms.

The first principle inevitably tends to increase the pressure on the system mechanisms implementing the resource sharing between users. In the general situation, the resources are belonging to institutions and are already used by some of their members. Thus, building virtual organizations on top of already used resources increases the number of potential users for these resources, leading to an increased requirement to share these resources fairly among the users. In large system within HPC centers, queues of jobs are already quite long. It is not uncommon to wait days before having large jobs done. If nothing is done in Grids associating tens of such sites, the waiting time would certainly evolves from days to weeks which in some circumstances will not be acceptable

for users. An other fairness issue concerns the capacity to establish several level of priority between parallel executions on the Grid. High priority execution should be able to preempt the resource of an low priority parallel application under execution.

The second principle, as examined in the light of the first one, implies the use of efficient coordination mechanisms ensuring A) parallel application performance closed to the one obtained in a dedicated system and B) limiting strictly the waste of computing resources by providing rapid parallel applications context switch. In this paper, we restrict our investigation domain to users running MPI applications on clusters shared within a Grid virtual organization. We focus on the following scenario: users submit their MPI jobs to a meta-scheduler (from a portal) which schedules them on dynamically selected clusters. Any MPI execution may span over several clusters or stay within a single cluster. In the first case, an efficient coordination mechanism should schedule simultaneously all the components involved in each MPI execution. The efficiency in fulfilling criteria A) and B) depends on the synchronization mechanism coordinating the scheduling of the MPI subparts over all the clusters involved in the execution and on the speed of context switch between the MPI executions on each cluster. In the second case, clusters are not synchronized and the efficiency of the Grid only depends on the capability of each cluster management system to meet them criteria A) and B).

Fairness and performance are in principle contradictory. Reaching top performance on parallel execution involves dedicated usage of the cluster resources. The less the operating system interrupts the application execution, the highest is the performance. Usually, fairness relies on context switch mechanisms enabling the resource sharing between users. Context switching adds an overhead on the total application execution time. The more context switches are experienced during an execution, the more the application is slowed down.

Sharing fairly the cluster resources between multiple users may be examined in two cases: 1) when the concurrent executions of all users fit in the memory of the cluster nodes and 2) conversely when disc storage should be used in addition to the memory in the cluster nodes to store all concurrent executions. We will call these two cases respectively in-core and out-of-core context switch.

The two principles lead to three main consequences: 1) a

fast mechanism for switching the context of MPI execution on a cluster is the corner stone to meet efficiency criteria, 2) because fairness and performance are contradictory objectives, we should consider a metric representing a tradeoff between them. For sake of simplicity, in this paper, we will consider application with the same execution time and the following ratio as the tradeoff metric : the execution time of a set of parallel applications over the standard deviation of their individual execution times and 3) in a Grid with lot of users, it is likely that out-of-core context switch will be the general case. Thus in this paper, will essentially focus on this context.

In this paper we study several MPI application context switching techniques, trying to discover which one has the lowest impact on application performance. We will demonstrate that the best technique for out-of-core context switching is a hybrid (two level) one mixing checkpoint based context switching of sets of MPI executions and uncoordinated scheduling (co-scheduling) of MPI executions within each set.

The second section presents the related works. Section III presents the general framework used to compare the different scheduling techniques. Section IV presents the different scheduling approaches compared in this paper. Section V presents the experimental results and section VI concludes and sums up what we learned from the experiences.

II. RELATED WORK

There are several main techniques to implement parallel application context switch in practice: one of the most used one is to queue the jobs submitted by the different users and apply selection, analyzing the queued jobs. In the general cases, after being elected for execution, parallel jobs are scheduled subsequently (one after the other). Thus, only one parallel execution runs on the cluster at a given time. Another way is to let the operating system schedules the processes of several parallel executions launched concurrently, according to their priorities. These techniques are called gang scheduling and co-scheduling respectively, depending on the scheduling coordination of parallel execution processes. There is no coordination in co-scheduling. In gang scheduling, processes of a given application are scheduled simultaneously, requiring some synchronization mechanism. In this technique all concurrent parallel applications resides in the cluster memory until their completion, generally leading to a huge usage of the virtual memory system. A third approach is to use a special mechanism capable of preempting a parallel execution, removing it from the memory, and scheduling another parallel execution, potentially retrieving its state from a storage system. This approach requires a checkpoint/restart system being able to save the state of a parallel execution and retrieve it. As for processes of a classical operating system, the notion of time slices can be used to apply the quotas and scheduling policies. More time slices or larger time slices are given for execution of higher priority. For this approach only one execution resides in memory at a given time. However, uploading and downloading

executions to and from the memory involves disc operations which can add a significant overhead.

Batch scheduling consists in queuing all jobs to run on a system and launching them in the queue order to fit free resources of the system. Example of existing implementation of batch schedulers are PBS[1], ... The main drawback of the batch scheduling approach is its lack of fairness in case of heterogeneous jobs. An application that needs a large number of nodes but for a short period may have to wait that all longer jobs running on fewer number of nodes end before being allow to run.

Gang scheduling consists in coordinate the execution of a whole single application, other applications are stalled. Thus all resources are used by the application, improving all synchronous operations.

[2], [3] has proposed the first implementation of gang scheduling, called SCore. SCore targets clusters and is based on a *Network preemption* procedure. SCore uses the PM communication library [4]. The gang scheduling itself is performed using UNIX signal mechanism. When an application have to be unscheduled, all it's processes on all nodes receive a SIGSTOP signal. Thus when another application is scheduled by the reception of the SIGCONT signal, it gets exclusive usage of computational power and network resources. However, the memory is shared between running and stopped applications. Memory sharing is resolved by the virtual memory mechanism of the operating system, as inactive applications may be transfered in swap memory. The gang scheduling strategy we present explicitly stores and reloads stopped processes, thus avoiding memory sharing and not relying on operating system swapping mechanism. Another difference induced by these approaches is the network flush algorithm. SCore does not use the Chandy&Lamport algorithm to flush the network, but a three phases synchronization algorithm using the PM flow-control protocol. The Chandy&Lamport algorithm we implemented uses only one synchronization phase. An example of gang scheduling evaluation on LLNL Cray T3D can be found in [5].

[6] evaluates different scheduling policies using the LSF batch scheduler. All policies are refinements of gang scheduling techniques allowing each application to run solely on the required processors. Three classes of applications are considered: short (5 minutes termination expected time), medium (60 minutes) and long running (no limit). The paper studies different policies when an application with a shorter termination expected time is queued. They demonstrate that preemption, implemented with checkpointing techniques, is crucial to obtain good response time.

The last approach, namely co-scheduling, consists in launching all applications on the system resources and letting each node operating system scheduling the different jobs. The lack of coordination for the simultaneous execution of all nodes job of an application is a major drawback for synchronous operations, but this approach allows better overlapping communication of a job by computation of an other one. Previous work [7] analyses experimentally that co-

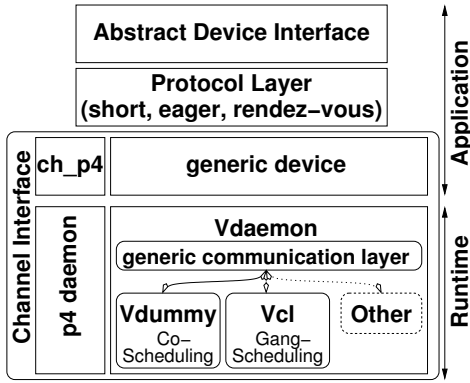


Fig. 1. Architecture of MPICH-V compared to architecture of MPICH-P4

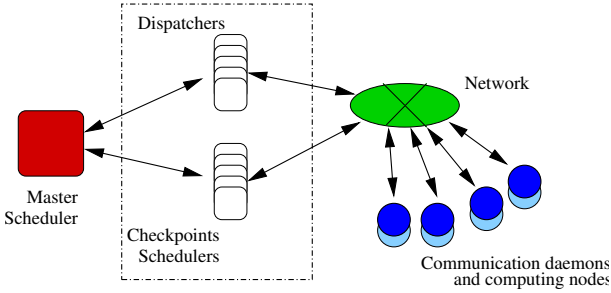


Fig. 2. Typical deployments of MPICH-V with many applications

scheduling outperform gang scheduling for clusters in-core running applications. Our study push this results in the case of out of core running applications, a major concern for Grid systems as the number of users and tasks expands.

[8], [9], [10] considers determined classes of application based on their communication and I/O characteristics. Thus they improve gang scheduling of such known applications by co-scheduling applications of different classes.

As the co-scheduling relies on the operating system memory scheduling, running out of core applications leads to high overhead due to the system swap policy. [11] introduces different paging techniques to reduce the number of page faults. We can also use other approaches

III. COMMON FRAMEWORK

We plan to compare two scheduling methods for MPI application time sharing. The first one is based on a checkpoint capable MPI implementation. At the end of a time slice, a set of applications is checkpointed on disk, and another set of applications is loaded from a previous checkpoint and run during the next time slice. The second one is based on a non checkpoint capable MPI implementation, running all applications simultaneously on the same set of nodes. This method relies on the operating system scheduler to perform time sharing and is called co-scheduling.

For our experiment, we use the MPICH-V framework. This framework allows to implement different kind of fault tolerance protocols for MPI. Deriving from this main purpose, we

developed two non fault tolerant protocols. One including distributed checkpointing facility based on the Chandy&Lamport algorithm (Vcl) [12], [13], and another one implementing basic MPI communication, without checkpoint capabilities (Vdummy). Using such an implementation is mandatory to perform a fair comparison: as the two implementations share the same framework, any performance difference is related to the scheduling itself, and not to implementation optimizations. Moreover, we compared the MPICH-V framework to the reference implementation MPICH-P4 in [13]. Figure 3 recalls these performances comparison between the MPICH-V framework and the reference implementation, and validates the checkpoint enabled MPICH-Vcl version performance compared to standard MPICH-Vdummy version.

MPICH-V is based on the MPICH library [14], which builds a full MPI library from a channel. A channel implements the basic communication routines for a specific hardware or for new communication protocols. MPICH-V consists in a set of runtime components and a channel (ch.v) for the MPICH library.

The different protocols are implemented in the MPICH-V framework at the same level of the software hierarchy, between a MPI high level protocol management layer (managing global operations, point to point protocols, etc.) and the low level network transport layer. Among the other benefits, this allows to keep unmodified the MPICH implementation of point to point and global operations, as well as complex concepts such as topologies and communication contexts. A potential drawback of this approach might be the necessity to implement a specific driver for all types of Network Interface (NIC). However, several NIC vendors provide low level, high performance (zero copy) generic socket interfaces such as Socket-GM for Myrinet, SCI-Socket for SCI and IPoIB for Infiniband. MPICH-V protocols typically seat on top of these low level drivers. So this is one of the most relevant layer for implementing fault tolerance if criteria such as design simplicity, high performance, heterogeneous network migration and portability are to be considered.

A. Dispatcher

The dispatcher of the MPICH-V environment has two main purposes: 1) to launch the whole runtime environment (encompassing the computing nodes and the auxiliary "special" nodes) on the pool of machines used for the execution, and 2) to monitor this execution, by detecting node disconnection and then stop the execution.

The Dispatcher is in charge of a single MPI application. If more than one application is running at a time on a cluster, each is controlled by it's own Dispatcher.

For all protocols, the launching phase is always divided in two parts: 1) a shell script sets up the execution, and 2) a C executable launches the different components involved in the MPI execution and monitors them.

The shell script creates configuration files, containing the paths to the different executables, the addresses of the machines used for the execution, as well as the ports they will

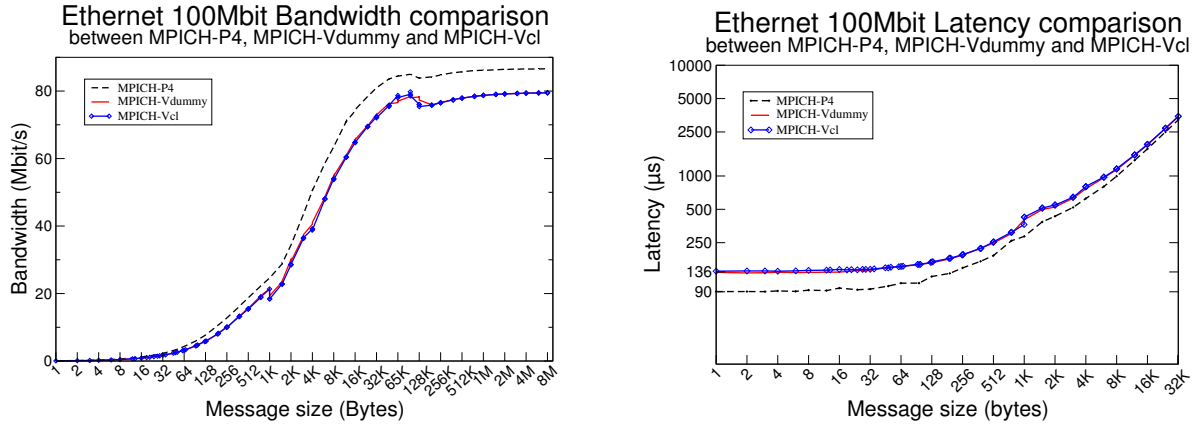


Fig. 3. Bandwidth and latency comparison between MPICH-P4, MPICH-Vdummy, and MPICH-Vcl on Fast-Ethernet network.

use to communicate. These files are read by the C program which launches the execution accordingly. A typical execution then consists in launching the auxiliary programs (checkpoint scheduler, checkpoint servers, depending on the protocol) and waiting for their connection to the dispatcher. When all auxiliaries have successfully connected, the computing nodes are then launched through their communication daemon. They connect to the dispatcher in order to establish a socket connection used to detect a crash or to send internal messages (such as the finalize message, stop, or continue scheduling messages).

During the execution of the program, a *select()* loop is used and waits for any node event. A node event can either be a disconnection or a finalize message. A pipe file, called the *dispatcher pipe* is used in order to receive "orders" from the user. This can be used to move an execution from a set of nodes to another one by simply launching the corresponding command in that pipe. Then the dispatcher kills the execution on the moving nodes and restarts it on the new nodes.

When the finalization is done, the dispatcher cleans the used nodes by removing all the auxiliary programs and the temporary files that have been created.

B. Driver

The driver is the part of the MPICH-V framework linked with the MPI application. It implements the Channel Interface of MPICH. Our implementation only provides synchronous functions (bsend, breceive, probe, initialize and finalize), as the asynchronism is delayed to another component of the architecture.

C. Communication daemon

The core of the communication daemon is a select loop: it manages one socket for every computing node and one socket for every specific components. Every send or receive operation is asynchronous. Thus, a communication is not blocked by another slower one. At the contrary, the communication across the Inter Process Communication mechanism to the MPI process is synchronous and its granularity is the whole

protocol message. The communication daemon is in charge of all distributed checkpoint mechanisms.

The checkpoint of the daemon uses an explicit serialization of its data and when a checkpoint is requested, some messages have to be logged on the receiver (considered as the in-transit messages of the Chandy-Lamport algorithm).

a) *The generic Communication daemon:* Daemons implements a generic communication layer to provide all the communication routines between the different kind of components involved in the MPICH-V architecture, independently of the protocols. Checkpoint enabled protocols are designed through the implementation of a set of hooks called in relevant routines of the generic layer and some specific components (figure 1).

The daemon handles the effective communications, namely sending, receiving, reordering messages, establishing connections with all components of the system and detects failures. In each of these routines, protocol dependent functions are called. The collection of all these functions is defined through a fault tolerance API and each protocol implements this API. In order to reduce the number of system calls, communications are packed using *iovec* related techniques by the generic communication layer. The different communication channels are multiplexed using a single thread and the *select()* system call. This common implementation of communications eases the implementation of protocols and allows a fair comparison between them.

D. Checkpoint scheduler

The checkpoint scheduler requests computation processes to checkpoint according to a given policy. The policy is protocol dependent. However, the checkpoint scheduler includes requesting information to computing nodes, useful to implement a smart policy for message logging protocols. It also implements a global checkpoint sequence number, usable in a coordinated checkpoint policy to ensure success of a global checkpoint. In this article, the checkpoint scheduler uses a coordinated checkpoint policy driven by the Master

Dispatcher, intended to checkpoint applications before it is stopped at the end of a time slice.

E. Master Scheduler

The Master Scheduler is a new component introduced in the MPICH-V architecture in order to target Grids. The master scheduler coordinates the scheduling of all the applications on the system. Given a list of applications to schedule, it first launches the dispatcher and checkpoint scheduler component of each application. The application deployment itself is performed by the Dispatcher of this application.

The number n of job to run simultaneously is given as parameter of the master scheduler. The Master Scheduler associates a time slice to each application. When a scheduled application has been running long enough to expire its time slice, the Master Scheduler requests the Dispatcher and the Checkpoint Scheduler to stop this application. When it is stopped, the Master Scheduler requests another Dispatcher to restart the associated application, and the time slice for this application begins.

The Master Dispatcher implements various policies for stopping/restarting applications. It is possible to Co-schedule k applications of a set of n . In addition to the usual SIGSTOP/SIGRESTART method, three checkpoint/restart overlap policies are implemented. These four policies are detailed in section IV-B.2.

IV. THREE TECHNIQUES OF SCHEDULING

We study three different kinds of scheduling for sharing multiple MPI application on a single set of nodes. Two of them, co-scheduling and gang scheduling, were compared in the context of clusters and using applications with low memory usage in [7]. These two techniques were not studied in the context of a large memory usage, leading to out-of-core computation when many applications runs simultaneously. We propose a new hybrid method based on the co-scheduling and the gang scheduling. The main idea is to limit the number of concurrent job co-scheduled using gang scheduling global time slice and checkpoint related techniques, so that physical memory would not be exhausted by co-scheduling all applications.

A. Literature techniques of time-sharing

1) *Co-Scheduling*: The first one, called co-scheduling consists in an uncoordinated approach. The processes of each application are launched simultaneously on all nodes. On each node, the local operating system scheduler is in charge of sharing resources between the processes of different applications. Thus processes of an application may not be scheduled at the same time on all nodes, which could impact performances of applications using a tightly synchronized communication scheme. Moreover, the frequent context switches between processes may introduce many cache faults. Nonetheless, This technique requires no specific implementation, and computational and network resources may be better used, like in the multi-threaded programming scheme.

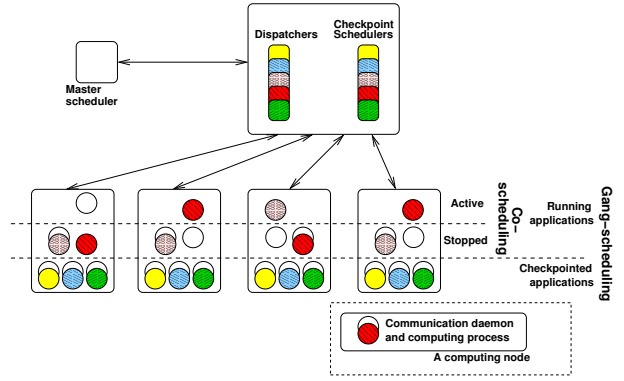


Fig. 4. Typical deployment of hybrid scheduling

To perform comparison of this technique to others, we used the MPICH-V environment. The MPICH-V dummy implementation does not include any checkpoint or time-sharing ability, and is a good candidate to perform a fair comparison with the other scheduling techniques we study. The Master Scheduler uses a policy which just spawn all applications on the nodes simultaneously, and then waits for termination, relying on local node's operating system to schedule processes of all applications.

2) *Gang-scheduling*: The second approach is called gang scheduling. It consists in synchronizing all local scheduler so that all processes of a single distributed application are scheduled simultaneously, while all other applications are stopped and sleeping. All resources are thus employed on the execution of a single application, avoiding the wait for messages from a process not scheduled on another node. However, no multi-thread effect is possible as only one application is running at a time. As discussed in [7], on many applications, the co-scheduling technique outperforms gang scheduling, as applications do not perfectly overlap communications with computation.

Gang scheduling is implemented in the Master Scheduler of the MPICH-V framework.

B. New hybrid technique of time-sharing for high memory requirement

In this method, we propose to use co-scheduling for in-core computation, as it has been proven to be more efficient than gang scheduling. When out-of-core computation would appear using co-scheduling, we use a gang scheduling related technique to enforce that only a subset of the applications is running on the nodes. As the overhead induced by the applications switch is related to time to store process memory on local disk, it is much higher than node operating system context switch overhead. As a consequence, time slices have to be much longer. Thus, gang scheduling is mandatory to reach good performance for any communicating application. Waiting for a message during the full time slice of an application would lead to very poor network performance.

The figure 4 presents a typical deployment of the hybrid architecture. In this example we suppose that physical memory

size allow to run 2 simultaneous co-scheduled applications without inducing out-of-core computation.

The Master Scheduler controls Dispatchers and Checkpoint Schedulers of each application, enforcing some applications to be stopped and swapped out of memory, while some others are restarted. We explain in section IV-B.2 different methods to swap out of memory an application. For the set of applications that are running, their processes are co-scheduled by each node operating system.

1) *Network management of application switching*: To stop an application, network status have to be saved. We use the Chandy&Lamport algorithm to flush network before the application is stopped. The Checkpoint Scheduler and the communication daemons are dedicated to save the network state. The checkpoint scheduler requests every communication daemon to take a global snapshot by sending a tag in each communication channel. On the reception of this tag, a daemon stops the computing process, and then sends the tag in every communication channel. When the checkpoint scheduler have received a tag from every communication daemon, the network flush is complete.

2) *Memory management of application switching*:

a) *SIGSTOP/SIGCONT policy, system memory management*: In previous implementations of gang scheduling, each process of an application are stopped by SIGSTOP UNIX mechanism. In the case of large memory consumption, memory sharing is managed by operating system swap policy. As running applications were swapped on disk during the previous time slice, the memory pages are reloaded on demand during the execution (thus swapping out some pages used by stopped applications). The overhead using this technique is very unpredictable as it relies on operating system swapping policy. The number of page faults have a major impact on overall performance. We plan to measure the number of page faults, and relying on SIGSTOP/SIGCONT mechanism would introduce perturbations in our measurements. Thus, we prefer to use a checkpoint based technique, which overhead is well bounded.

b) *Checkpoint policies, explicit memory management*: In our implementation, we use checkpoint/restart to explicitly manage memory swapping of applications memory. When a computing process is requested to be stopped, it performs a checkpoint to local disk, thus freeing all memory it uses. Complete memory of process to be scheduled in the next time slice is reloaded from checkpoint, thus no page fault occurs during the time slice.

Many policies may be used to overlap checkpoint, restart, and computations. Figure 5 presents the three techniques implemented in the comparison framework. The first method (called sequential checkpoint/restart) does not overlap checkpoint and restart. It avoids to load the two applications simultaneously in memory at the expense of serializing checkpoint, restart and computation during the context switch.

The second technique (called overlapped checkpoint/restart) overlaps checkpoint and restart, trying to reduce context switch cost. It requires more memory as one application is reloaded

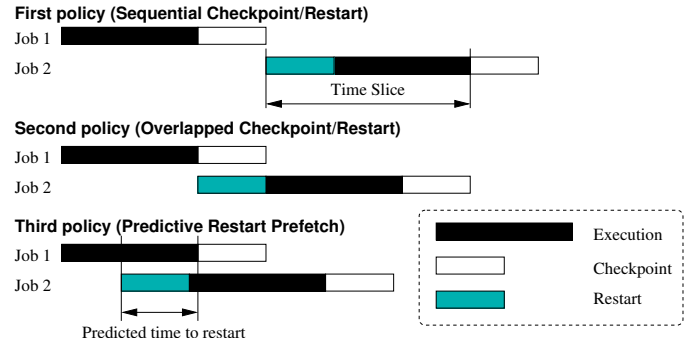


Fig. 5. Policy for checkpointing a job and restarting the next scheduled one on the same node

before the previous is fully flushed out of memory, and induces simultaneous disk accesses.

The third technique (called predictive restart prefetch) tries to prefetch restart during the end of the time slice of the previous application, so that restart is overlapped by computation of the application to be stopped, and checkpoint is overlapped by computation of the next application. It has to rely on an oracle predicting time to restart and induces simultaneous memory occupation during the end of the time slice of each application.

These three techniques are compared in section V-B.

V. PERFORMANCE EVALUATION

A. Experimental conditions

We present a set of experiments in order to evaluate the different components of the system.

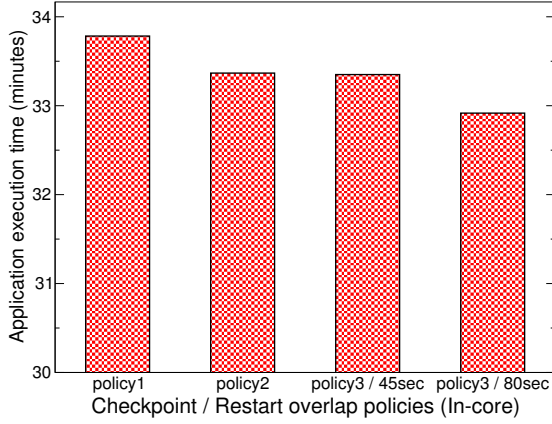
Experiments are run on a 32-nodes cluster. Each node is equipped with an AthlonXP 2800+ processor, running at 2GHz, 1GB of main memory (DDR SDRAM), and a 70GB IDE ATA100 hard drive and a 100Mbit/s Ethernet Network Interface card. Swap space is set to 10GB. All nodes are connected by a single Fast Ethernet Switch.

All these nodes use Linux 2.4.21 as operating system. The tests and benchmarks are compiled with GCC (with flag -O3) and the PGI Fortran77 compilers. All tests are run in dedicated mode. Each measurement is repeated 5 times and we present a mean of them.

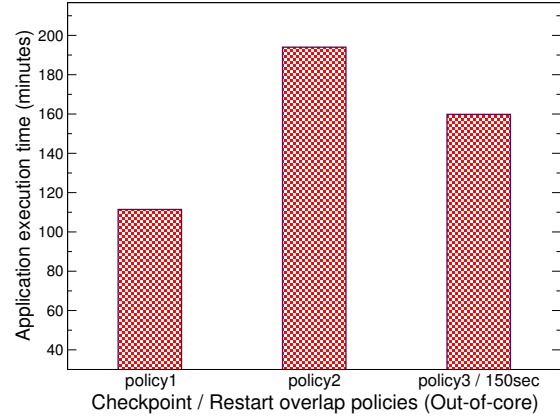
The first experiments are synthetic benchmarks analyzing the individual performance of the subcomponents. We use the NetPIPE [15] utility to measure bandwidth and latency of the MPICH-V framework. This is a ping pong test for several message sizes and small perturbations around these sizes. The second set of experiments is the set of kernels and applications of the NAS Parallel Benchmark suite [16], written by the NASA NAS research center to test high performance parallel computers.

B. Checkpoint/restart overlap scheduling policy

Figure 7 presents NAS benchmark BT class C on 25 nodes performance, for the three Checkpoint/Restart application context switch policies presented in section IV-B.2.b.



(a) Memory occupation fits in physical memory (in-core).



(b) Memory occupation does not fit un physical memory (out-of-core).

Fig. 7. Checkpoint / Restart application context switch policies comparison, using NAS benchmark BT class C on 25 nodes on Ethernet performance criteria.

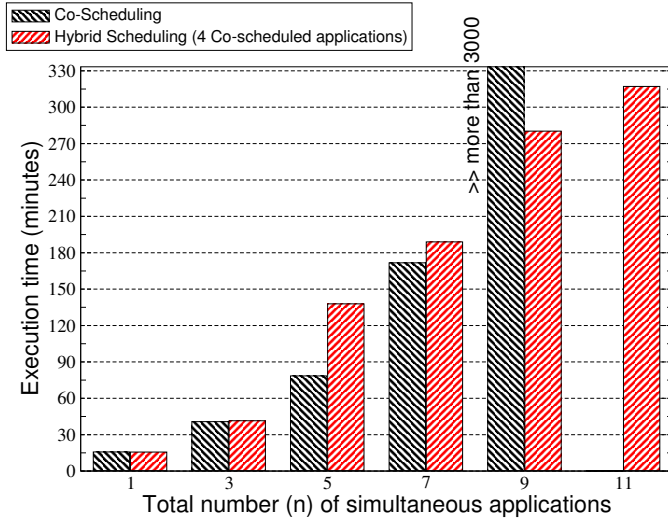


Fig. 6. Computation time for running n BT class C on 25 nodes

VI. CONCLUSION

REFERENCES

- [1] R. L. Henderson, "Job scheduling under the portable batch system," in *Workshop on Job Scheduling Strategies for Parallel Processing*, 1995, pp. 279–294.
- [2] A. Hori, H. Tezuka, Y. Ishikawa, N. Soda, H. Konaka, and M. Maeda, "Implementation of gang-scheduling on workstation cluster," in *Proceeding of IPPS'96 Workshop on Job Scheduling Strategies for Parallel Processing*, D. G. Feitelson and L. Rudolph, Eds. Springer-Verlag, April 1996, pp. 126–139.
- [3] A. Hori, H. Tezuka, and Y. Ishikawa, "Overhead analysis of preemptive gang scheduling," *Lecture Notes in Computer Science*, vol. 1459, pp. 217–230, April 1998.
- [4] H. Tezuka, A. Hori, Y. Ishikawa, and M. Sato, "PM: An operating system coordinated high performance communication library," in *Proceedings of the international conference and Exhibition on high-Performance Computing and Networking*. Springer-Verlag, April 1997, pp. 708–717.
- [5] D. G. Feitelson and M. A. Jette, "Improved utilization and responsiveness with gang scheduling," in *Job Scheduling Strategies for Parallel Processing*, D. G. Feitelson and L. Rudolph, Eds., vol. 1291. Springer Verlag, 1997, pp. 238–261.
- [6] E. W. Parsons and K. C. Sevcik, "Implementing multiprocessor scheduling disciplines," in *Proceeding of IPPS'97 Workshop on Job Scheduling Strategies for Parallel Processing*, D. G. Feitelson and L. Rudolph, Eds. Springer Verlag, April 1997, pp. 166–192.
- [7] P. Strazdins and J. Uhlmann, "Local scheduling out-performs gang scheduling on a beowulf cluster," Department of Computer Science, Australian National University, Tech. Rep., January 2004.
- [8] Y. Wiseman and D. G. Feitelson, "Paired gang scheduling," in *IEEE Transactions on Parallel and Distributed Systems*, vol. 14, June 2003, pp. 581–592.
- [9] W. Lee, M. Frank, V. Lee, K. Mackenzie, and L. Rudolph, "Implications of I/O for gang scheduled workloads," *Proceedings of the Job Scheduling Strategies for Parallel Processing*, pp. 215–237, 1997.
- [10] F. A. B. da Silva and I. D. Scherson, "Improving parallel job scheduling using runtime measurements," *Proceedings of the 6th Workshop on Job Scheduling Strategies for Parallel Processing*, pp. 18–38, May 2000.
- [11] K. D. Ryu, N. Pachapurkar, and L. L. Fong, "Adaptive memory paging for efficient gang scheduling of parallel applications," in *18th International Parallel and Distributed Processing Symposium IPDPS'04*, April 2004.
- [12] A. Bouteiller, P. Lemarinier, G. Krawezik, and F. Cappello, "Coordinated checkpoint versus message log for fault tolerant MPI," in *IEEE International Conference on Cluster Computing (Cluster 2003)*. IEEE CS Press, december 2003.
- [13] P. Lemarinier, A. Bouteiller, T. Herault, G. Krawezik, and F. Cappello, "Improved message logging versus improved coordinated checkpointing for fault tolerant MPI," in *to appear in IEEE International Conference on Cluster Computing (Cluster 2004)*, 2004.
- [14] W. Gropp, E. Lusk, N. Doss, and A. Skjellum, "High-performance, portable implementation of the MPI message passing interface standard," *Parallel Computing*, vol. 22, no. 6, pp. 789–828, September 1996.
- [15] Q. Snell, A. Mikler, and J. Gustafson, "Netpipe: A network protocol independent performance evaluator," in *IASTED International Conference on Intelligent Information Management and Systems*, June 1996.
- [16] D. Bailey, T. Harris, W. Saphir, R. V. D. Wijngaart, A. Woo, and M. Yarrow, "The NAS Parallel Benchmarks 2.0," Numerical Aerodynamic Simulation Facility, NASA Ames Research Center, Report NAS-95-020, 1995.