

Performance Analysis of Mapping Strategies for Parallel Simulation

Yong Meng TEO and Seng Chuan TAY

Department of Information Systems & Computer Science

National University of Singapore

Lower Kent Ridge Road

Singapore 119260

email: {teoym,taysengc}@iscs.nus.sg

Abstract

This paper presents an analytical model to evaluate the performance of conservative parallel simulation using a finite-buffered multistage interconnection network as an example. The model is formalized based on two important time components in parallel and distributed processing: computation time and communication time. Our mathematical analysis identifies the major constituents of simulation overheads that affect the performance of parallel simulation. We also show that a perfectly balanced workload distribution may not necessarily translate into better performance. On the contrary, we have shown that a balanced mapping of workload may increase communication overheads resulting in a longer simulation elapsed time.

Keywords : simulation modeling, performance modeling, parallel simulation, multistage interconnection network, workload balancing

1 Introduction

Parallel discrete-event simulation (PDES) technique provides a promising avenue to reducing the runtime of large applications. Two related paradigms, *conservative approach* and *optimistic approach*, have been widely discussed [1, 2, 3]. In brief, both approaches decompose the model to be simulated into loosely coupled components and simulate each component by a process. Interaction among the processes in PDES is performed by passing timestamped messages, which usually represent events scheduled, or to be scheduled if feasible, by one process at another.

Many factors contribute to the performance of a PDES program. Active research areas within this arena include new protocols, mathematical performance analysis, time parallelism, hardware support, load balancing and dynamic memory management [4]. Among them

performance analysis provides an insight to the simulation modeling and identifies the cause of unsatisfactory implementations. The work done in this area is useful to simulationists as it provides the guidance necessary to improve, or even redesign, their programs. This paper focuses on the performance modeling of the conservative parallel simulation. Instead of using self-contrived queuing network models, our mathematical analysis is based on a realistic multistage interconnection network (MIN). The proposed simulation and performance models can easily be applied to study other simulation applications. This paper is organized as follows. Section 2 gives an overview of conservative parallel simulation mechanism for modeling and simulating interconnection networks. Section 3 abstracts the performance model for each component of MIN simulation. Section 4 analyzes the performance of different partitioning schemes for mapping a simulation application onto a finite number of processors. Section 5 validates the derived metrics from the proposed analytical performance model with results collected from the simulation model. We comment in particular the effect of load balancing on simulation performance. Based on the analytical results, section 6 discusses possible improvements in simulation modeling. Lastly, section 7 contains our concluding remarks.

2 Conservative Parallel Simulation Model Overview

A detailed description of the MIN simulation is available in [5]. Briefly, the process-oriented parallel simulation model consists of three types of logical processes: *packet generator (G)*, *switching element (SW)* and *sink (S)* (see figure 1). Based on the statistical distributions for inter-arrival time and packet destination, each *G* process generates access request packets and sends them to a *SW* process on the $(\log_2 n - 1)$ th stage. Each *SW* process serves as a relaying agent to forward the

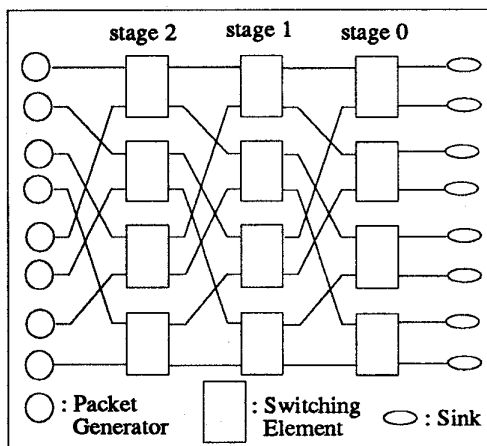


Figure 1: Components of MIN Simulation

packets to their destinations. Two buffers of finite size are embedded in each SW process to model a blocking switch. A constant switch delay is used in each SW process. Each S process serves as a destination for access packets and its service time is assumed to be instantaneous. Each process has its own local clock to indicate its simulation progress. A blocked message due to the unavailability of free space in its receiving buffer will have to wait until the buffer space become available.

2.1 Ensuring the Correctness of Simulation Results

The simulation of a finite-buffered interconnection network is more sophisticated as compared to that of infinite buffered network. For example, the execution of each departure event in G and SW processes should not send the departing packet to their successor unless the receiving buffer contains at least one unoccupied buffer. Simulation processes therefore must adhere to some *hand-shaking protocols* during each packet transmission [5]. The following signals and messages are used in the MIN simulation. The underlined indicates that a signal/message needs to be time-stamped.

- Request Progress (REQ) - This signal is used by a generator to find out if the message buffer of the service station can accommodate additional message.
- Slot Available (AVA) - This signal is replied by a service station to notify that the message buffer of the desired in-coming link is not full.
- Slot Not Available (NAV) - This signal is replied by a service station to notify that the message buffer of the desired in-coming link is full.
- Access Request (ACC) - This message contains the source and destination addresses of an access re-

quest. It is sent to the service station upon receiving the reply signal "Slot Available".

- Transmit Time (TIM) - This message contains the time for a generator to re-transmit a blocked message. It is sent by a service station when a message is removed from its message buffer where a NAV signal was previously replied to the generator on that link.
- Stop (STO) - This signal is used by a service station/sink to inform the generator/service station that its process is going to terminate.

An illustration of the hand-shaking protocols is given in figure 2.

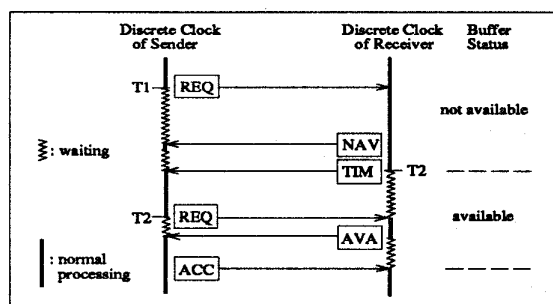


Figure 2: Hand-Shaking Protocols for Packet Transmission

2.2 Resolving Deadlock Problem

In the conservative parallel MIN simulation, the execution of an arrival event in a SW process must be followed by the scheduling of a new arrival event so that causality can be ensured in each iteration. After an arrival event execution, the SW process will therefore have to enter a waiting state until the new packet arrives. If the waiting processes happen in a loop, *deadlock* will occur.

The following *relaying protocols* are proposed to resolve the deadlock problem in the MIN simulation. In general, a process on receiving a null message or requesting a buffer status should also send null messages, with lookahead timestamps if available, to all the other three links. In this way the null messages are circulated among the processes. If a received null message indicator is greater than all other outstanding event times, a process is safe to proceed on with the simulation, thereby resolving a deadlock. In the following descriptions, $toa[i]$ and $tosc[j]$ refer to the time of arrival on incoming link $[i]$, $i \in \{0, 1\}$ and time of service completion on out-going link $[j]$, $j \in \{0, 1\}$, respectively where 0 refers to the event occurred on upper link and 1 lower link. α, β , and γ are the time indicators to inform the sender/receiver on the various links that the process will not send them any useful (hand-shaking)

signals and messages before the specified time. We let $i' = i \oplus_2 1$, and $j' = j \oplus_2 1$, where \oplus_2 is an additive operator of modulo two. The underlined refers to in-coming link, and overlined out-going link.

1. When a REQ signal of timestamp Θ is sent to an out-going link[j], three NULL messages containing the time indicators α , β and γ respectively, where

$$\alpha = \begin{cases} \text{toa}[0] & \text{if } \overline{\text{link}[0]} \text{ full} \\ \Theta & \text{otherwise} \end{cases}$$

$$\beta = \begin{cases} \text{toa}[1] & \text{if } \overline{\text{link}[1]} \text{ not full} \\ \Theta & \text{otherwise} \end{cases}$$

$$\gamma = \begin{cases} \text{tosc}[j'] & \text{if } \overline{\text{link}[j']} \text{ not empty} \\ \Theta + \text{switch delay} & \text{otherwise} \end{cases}$$

are also sent to in-coming link[0], in-coming link[1] and out-going link[j'] respectively.

2. When a NULL message is received while a process is waiting on an out-going link[j] for the reply to a REQ signal, the process should also forward the NULL message on the other three links with the time indicators specified in clause 1.
3. When scheduling a new arrival event on an in-coming link[i], a process should wait for a REQ signal and also process a NULL message if it arrives.
4. If a NULL message is received in clause 3, three NULL messages containing the time indicators α , β and γ respectively, where

$$\alpha = \begin{cases} \text{toa}[i'] & \text{if } \overline{\text{link}[i']} \text{ not full} \\ \min(\text{tosc}[0], \text{tosc}[1]) & \text{otherwise} \end{cases}$$

$$\beta = \begin{cases} \text{tosc}[0] & \text{if } \overline{\text{link}[0]} \text{ not empty} \\ \Theta + \text{switch delay} & \text{otherwise} \end{cases}$$

$$\gamma = \begin{cases} \text{tosc}[1] & \text{if } \overline{\text{link}[1]} \text{ not empty} \\ \Theta + \text{switch delay} & \text{otherwise} \end{cases}$$

are also sent to in-coming link[i'], out-coming link[0] and out-going link[1] respectively.

2.3 Reducing Synchronization Messages

A major drawback of using null messages in parallel simulation is that the messages may cause combinatoric explosion in memory utilization. In the worst case, simulation cannot proceed because the machine runs out of memory space required for storing null messages. A flushing mechanism proposed in [5] is used in our implementation to handle this problem. When a null message is received, the process will flush in all the null messages that have already arrived on that link and select the *largest* time indicator. The forwarding of null messages is only on the largest time indicator and the rest are discarded. Such a flushing mechanism is also able to maintain the liveness of simulation as at least one null message is circulating in the network throughout the simulation duration.

3 Analytical Performance Model

Table 1 contains a list of performance parameters used in our analysis. All logarithmic representations are base 2. Based on problem size, MIN's characteristics, data transmission protocols and simulation mechanism, the proposed model formalizes the workload of each simulation component in terms of their elapsed time.

We assume that the destination of each packet transmission is uniformly distributed. Each data transmission among the simulation processes is modeled by two main components: *buffer access time* and *transit time*. The reception of data is modeled by buffer access time only and transit time is excluded to prevent double accounting. We divide the analysis of simulation performance into two parts: *light traffic* ($\lambda \leq \mu$) and *heavy traffic* ($\lambda > \mu$) conditions. In the speedup analysis and simulation bandwidth analysis, we assume that $p = 2 \times n + \frac{n \times \log n}{2}$, i.e., the number of processors is sufficient to achieve the maximum degree of parallelism for the MIN simulation. While modeling T_p , we assume that the elapsed time of each SW process outweighs those of G and S processes.

3.1 Light Traffic

In this case we let $\lambda_{eff} = \lambda$. It is highlighted that the actual number of packets sent by each generator during simulation may be smaller due to unavailability of buffer space. To simplify the analytical model we assume that the probability of such occurrences is negligible, and we further assume that the hand-shaking protocols used in each packet transmission include only REQ and AVA signals. In other words, when a process sends a REQ signal to its receiver to find out the

parameter	description
t	duration of simulation
λ	arrival rate used in generator
λ_{eff}	effective arrival rate used in each generator
μ	service rate used in switching element
g	total number of packets sent per generator ($\lambda_{eff} \times t$)
n	number of inputs/outputs of MIN, i.e. network size
p	number of processors
T_{ev}	event (arrival or departure) execution time
T_{bf}	buffer (receive or transmit) access time
T_{tx}	packet/message/signal transmission time
T_{ge}	packet generation time
T_{ac}	packet accounting time
T_{total}^G	total elapsed time of packet generator process
T_{total}^{SW}	total elapsed time of switching element process
T_{total}^S	total elapsed time of sink process
T_1	total elapsed time when one processor is used to execute simulation processes
T_p	total elapsed time when p processors are used to execute simulation processes
$S_p(scheme)$	speedup metric when p processors are used for a particular mapping scheme
$BW_p(scheme)$	bandwidth (number of events simulated per unit time) for a particular mapping scheme

Table 1: Performance Model Parameters

buffer status, the reply is always “available” (see figure 2). It follows that the time incurred in each set of hand-shaking protocols is $2 \times T_{bf} + T_{tx}$, and the time incurred in each packet transmission is $T_{bf} + T_{tx}$.

3.1.1 Packet Generator

The elapsed time of each G process is modeled by three components: *packet generation* (T_1^G), *hand-shaking protocols* (T_2^G) and *packet transmission* (T_3^G). Thus,

$$T_1^G = g \times T_{ge} \quad (1)$$

$$T_2^G = g \times (2T_{bf} + T_{tx}) \quad (2)$$

$$T_3^G = g \times (T_{bf} + T_{tx}) \quad (3)$$

By equations (1), (2) and (3), the total elapsed time of each G process becomes

$$T_{total}^G = g \times (T_{ge} + 3T_{bf} + 2T_{tx}) \quad (4)$$

3.1.2 Switching Element

We model the total elapsed time of each SW process by four components: *event execution* (T_1^{SW}), *hand-shaking protocols* (T_2^{SW}), *packet transmission and reception* (T_3^{SW}) and *null-message transmissions* (T_4^{SW}).

As there are two input links to each SW process, the number of received packets is $2 \times g$. Each packet,

corresponding to one arrival event, will generate one departure event and therefore the total number of events to be executed in a SW process is $2 \times 2g$. It follows that

$$T_1^{SW} = 4g \times T_{ev} \quad (5)$$

$$\begin{aligned} T_2^{SW} &= 2g \times (2T_{bf} + T_{tx}) + \\ &\quad 2g \times (2T_{bf} + T_{tx}) \\ &= 4g \times (2T_{bf} + T_{tx}) \end{aligned} \quad (6)$$

For each SW process, $2g$ packets are received and $2g$ packets are transmitted. Therefore

$$\begin{aligned} T_3^{SW} &= 2g \times T_{bf} + 2g \times (T_{bf} + T_{tx}) \\ &= 4g \times T_{bf} + 2g \times T_{tx} \end{aligned} \quad (7)$$

For each transmission of AVA signal and, REQ signal and ACC message on an in-coming link and out-going link respectively, null messages will be sent on the other three links of the SW process to prevent deadlock. We assume that the null messages are then discarded by their receivers. Since there are $2g$ packets to be received and $2g$ packets to be transmitted in each SW process, we have

$$\begin{aligned} T_4^{SW} &= 2g \times 1 \times 3 \times (T_{bf} + T_{tx}) + \\ &\quad 2g \times 2 \times 3 \times (T_{bf} + T_{tx}) \\ &= 18g \times (T_{bf} + T_{tx}) \end{aligned} \quad (8)$$

It is noteworthy that the actual elapsed time in processing the null messages is larger as the receiving *SW* processes will also select the largest time indicator among all the arrived null message and relay it to the other three *SW* processes. To simplify the analytical model, we assume that the workload in relaying the null messages is negligible.

By equations (5), (6), (7) and (8), the total elapsed time of each *SW* process becomes

$$T_{total}^{SW} = 2g \times (2T_{ev} + 15T_{bf} + 12T_{tx}) \quad (9)$$

3.1.3 Sink

The elapsed time of each *S* process is modeled by two components: *packet accounting* (T_1^S) and *hand-shaking protocols* (T_2^S). Since a *S* process is connected to each output link of the MIN, the number of received packets is g . It follows that

$$T_1^S = g \times T_{ac} \quad (10)$$

$$T_2^S = g \times (2T_{bf} + T_{tx}) \quad (11)$$

By equations (10) and (11), the total elapsed time of each *S* process becomes

$$T_{total}^S = g \times (T_{ac} + 2T_{bf} + T_{tx}) \quad (12)$$

3.1.4 Speedup

Speedup for p processors is defined as the execution time for the best serial algorithm in a single processor (T_1) divided by the execution time for the parallel algorithm using p processors (T_p). In our analysis, T_1 is modeled by the elapsed time of MIN simulation incurred by a sequential program. As causality correctness is easily ensured in sequential simulation, the use of transmission protocols is not necessary in the simulation program and is discarded in our derivation of T_1 . This presents a fairer speedup measure than sometime reported in the literature whereby the value of T_1 is measured by running the parallel algorithm on one processor. This later instance certainly gives better speedup. Given that the number of simulated packets transmitted across the MIN is $n \times g$, the elapsed time of the simulation program becomes

$$\begin{aligned} T_1 &= n \times g \times T_{ge} + \\ &\quad \frac{n \log n}{2} \times 2g \times 2T_{ev} + n \times g \times T_{ac} \\ &= ng \times (T_{ge} + 2 \log n \times T_{ev} + T_{ac}) \end{aligned} \quad (13)$$

Assuming that the elapsed time of each *SW* process outweighs those of *G* and *S* processes, we have

$$T_p = T_{total}^{SW} = 2g \times (2T_{ev} + 15T_{bf} + 12T_{tx}) \quad (14)$$

By equations (13) and (14), under light traffic condition the speedup when p processors (where $p = 2 \times n + \frac{n \times \log n}{2}$) are used becomes

$$\begin{aligned} S_p(nil) &= \frac{ng \times (T_{ge} + 2 \log n \times T_{ev} + T_{ac})}{2g \times (2T_{ev} + 15T_{bf} + 12T_{tx})} \\ &\approx \frac{p}{1 + \frac{15T_{bf} + 12T_{tx}}{2T_{ev}}} \end{aligned} \quad (15)$$

3.1.5 Simulation Bandwidth

Simulation bandwidth, a measure of simulation implementation efficiency, is defined as the number of events executed per second. The total number of events to be executed is

$$2 \times 2g \times \frac{n \times \log n}{2} \quad (16)$$

By equations (14) and (16), the simulation bandwidth under light traffic condition becomes

$$BW_p(nil) \approx \frac{p}{T_{ev} + \frac{15}{2}T_{bf} + 6T_{tx}} \quad (17)$$

3.2 Heavy Traffic

When the arrival rate is greater than the service rate of the *SW* processes connected to the *G* processes, the buffer spaces of *SW* processes will be filled up progressively until all slots are occupied. Subsequently, the packet arrival rate is moderated by the *SW* processes to prevent packet loss. When the steady stage is reached $\lambda_{eff} = \mu$. Our worst case analysis assumes that the time interval before the steady stage is negligible. The hand-shaking protocols used by generators for packet transmission include REQ, NAV and AVA signals and TIM message to delay the packet arrival times (see figure 2). It follows that the elapsed times incurred in each set of hand-shaking protocols used by *G* processes for packet transmission and in *SW* process on the $(n-1)$ th stage for packet reception are $5T_{bf} + 2T_{tx}$ and $5T_{bf} + 3T_{tx}$ respectively. The hand-shaking protocols used in the remaining *SW* processes, on the other hand, remain unchanged as the arrival rate to these processes has been moderated by their service rate. In this aspect, we assume that the buffer space from the 0th to $(\log n - 2)$ th stages are available throughout the simulation duration. The derivations of the remaining performance metrics are available in [6] and we summarize the results in tables 2a and 2b due to space limitation.

4 Analysis for Different Mapping Schemes

Our conservative simulation model requires n *G* processes, $\frac{n \log n}{2}$ *SW* processes and n *S* processes to simulate a $n \times n$ MIN. Mapping of such processes onto a

traffic	elapsed time
Light	$2g \times (2T_{ev} + 15T_{bf} + 12T_{tx})$
Heavy	$2g \times (2T_{ev} + 21T_{bf} + 17T_{tx})$

(a) Elapsed Time

traffic	speedup	simulation bandwidth
Light	$\frac{p}{1 + \frac{15T_{bf} + 12T_{tx}}{2T_{ev}}}$	$\frac{p}{T_{ev} + \frac{15}{2}T_{bf} + 6T_{tx}}$
Heavy	$\frac{p}{1 + \frac{21T_{bf} + 17T_{tx}}{2T_{ev}}}$	$\frac{p}{T_{ev} + \frac{21}{2}T_{bf} + \frac{17}{2}T_{tx}}$

(b) Asymptotic Speedup and Simulation Bandwidth

Table 2: Analytical Results for $p = 2 \times n + \frac{n \times \log n}{2}$

particular processors configuration for efficient implementation is a NP-complete problem [7]. In this paper, we introduce and analyze the performance of *three mapping schemes*, namely, *horizontal partitioning scheme (HPS)*, *vertical partitioning scheme (VPS)* and *modular partitioning scheme (MPS)*. Details of the formulation and mathematical transformation used can be found in [5]. Figures 3a, 3b and 3c show the three schemes respectively. A rectangular block denotes the mapping of a number of simulated processes onto *one* processor. In all mapping schemes each *G* process is placed together with the *SW* process on its out-going link in a same partition, and each *S* process is placed together with the *SW* process on its in-coming link. It is noteworthy that in MPS each module containing *four* *SW* processes where each process corresponding to a 2×2 switch, is executed/mapped on the same processor. However, in VPS each partition containing $\log n$ *SW* processes is executed on a processor. Therefore, in the following speedup and simulation bandwidth comparison analysis we let $n = 2^i$ for some integer i . For ease of discussion, we further assume that $p = \frac{n}{2}$ and the transit time for intra-processor communication is negligible. Due to space constraint, we assume light traffic condition in the following elapsed time analysis, speedup analysis and simulation bandwidth analysis. The performance metrics for heavy traffic condition can be derived using the same steps.

4.1 Horizontal Partitioning Scheme

In HPS, the *G*, *SW* and *S* simulated processes on one horizontal rectangular box (see figure 3a) are mapped onto one processor. Parallel simulation is terminated only after the execution of *S* processes in all processors are completed. Therefore, T_p is modeled by the maximum workload among processors. Before proceeding further, we re-compute the elapsed times of the affected processes due to this scheme.

4.1.1 Packet Generator

As the intra-processor transit time between the generator process and the *SW* process on the $(\log n - 1)$ th stage is negligible, the total elapsed time for each *G* process (see equation (4)) becomes

$$\begin{aligned} T_{total}^G &= g \times (T_{ge} + 3T_{bf} + 2 \times 0) \\ &= g \times (T_{ge} + 3T_{bf}) \end{aligned} \quad (18)$$

4.1.2 SW Process on the $(\log n - 1)$ th Stage

While modeling the maximum workload among processors, we do not consider the top and the bottom partitions due to the diminished inter-processor links which reduce the total elapsed time. The elapsed times due to event execution (equation (5)) and packet transmission and reception (equation (7)) remain unchanged. The elapsed time due to hand-shaking protocols (equation (6)) for the *SW* processes on the $(\log n - 1)$ th stage is reduced to

$$\begin{aligned} T_2^{SW} &= 2g \times (2T_{bf} + 0) + \\ &\quad 2g \times (2T_{bf} + T_{tx}) \\ &= 2g \times (4T_{bf} + T_{tx}) \end{aligned} \quad (19)$$

For the same reason, the elapsed time due to null-message transmission (equation (8)) is reduced to

$$\begin{aligned} T_4^{SW} &= 2g \times 1 \times 3 \times (T_{bf} + 0) + \\ &\quad 2g \times 2 \times 3 \times (T_{bf} + T_{tx}) \\ &= 6g \times (3T_{bf} + 2T_{tx}) \end{aligned} \quad (20)$$

By equations (5), (7), (19) and (20), the total elapsed time of each *SW* process on the $(\log n - 1)$ th stage becomes

$$T_{total}^{SW} = 2g \times (2T_{ev} + 15T_{bf} + 8T_{tx}) \quad (21)$$

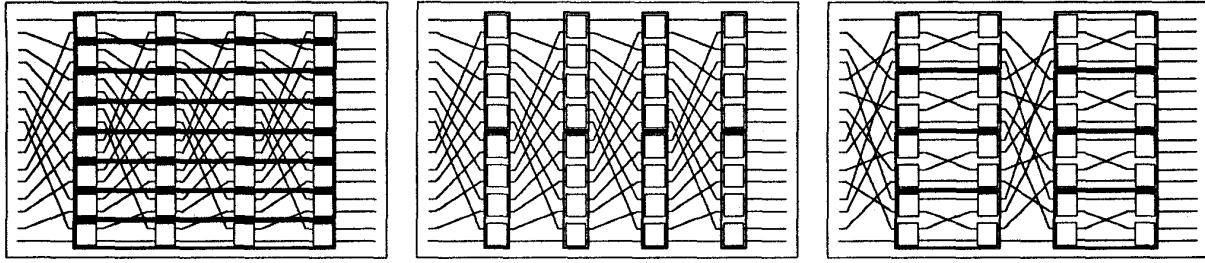
4.1.3 SW Process on the 0th Stage

The elapsed time due to event execution (equation (5)) remain unchanged. The elapsed time due to hand-shaking protocols (equation (6)) is reduced to

$$\begin{aligned} T_2^{SW} &= 2g \times (2T_{bf} + T_{tx}) + 2g \times (2T_{bf} + 0) \\ &= 2g \times (4T_{bf} + T_{tx}) \end{aligned} \quad (22)$$

The elapsed time due to packet transmission and reception (equation (7)) becomes

$$\begin{aligned} T_3^{SW} &= 2g \times T_{bf} + 2g \times (T_{bf} + 0) \\ &= 4g \times T_{bf} \end{aligned} \quad (23)$$



(a) Horizontal Partition

(b) Vertical Partition

(c) Modular Partition

Figure 3: Mapping Schemes for a Finite Number of Workstations

The elapsed time due to null-message transmission for the SW processes (equation (8)) is reduced to

$$\begin{aligned} T_4^{SW} &= 2g \times 1 \times (3T_{bf} + T_{tx}) + \\ &\quad 2g \times 2 \times (3T_{bf} + 2T_{tx}) \\ &= 2g \times (9T_{bf} + 5T_{tx}) \end{aligned} \quad (24)$$

By equations (5), (22), (23) and (24), the total elapsed time of each SW process on the 0 th stage becomes

$$T_{total}^{SW} = 2g \times (2T_{ev} + 15T_{bf} + 6T_{tx}) \quad (25)$$

4.1.4 Sink Process

The elapsed time of each S process (see equation (12)) is reduced to

$$T_{total}^S = g \times (T_{ac} + 2T_{bf}) \quad (26)$$

By multiplying the respective numbers of processes to equations (9), (18), (21), (25) and (26), we have

$$\begin{aligned} T_p &= 2 \times g \times (T_{ge} + 3T_{bf}) + \\ &\quad 1 \times 2g \times (2T_{ev} + 15T_{bf} + 8T_{tx}) + \\ &\quad (\log n - 2) \times 2g \times (2T_{ev} + 15T_{bf} + 12T_{tx}) + \\ &\quad 1 \times 2g \times (2T_{ev} + 15T_{bf} + 6T_{tx}) + \\ &\quad 2 \times g \times (T_{ac} + 2T_{bf}) \\ &= 2g \times (T_{ge} + 2 \log n \times T_{ev} + T_{ac} + \\ &\quad (5 + 15 \log n) \times T_{bf} + \\ &\quad 2(6 \log n - 5) \times T_{tx}) \end{aligned} \quad (27)$$

4.1.5 Speedup and Bandwidth

By equations (13) and (27), the speedup when p processors (where $p = \frac{n}{2}$) are used becomes

$$S_p(HPS) \approx \frac{p}{1 + \frac{15T_{bf} + 12T_{tx}}{2T_{ev}}} \quad (28)$$

By equations (16) and (27), the simulation bandwidth of HPS becomes

$$BW_p(HPS) \approx \frac{2p}{2T_{ev} + 15T_{bf} + 12T_{tx}} \quad (29)$$

4.2 Vertical Partitioning Scheme

In VPS, each processor executes $\log n$ SW processes belonging to the same MIN stage (see figure 3b). During simulation, ACC packets are created by the generators, relayed by SW processes in each stage, and absorbed by the sinks. As the mapping is performed on a stage basis, parallel simulation is completed only when the last (rightmost) processor terminates. T_p is therefore modeled by two components: *elapsed time of the first (leftmost) processor*, and the *traveling time for the last transmission to reach the sink process*. The analytical results are summarized in tables 3 and 4.

4.3 Modular Partitioning Scheme

The MPS is based on partitions of *four* SW processes spread over two MIN stages. A proof of the transformation from the Omega MIN to its equivalent modular topology can be found in [5]. Each processor executes $\log n$ SW processes arranged in $\frac{\log n}{4}$ blocks. For the same reason that sink processes are contained in the rightmost module only, we also model T_p by the same two components as in VPS.

A summary of the analytical results is also given in tables 3 and 4.

5 Model Validation and Performance Analysis

We implemented the conservative parallel MIN simulation model in C language on a network of workstations which mimics a distributed-memory parallel computer.

scheme	elapsed time
HPS	$2g \times (T_{ge} + 2 \log n \times T_{ev} + T_{ac} + (5 + 15 \log n) \times T_{bf} + 2(6 \log n - 5) \times T_{tx})$
VPS	$2g \log_2 n \times T_{ge} + 2(\log n \times (2g + 1) - 1)T_{ev} + T_{ac} + (3 \log n \times (12g + 5) - 13)T_{bf} + 2(2 \log n \times (4g + 3) - 9)T_{tx}$
MPS	$g \log n \times T_{ge} + 2((2g + 1) \log n - 2) \times T_{ev} + T_{ac} + (3(11g + 5) \log n - 28) \times T_{bf} + ((8g + 7) \log n - 22) \times T_{tx}$

 Table 3: Elapsed Time for $p = \frac{n}{2}$

scheme	speedup	simulation bandwidth
HPS	$\frac{p}{1 + \frac{15T_{bf} + 12T_{tx}}{2T_{ev}}}$	$\frac{2p}{2T_{ev} + 15T_{bf} + 12T_{tx}}$
VPS	$\frac{p}{1 + \frac{T_{ge} + 18T_{bf} + 8T_{tx}}{2T_{ev}}}$	$\frac{2p}{T_{ge} + 2T_{ev} + 18T_{bf} + 8T_{tx}}$
MPS	$\frac{p}{1 + \frac{T_{ge} + 33T_{bf} + 4T_{tx}}{2T_{ev}}}$	$\frac{2p}{T_{ge} + 2T_{ev} + 33T_{bf} + 4T_{tx}}$

 Table 4: Asymptotic Speedup and Simulation Bandwidth for $p = \frac{n}{2}$

PVM software [8] was used for spawning the simulation processes and for handling message passing. The simulation model results presented are based on simulating a 16×16 Omega MIN on a network of *eight* homogeneous SUN workstations connected using a 10Mbps Ethernet network. Altogether 64 PVM processes are spawned on the processors.

Five parameter values used in the analytical model are obtained by taking measurements on the workstation network (see table 5). Additional statistical mea-

parameter	time (μsec)
T_{ge}	926
T_{ev}	2990
T_{ac}	479
T_{bf}	3238
T_{tx}	1821

Table 5: Performance Parameter Measurement

asures used in the analysis are defined. Let s_i be a simulation process and $|s_i|$ denotes its workload. In the following analysis, the workload for each simulation process is represented by the size of its object codes, and normalized to that of the S process. For simplicity, we let $|S| = 1$. Thus, the workloads for $|G|$ and $|SW|$ are 1.93 and 6.24 respectively. The normalized coefficient of T_{tx} is denoted by \mathcal{T} . The percentage deviation between the measured and the predicted values δ is computed as $\frac{\text{measured} - \text{predicted}}{\text{predicted}} \times 100$. The average of δ over six different sets of simulation parameters is denoted by $\bar{\delta}$.

Let M be a set containing all s_i required in a MIN simulation. Let \hat{P} be the processor that incurs the largest workload among all processors. We define the *workload distribution coefficient* (ω) as $\omega = \frac{\sum_{s_i \in M} |s_i|}{\sum_{s_j \in \hat{P}} |s_j|}$. The load balancing factor (η) is then define as $\eta = \frac{p}{\omega}$. Thus, the larger the value of η , the more unbalance is the workload distribution. The smallest (best) value of η equals to 1, and corresponds to a perfectly balanced distribution of workload.

We first compare the elapsed times among the different schemes, and between the simulation model (*measured*) and the analytical model (*predicted*). For VPS and MPS, the latter outperforms in terms of elapsed time as shown in table 6. This phenomenon is due to the poor workload distribution and larger inter-processors communication overheads of VPS ($\mathcal{T} = 2$) compared to those of MPS ($\mathcal{T} = 1$). The comparisons of VPS with HPS is more interesting. Although the interprocessor communication overheads of VPS is only half of that of HPS, it seems that the imbalance workload factor is more dominant, making the elapsed time of VPS ($\eta = 1.31$) longer than that of HPS ($\eta = 1$). In this aspect, our explanation is that for VPS the worst number of processes in the processors is larger than that of HPS. Therefore, the elapsed time of VPS consists of a larger number of context switching time incurred by the task scheduler of the UNIX operating system. The comparison of HPS ($\eta = 1$) and MPS ($\eta = 1.06$) is rather counter-intuitive. We observe that the HPS, a perfect load balancing scheme, incurs a significantly larger elapsed time as compared to that of a less balance scheme. Such an observation *contradicts* to the common belief, that the workload distribution among processors must be balanced in order to improve elapsed time. We note here that other factors such as inter-processors communication overheads incurred by a mapping scheme should also be taken into account while balancing the workload distribution. As the inter-processors communication overheads incurred by HPS is larger than that of MPS, its elapsed time is therefore aggravated despite the balanced workload distribution.

The average percentage deviations $\bar{\delta}$ show that accu-

scheme		no. of messages per generator						$\bar{\delta}$ (%)	\mathcal{T}	η
		120	240	360	480	600	720			
HPS	measured	76.79	151.15	226.49	304.09	377.68	454.93	3.66	4	1.00
	predicted	73.19	146.38	219.60	292.76	365.95	439.14	-	-	-
VPS	measured	79.61	160.14	239.93	322.23	403.43	482.46	4.78	2	1.31
	predicted	76.75	153.28	229.81	306.33	382.86	459.39	-	-	-
MPS	measured	67.29	134.00	200.25	269.73	337.18	406.74	4.31	1	1.06
	predicted	64.58	129.03	193.48	257.93	322.38	386.83	-	-	-

Table 6: Comparison of Elapsed Times (in second)

racy of the analytical model for predicting elapsed time is good for all the three mapping schemes. However, the predicted timings show slightly better performance than the measured simulation model timings because the context switching time incurred by the scheduler of the UNIX Operating System is not accounted in our analytical model.

Tables 7 and 8 show that our performance model closely predicts the parallel simulation speedups and simulation bandwidths. The predicted values for all the three mapping schemes are slightly better than the measured values for the same reason. The speedup efficiency is unacceptably low, i.e. around 10%. Our analysis reveals that this is due to the higher than expected cost of PVM communication. Better speedup can be obtained on parallel machine where communication is much more efficient and less costly. In the following section, we present a breakdown of various communication costs, and analyze the effect of reducing these overheads on simulation performance.

In summary, the measured and predicted values for elapsed time, speedup and simulation bandwidth show the following important findings:

- A perfectly balanced workload distribution, such as HPS, may not necessary translate into better performance.
- Inter-processors communication overheads may cause dominant aggravation to the program elapsed time.
- A good mapping scheme should take into account the load balancing as well as communication overheads. Ignoring either factor may result in poor implementation performance.

6 Performance Improvement

A detailed study on using the results of performance modeling to improve the efficiency of simulation modeling is described in [6]. Due to space limitation we only highlight the speedup aspects. Using the MPS as an

example, we analyze how the overheads affect parallel simulation performance. We compared the implications of setting overheads to zero (*no overheads*), buffer access time to zero (*zero T_{buffer}*), reducing T_{buffer} to half ($T_{buffer}/2$), and setting message transmission time to zero (*zero $T_{transit}$*). Figure 4 shows that good parallel simulation speedup can be achieved even with the imposed transmission protocols, provided the communication overheads are negligible in comparison with the event grain size. We observed that the worst bottleneck is attributed to T_{buffer} followed by $T_{transit}$. Thus, more effort should be placed on reducing PVM-based memory allocation time, and improving access protocols to reduce overall message buffering time than on improving the network transmission speed.

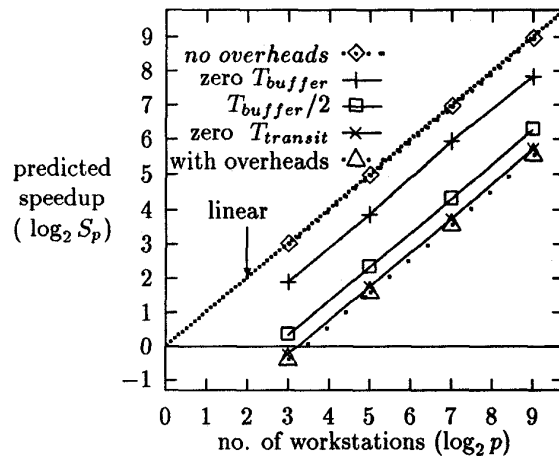


Figure 4: Speedup Analysis (using MPS Mapping Scheme)

7 Conclusions

We have developed an analytical model that characterized the performance of a conservative parallel simulation model using five timing parameters. Performance measures are derived for simulation elapsed time, speedup and bandwidth for both light and heavy

scheme		no. of messages per generator						$\bar{\delta}$ (%)
		120	240	360	480	600	720	
HPS	measured	0.61	0.62	0.62	0.61	0.62	0.63	6.31
	predicted	0.66	0.66	0.66	0.66	0.66	0.66	-
VPS	measured	0.59	0.58	0.58	0.59	0.60	0.59	6.61
	predicted	0.63	0.63	0.63	0.63	0.63	0.63	-
MPS	measured	0.72	0.71	0.72	0.71	0.71	0.71	4.89
	predicted	0.75	0.75	0.75	0.75	0.75	0.75	-

Table 7: Comparison of Speedup

scheme		no. of messages per generator						$\bar{\delta}$ (%)
		120	240	360	480	600	720	
HPS	measured	200.02	203.24	203.45	202.05	203.35	202.58	2.98
	predicted	209.85	209.87	209.86	209.86	209.87	209.87	-
VPS	measured	191.68	191.07	192.97	191.94	191.54	192.73	4.26
	predicted	200.12	200.41	200.51	200.57	200.59	200.61	-
MPS	measured	228.26	229.25	230.11	227.78	227.77	226.58	4.13
	predicted	237.85	238.09	238.16	238.20	238.23	238.24	-

Table 8: Comparison of Simulation Bandwidth (events/second)

traffic conditions. Validation experiments comparing the performance metrics from the analytical model and the simulation model show an acceptable 5% difference. We analyzed the performance for three process-to-processor mapping schemes to identify the main causes of poor performance. The performance model can be easily extended to analyze other applications using the same parallel mechanism. Our work reveals that in the PVM-based implementation reducing buffer access time is more critical than network transmission time. In addition, we dispel the common belief that to reduce the elapsed time of a parallel program the workload distribution among the processors must be balanced. We observed that other factors such as inter-processors communication overheads may also cause dominant aggravation to the simulation performance despite the balanced workload distribution. Therefore, a balanced workload distribution may not necessary translate into better performance. Both the analytical and implementation results confirm such an exceptional phenomenon.

References

- [1] P.A. Fishwick, "Simulation Model Design and Execution", Prentice-Hall, 1995.
- [2] R.M. Fujimoto, "Parallel Discrete Event Simulation", Comm. of ACM, Vol. 33, No. 10, pp. 31-53, October 1990.
- [3] D.R. Jefferson, "Virtual Time", ACM Trans. on Programming Languages and Systems, Vol. 7, No. 3, pp. 404-425, July 1985.
- [4] D. Nicol and R. Fujimoto, "Parallel Simulation Today", in Annals of Operations Research: Simulation and Modeling, edited by O. Balci, Vol. 53, 1994.
- [5] Y.M. Teo and S.C. Tay, "Modeling and Efficient Distributed Simulation of Multistage Interconnection Network", Proc. of International Conference on Algorithms and Architectures for Parallel Processing, Australia, IEEE Computer Society Press, pp. 83-92, April 1995.
- [6] S.C. Tay and Y.M. Teo, "Mapping Parallel Simulation onto Distributed-Memory Systems: Model and Performance", Technical Report TRB3/95, Department of Information Systems and Computer Science, National University of Singapore, pp. 24, March 1995.
- [7] S. Bokhari, "On Mapping Problem", IEEE Trans. of Computers, Vol. C-30, No. 3, pp. 207-214, March 1981.
- [8] A. Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, and V. Sunderan, "PVM 3 User's Guide and Reference Manual", Oak Ridge National Laboratory, Oak Ridge, Tennessee 37831, May 1993.