# Contention-Aware Metrics for Distributed Algorithms: Comparison of Atomic Broadcast Algorithms\*

Péter Urbán, Xavier Défago, and André Schiper

Département de Systèmes de Communication École Polytechnique Fédérale de Lausanne 1015 Lausanne EPFL, Switzerland E-mail: {peter.urban, xavier.defago, andre.schiper}@epfl.ch

#### Abstract

Resource contention is widely recognized as having a major impact on the performance of distributed algorithms. Nevertheless, the metrics that are commonly used to predict their performance take little or no account of contention. In this paper, we define two performance metrics for distributed algorithms that account for network contention as well as CPU contention. We then illustrate the use of these metrics by comparing four Atomic Broadcast algorithms, and show that our metrics allow for a deeper understanding of performance issues than conventional metrics.

# **1** Introduction

Performance prediction and evaluation are a central part of every scientific and engineering activity, including the construction of distributed applications. Engineers of distributed systems rely heavily on various performance evaluation techniques and have developed the necessary techniques for this activity. In contrast, algorithm designers invest considerable effort in proving the correctness of their algorithms (which they of course should do!), but often oversee the importance of predicting the performance of their algorithms, i.e., they rely on simplistic metrics. As a result, there is a serious gap between the prediction and the evaluation of performance of distributed algorithms.

**Performance Prediction vs. Evaluation of Algorithms.** When analyzing performance, one has to make a distinction between prediction and evaluation. Performance prediction gives an indication of the expected performance of an algorithm, *before* it is actually implemented. Performance prediction techniques give fairly general yet imprecise information, and rely on the use of various metrics. Conversely, performance evaluation is an *a posteriori* analysis of an algorithm, once it has been implemented and run in a given environment (possibly a simulation). While the information obtained is usually very accurate and precise, the results depend on specific characteristics of the environment and thus lack generality. Performance prediction and evaluation are complementary techniques. Performance prediction is used to orient design decisions, while performance evaluation can confirm those decisions and allows the dimensioning of the various system parameters.

Definition of a Metric. In this paper, we focus on the problem of predicting and comparing the performance of distributed algorithms. The goal is hence to investigate metrics that answer typical questions such as choosing the best algorithm for a particular problem, or identifying the various performance tradeoffs related to the problem. We define a metric as a value associated with an algorithm that has no physical reality and is used to define an order relation between algorithms. A good metric should provide a good approximation of the performance of an algorithm, regardless of the implementation environment<sup>1</sup>. Even though some performance evaluation techniques are also based on an abstract model of the system (e.g., analytical approaches, simulation), a metric must be a computable value. This is in contrast with simulation techniques that can model details of the system and the environment, and thus use complex models.

**Existing Metrics for Distributed Algorithms.** As mentioned earlier, performance prediction of distributed algorithms is usually based on two rather simplistic metrics: time and message complexity. These metrics are indeed useful, but there is still a large gap between the accuracy of the information they provide, and results obtained with more environment specific approaches.

The first commonly used metric, called *time complexity*, measures the latency of an algorithm. There exist many definitions of time complexity that are more or less equivalent. A common way to measure the time complexity of an algorithm (e.g., [1, 26, 22, 18, 14, 24]) consists in considering the algorithm in a model where the message delay has a known

<sup>\*</sup>Research supported by a grant from the CSEM Swiss Center for Electronics and Microtechnology, Inc., a company based in Neuchâtel.

<sup>&</sup>lt;sup>1</sup>A combination of several metrics, each focusing on different aspects of performance, might yield better insight into the behavior of a given algorithm than a single metric.

<sup>©2000,</sup> IEEE Computer Society Press. Appeared in Proc. of the 7<sup>th</sup> Int'l Conference on Computer Communications and Networks (IC3N'2000).

upper bound  $\delta$ . The efficiency of the algorithm is measured as the maximum time needed by the algorithm to terminate. This efficiency is expressed as a function of  $\delta$ , and is sometimes called the latency degree. This metric is *latency-oriented*, i.e., measures the cost of *one* execution of the algorithm.

The second metric, called *message complexity*, consists in counting the total number of messages generated by the algorithm [22, 14, 1]. This metric is useful when combined with time complexity, since two algorithms that have the same time complexity can generate a different volume of messages. Knowing the number of messages generated by an algorithm gives a good indication of its scalability and the amount of resources it uses. Furthermore, an algorithm that generates a large number of messages is likely to generate a high level of network contention.

**Resource Contention.** Resource contention is often a limiting factor for the performance of distributed algorithms. In a distributed system, the key resources are (1) the CPUs and (2) the network, any of which is a potential bottleneck. The major weakness of the time and message complexity metrics is that neither attaches enough importance to the problem of resource contention. While the message complexity metric ignores the contention on the CPUs, the time complexity metric ignores contention completely.

**Contribution and Structure.** In this paper, we define two metrics (one latency-oriented, the other throughput-oriented) which account for resource contention, both on the CPUs and the network. The use of those metrics is then illustrated by comparing Atomic Broadcast algorithms. The rest of the paper is structured as follows. Section 2 presents related work. In Section 3, we present the system model on which our metrics are based. Section 4 presents a *latency-oriented* and a *throughput-oriented* metric. We then compare some algorithms using our metrics in Section 5, and discuss the results. Finally, Section 6 concludes the paper.

# 2 Related Work

**Resource Contention in Network Models.** Resource contention (also sometimes called congestion) has been extensively studied in the literature. The bulk of the publications about resource contention describe strategies to either avoid or reduce resource contention (e.g. [15, 16]). Some of this work analyze the performance of the proposed strategies. However, these analyses consist in performance *evaluation*, and use models that are often specific to a particular network (e.g., [21]). Distributed algorithms are normally developed assuming the availability of some transport protocol. A metric that compares these algorithms must abstract out details that are only relevant to some implementations of a transport layer. In other words, it is necessary to relinquish precision for the sake of generality.

**Resource Contention in Parallel Systems.** Dwork, Herlihy, and Waarts [12] propose a complexity model for sharedmemory multiprocessors that takes contention into account. This model is very interesting in the context of shared memory systems but is not well suited to the message passing model that we consider here. The main problem is that the shared memory model is a high-level abstraction for communication between processes. As such, it hides many aspects of communication that are important in distributed systems. Dwork, Herlihy, and Waarts associate a unit cost based on the access to shared variables, which has a granularity too coarse for our problem.

**Computational Models for Parallel Algorithms.** Unlike distributed algorithms, many efforts have been directed at developing performance prediction tools for parallel algorithms. However, the execution models are not adapted to distributed algorithms: for instance, the PRAM model (e.g., [19]) requires that processors evolve in lock-steps and communicate using a global shared memory; the BSP model [31] requires that processors communicate using some global synchronization operation; the LogP model [9] assumes that there is an absolute upper bound on the transmission delay of messages. These models are not adequate to predict the performance of distributed algorithms. The main reason is that they do not naturally suit *asynchronous* distributed algorithms, which do not assume any form of global synchronization nor any restriction on communication delays.

**Competitive Analysis.** Other work, based on the method of competitive analysis proposed by Sleator and Tarjan [28], has focused on evaluating the competitiveness of distributed algorithms [5, 6]. In this work, the cost of a distributed algorithm is compared to the cost of an optimal centralized algorithm with a global knowledge. This work has been refined in [2, 3, 4] by considering an optimal *distributed* algorithm as the reference for the comparison. This work assumes an asynchronous shared-memory model and predicts the performance of an algorithm by counting the number of steps required by the algorithms to terminate. The idea of evaluating distributed algorithms against an optimal reference is appealing, but this approach is orthogonal to the definition of a metric. The metric used is designed for the shared-memory model, and still ignores the problem of contention.

### **3** Distributed System Model

The two metrics that we define in this paper are based on an abstract system model which introduces two levels of resource contention: *CPU contention* and *network contention*. First, we define a basic version of the model that leaves some aspects unspecified, but is sufficient to define our throughput oriented metric (see Definition 5). Second, we define an extended version of the model by removing the ambiguities left in the basic version. This extended model is used in Sect. 4 to define our latency oriented metric (see Definition 3).

#### 3.1 Basic Model

The model is inspired from the models proposed in [27, 29]. It is built around two types of resources: CPU and network. These resources are involved in the transmission of messages between processes. There is only one network that is shared among processes, and it is used to transmit a message from one process to another. Additionally, there is one CPU resource attached to each process in the system. These CPU resources represent the processing performed by the network controllers and the communication layers, during the emission and the reception of a message. In this model, the cost of running the distributed algorithm is neglected, and hence this does not require any CPU resource.

The transmission of a message m from a sending process  $p_i$  to a destination process  $p_j$  occurs as follows (see Fig. 1):

- 1. *m* enters the *sending queue*<sup>2</sup> of the sending host, waiting for  $CPU_i$  to be available.
- 2. *m* takes the resource  $CPU_i$  for  $\lambda$  time units, where  $\lambda$  is a parameter of the system model ( $\lambda \in \mathbb{R}_0^+$ ).
- 3. *m* enters the *network queue* of the sending host and waits until the network is available for transmission.
- 4. m takes the network resource for 1 time unit.
- 5. *m* enters the *receiving queue* of the destination host and waits until  $CPU_j$  is available.
- 6. *m* takes the resource  $CPU_j$  of the destination host for  $\lambda$  time units.
- 7. Message m is received by  $p_j$  in the algorithm.

#### 3.2 Extended Model

The basic model is not completely specified. For instance, it leaves unspecified the way some resource conflicts are resolved. We now extend the definition of the model in order to specify these points. As a result, the execution of a (deterministic) distributed algorithm in the extended system model is *deterministic*.

**Network** Concurrent requests to the network may arise when messages at different hosts are simultaneously ready for transmission. The access to the network is modeled by a round-robin policy,<sup>3</sup> illustrated by Algorithm 1.

**CPU** CPU resources also appear as points of contention between a message in the sending queue and a message in the receiving queue. This issue is solved by giving priority on every host to outgoing messages over incoming ones. Algorithm 1 Network access policy (executed by network).

Send to all Distributed algorithms often require to send a message m to all processes, using a "send to all" primitive. The way this is actually performed depends on the model (see below).

**Definition 1 (point-to-point)** Model  $\mathcal{M}_{pp}(n, \lambda)$  is the extended model with parameters  $n \in \mathbb{N}$  and  $\lambda \in \mathbb{R}_0^+$ , where n > 1 is the number of processes and  $\lambda$  is the relative cost between CPU and network. The primitive "send to all" is defined as follows: If p is a process that sends a message m to all processes, then p sends the message m consecutively to all processes in the lexicographical order  $(p_1, p_2, \ldots, p_n)$ .

Nowadays, many networks are capable of broadcasting information in an efficient manner, for instance, by providing support for IP multicast [10]. For this reason, we also define a model that integrates the notion of a broadcast network.

**Definition 2 (broadcast)** Model  $\mathcal{M}_{br}(n, \lambda)$  is defined similarly to Definition 1, with the exception of the "send to all" primitive, which is defined as follows: If p is a process that sends a message m to all, then p sends a single copy of m, the network transmits a single copy of m, and each process (except p) receives a copy of m.

# 3.3 Illustration

Let us now illustrate the model with an example. We consider a system with three processes  $\{p_1, p_2, p_3\}$  which execute the following simple algorithm. Process  $p_1$  starts the algorithm by sending a message  $m_1$  to processes  $p_2$  and  $p_3$ . Upon reception of  $m_1$ ,  $p_2$  sends a message  $m_2$  to  $p_1$  and  $p_3$ , and  $p_3$  sends a message  $m_3$  to  $p_1$  and  $p_2$ .

Figure 2 shows the execution of this simple algorithm in model  $\mathcal{M}_{pp}(3, 0.5)$ . The upper part of the figure is a timespace diagram showing the exchange of messages between the three processes. The lower part is a more detailed diagram that shows the activity (send, receive, transmit) of each resource in the model. For instance, process  $p_3$  sends a copy of message  $m_3$  to process  $p_1$  (denoted  $m_{3,1}$ ) at time 3. The message takes the CPU resource of  $p_3$  at time 3, takes the network resource at time 4.5, and takes the CPU resource of  $p_1$  at time 5.5. Finally,  $m_3$  is received by  $p_1$  at time 6.

 $<sup>^2\</sup>mbox{All}$  queues in the model use a FIFO policy (sending, receiving, and network queues).

<sup>&</sup>lt;sup>3</sup>Many thanks to Jean-Yves Le Boudec for this suggestion.



Figure 1. Decomposition of the end-to-end delay (tu=time unit).



Figure 2. Simple algorithm in model  $\mathcal{M}_{pp}(3, 0.5)$  ( $m_{i,j}$  denotes the copy of message  $m_i$  sent to process  $p_j$ ).

#### 4 Contention-Aware Metrics

### 4.1 Latency Metric

The definition of the latency metric uses the terms: "start" and "end" of a distributed algorithm. These terms are supposed to be defined by the problem  $\mathcal{P}$  that an algorithm  $\mathcal{A}$  solves. They are not defined as a part of the metric.

**Definition 3 (latency metric, point-to-point)** Let  $\mathcal{A}$ be a distributed algorithm. The latency metric Latency<sub>pp</sub>( $\mathcal{A}$ )( $n, \lambda$ ) is defined as the number of time units that separate the start and the end of algorithm  $\mathcal{A}$  in model  $\mathcal{M}_{pp}(n, \lambda)$ .

**Definition 4 (latency metric, broadcast)** *Let*  $\mathcal{A}$  *be a distributed algorithm. The latency metric*  $\text{Latency}_{\text{br}}(\mathcal{A})(n, \lambda)$  *is defined as the number of time units that separate the start and the end of algorithm*  $\mathcal{A}$  *in model*  $\mathcal{M}_{\text{br}}(n, \lambda)$ .

### 4.2 Throughput Metric

The throughput metric of an algorithm  $\mathcal{A}$  considers the utilization of system resources in one run of  $\mathcal{A}$ . The most

heavily used resource constitutes a bottleneck, which puts a limit on the *maximal throughput*, defined as an upper bound on the frequency at which the algorithm can be run.

**Definition 5 (throughput metric, point-to-point)** Let  $\mathcal{A}$  be a distributed algorithm. The throughput metric is defined as follows:

Thput<sub>pp</sub>(
$$\mathcal{A}$$
) $(n, \lambda) \stackrel{\text{def}}{=} \frac{1}{\max_{r \in \mathcal{R}_n} T_r(n, \lambda)}$ 

where  $\mathcal{R}_n$  denotes the set of all resources (i.e.,  $CPU_1, \ldots, CPU_n$  and the network), and  $T_r(n, \lambda)$  denotes the total duration for which resource  $r \in \mathcal{R}_n$  is utilized in one run of algorithm  $\mathcal{A}$  in model  $\mathcal{M}_{pp}(n, \lambda)$ .

Thput<sub>pp</sub>( $\mathcal{A}$ ) $(n, \lambda)$  can be understood as an upper bound on the frequency at which algorithm  $\mathcal{A}$  can be started. Let  $r_b$  be the resource with the highest utilization time:  $T_{r_b} = \max_{r \in \mathcal{R}_n} T_r$ . At the frequency given by Thput<sub>pp</sub>( $\mathcal{A}$ ) $(n, \lambda)$ ,  $r_b$  is utilized at 100%, i.e., it becomes a bottleneck.

**Definition 6 (throughput metric, broadcast)** Let  $\mathcal{A}$  be a distributed algorithm. The definition of the throughput metric Thput<sub>br</sub>( $\mathcal{A}$ ) $(n, \lambda)$  is the same than Definition 5, but in model  $\mathcal{M}_{br}(n, \lambda)$ .

**Relation with Message Complexity.** The throughput metric can be seen as a generalization of the message complexity metric. While our metric considers different types of resources, message complexity only considers the network. It is easy to see that  $T_{network}$ , the utilization time of the network in a single run, gives the number of messages exchanged in the algorithm.

# 5 Comparison of Atomic Broadcast Algorithms

We now illustrate the use of our two metrics by comparing four different algorithms that solve the problem of *Atomic Broadcast*. These examples show that our metrics yield results that are more precise than what can be obtained by relying solely on time and message complexity. This confirms the observation that contention is a factor that cannot be overlooked.

A more extensive analysis of total order broadcast algorithms (using the same metrics) appears in [11].

#### 5.1 Atomic Broadcast Algorithms

Atomic Broadcast is a fundamental problem in the context of distributed systems [14]. Informally, the problem consists in broadcasting messages to other processes, in such a way that all messages are delivered in the same order by all destination processes. The problem is defined in terms of the two events *A*-*Broadcast* and *A*-*Deliver*. When a process wants to atomically broadcast a message *m* it executes A-Broadcast(*m*), and A-Deliver(*m*) executed by process *q* corresponds to the delivery of message *m* by *q*. The latency of the algorithm with respect to message *m* is then defined as follows. We consider a run in which no other message is A-Broadcast; the algorithm starts when a process executes A-Broadcast(*m*) and ends when the last process executes A-Deliver(*m*).

We briefly describe four different Atomic Broadcast algorithms for a system with no failures, and compare them using our metrics. Figure 3 shows the communication pattern associated with the broadcast of a single message m for each of the four algorithms. Note that the communication pattern is enough to compute our metrics. For this reason, we have omitted to give the details of each algorithm.

**Lamport.** In Lamport's algorithm [20], every message carries a logical time-stamp. To atomically broadcast a message m, the sender process first sends m to all other processes (Fig. 3(a)). Upon reception of m, a process p sends a time-stamped "null message" to all others, thus informing them that it has no other message that may have to be delivered before m. These null messages appear only when a process has no message to broadcast.

**Skeen.** Skeen's algorithm (described in [8, 13, 25]) is a twophase protocol that can use Lamport's logical clocks [20]. To atomically broadcast a message m, a process p first sends m to all processes (Fig. 3(b)). Upon reception of m, the processes send a time-stamped acknowledgment message to p. Once p has received all acknowledgments, it takes the maximum of the time-stamps received, and sends this information to all processes. Processes deliver m after they receive this message (the details of the delivery condition are irrelevant here).

**Token.** In Rajagopalan and McKinley's token-based algorithm [23], a token circulates in the system and a process is allowed to broadcast messages only when it holds the token. To atomically broadcast a message m, a process p must first wait for the token<sup>4</sup> (Fig. 3(c)). When it holds the token, p broadcasts m to the other processes and passes the token to the next process. The message m can be delivered only after it has been acknowledged by all processes. The acknowledgments of messages are carried by the token. So m is delivered by the last process only after two round-trips of the token.

**Sequencer.** Many Atomic Broadcast algorithms are based on the principle that one process is designated as a sequencer and constructs the order (e.g., [7, 17]). In the version that we consider here (Fig. 3(d)), a process atomically broadcasts a message m by sending m to the sequencer. Upon reception of m, the sequencer attaches a sequence number to m and sends it to all other processes. Messages are then delivered according to their sequence number.

These algorithms are interesting to illustrate our metrics because they take contrasting approaches to solve the problem of Atomic Broadcast. Although they all deliver messages according to some total order, these algorithms actually provide varying levels of guarantees, and are hence not equivalent. An actual comparison must take these issues into account.

### 5.2 Latency Metric

We now analyze the latency of the four Atomic Broadcast algorithms: Lamport, Skeen, Token, and Sequencer. For each algorithm, we compute the value of the latency metric in model  $\mathcal{M}_{\rm PP}(n, \lambda)$ . The results are summarized in Table 1 and compared in Fig. 4(a).<sup>5</sup> Table 1 also shows the time complexity of the algorithms. For time complexity, we use the *latency degree* [26]: roughly speaking, an algorithm with latency degree *l* requires *l* communication steps.

Figure 4(a) represents the results of the comparison between the four algorithms with respect to the latency metric. The area is split into three zones in which algorithms perform differently with respect to each other (e.g., in Zone I, we have Sequencer > Lamport > Skeen > Token, where > means "better than"). The latency metric and time complexity yield the same results

<sup>&</sup>lt;sup>4</sup>In our analysis, we take the average case where the token is always halfway on its path toward p.

<sup>&</sup>lt;sup>5</sup>For reasons of clarity, we choose to give approximate formulas for Latency<sub>pp</sub>(Lamport) $(n, \lambda)$  and Latency<sub>pp</sub>(Skeen) $(n, \lambda)$ . The expressions given for these two algorithms ignore a factor that is negligible compared to the rest of the expression. The exact expressions, as well as a description of the analysis are given in [30].



Figure 3. Communication patterns of Atomic Broadcast algorithms.

Table 1. Latency metric	evaluation of Atomic	Broadcast algorithms	(in model $\mathcal{M}$	$l_{\rm pp}(n)$	$,\lambda)$	))
-------------------------	----------------------	----------------------	-------------------------	-----------------	-------------	----

Algorithm $\mathcal{A}$	$\operatorname{Latency}_{\mathrm{pp}}(\mathcal{A})(n,\lambda)$		Time complexity
	$\approx 3(n-1)\lambda + 1$	$\text{if } n \leq \lambda + 2$	
Lamport	$\approx \frac{1}{2}n(n-3) + 2\lambda n + \frac{1}{2}\lambda^2 - \frac{3}{2}\lambda$	$\text{if } n \leq 2\lambda + 3$	2
	$\approx \frac{1}{2}n(n-1) + 2\lambda n + \lambda^2 - \frac{7}{2}\lambda - 3$	$\text{if } n \leq 4\lambda - 4$	2
	$\approx n(n-1) + \lambda^2 + \lambda + 5$	otherwise	
Skeen	$\approx 3(n-1) + 4\lambda$	$\text{if }\lambda <1$	3
	$\approx (3n-2)\lambda + 1$	$\text{if }\lambda\geq 1$	3
Token	$(2.5n-1)(2\lambda+1) + \max(1,\lambda)(n-1)$		2.5n - 1
Sequencer	$4\lambda + 2 + \max(1,\lambda)(n-2)$		2

for three of the four algorithms: Token, Skeen, and Sequencer. Both metrics yield that Sequencer performs better than Skeen, which in turn performs better than Token. For Lamport, time complexity (Table 1) suggests that it always performs better than the other algorithms. This comes in contrast with our latency metric which shows that the relative performance of Lamport are dependent on the system parameters n and  $\lambda$ . The reason is that Lamport generates a quadratic number of messages and is hence subject to network contention to a greater extent. Time complexity is unable to predict this as it fails to account for contention.

#### 5.3 Throughput Metric

We now analyze the throughput of the four algorithms. In a throughput analysis, one run of the algorithm should not be considered in isolation. Indeed, many algorithms behave differently whether they are under high load or not (e.g., Lamport does not need to generate null messages under high load). For this reason, the throughput metric is computed by considering a run of the algorithm *under high load*. We also assume that every process atomically broadcasts messages, and that the emission is fairly distributed among them. For each algorithm, we compute the value of the throughput metric in model  $\mathcal{M}_{pp}(n, \lambda)$ . The results are summarized<sup>6</sup> in Table 2. The algorithms are then compared in Fig. 4(b).

Table 2. Throughput metric: evaluation of Atomic Broadcast algorithms (in model  $\mathcal{M}_{pp}(n, \lambda)$ )

Algorithm ${\cal A}$	$(\mathrm{Thput}_{\mathrm{pp}}(\mathcal{A})(n,\lambda))^{-1}$	Message complexity
Lamport	$(n-1) \cdot \max(1, \frac{2\lambda}{n})$	n-1
Skeen	$3(n-1) \cdot \max(1, \frac{2\lambda}{n})$	3(n-1)
Token	$n \cdot \max(1, \frac{2\lambda}{n})$	n
Sequencer	$(n - \frac{1}{n}) \cdot \max(1, \lambda)$	$n-\frac{1}{n}$

Figure 4(b) illustrates the relative throughput of the four algorithms. The graph is split into three zones in which algorithms perform differently with respect to each other. The throughput metric and message complexity both yield that Lamport performs better than Token which in turn performs better than Skeen. However, the two metrics diverge when considering Sequencer. Indeed, while message complexity (Table 2) suggests that Sequencer always performs better than Skeen and Token, our throughput metric shows that it is not always the case. In fact, Sequencer is more subject to CPU contention than the other three algorithms. This type of contention is especially noticeable in systems with large values of  $\lambda$ . Message complexity fails to pinpoint this, as it does not take CPU contention into account.

<sup>&</sup>lt;sup>6</sup>The full description of the analysis is given in [30].



Figure 4. Comparison of Atomic Broadcast algorithms (A > A' means A "better than" A').

### 5.4 Latency and Throughput in Broadcast Networks

The analyses in model  $\mathcal{M}_{br}(n, \lambda)$  are not much different. In fact, there are less messages and less contention<sup>7</sup>. Table 3

Table	3.	Latency <sub>br</sub> $(\mathcal{A})(n, \lambda)$ :	evaluation	of
Atomie	c Br	oadcast algorithms.		

Algorithm $\mathcal{A}$	$Latency_{br}(\mathcal{A})(n,\lambda)$		
Lamport	$4\lambda + n$		
Skeen	$6\lambda + 3 + (n-2) \cdot \max(1,\lambda)$		
Token	$\left(\frac{5n}{2}-1\right)(2\lambda+1) + \max(1,\lambda)$		
Sequencer	$4\lambda + 2$		

Table 4.  $\mathrm{Thput}_{\mathrm{br}}(\mathcal{A})(n,\lambda)$ : evaluation of Atomic Broadcast algorithms.

Algorithm $\mathcal{A}$	$(\mathrm{Thput}_{\mathrm{br}}(\mathcal{A})(n,\lambda))^{-1}$	Msg complexity
Lamport	$\max(1,\lambda)$	1
Skeen	$\max(n+1, \frac{4n+1}{n}\lambda)$	n+1
Token	$\max(2, \frac{n+2}{n}\lambda)$	2
Sequencer	$\frac{2n-1}{n} \max(1,\lambda)$	$2 - \frac{1}{n}$

and Table 4 show the results of the two metrics in a broadcast network (Latency<sub>br</sub>( $\mathcal{A}$ )( $n, \lambda$ ) and Thput<sub>br</sub>( $\mathcal{A}$ )( $n, \lambda$ )). Apart from the fact that these results are simpler than in a model with point-to-point communication, there are interesting differences.

According to the latency metric, for any "realistic" value<sup>8</sup> of  $\lambda$  and n, the algorithms are always ordered as follows:

Sequencer > Lamport > Skeen > Token

Unlike the results obtained with  $\text{Latency}_{pp}(\mathcal{A})(n, \lambda)$ , there is only one single zone with a broadcast network. This zone corresponds to zone I depicted on Figure 4(a) but, in model  $\mathcal{M}_{br}(n, \lambda)$ , the algorithms are not ordered differently as *n* increases. This is easily explained by the fact that Lamport is quadratic in model  $\mathcal{M}_{pp}(n, \lambda)$  while it is linear in model  $\mathcal{M}_{br}(n, \lambda)$ . The latency of the three other algorithms is not so different because they are linear in both models.

Similarly, Thput<sub>br</sub>( $\mathcal{A}$ ) $(n, \lambda)$  yields simpler results than Thput<sub>pp</sub>( $\mathcal{A}$ ) $(n, \lambda)$ . As shown in Figure 4(c), the parameter space is cut into two zones (instead of three for Thput<sub>pp</sub>( $\mathcal{A}$ ) $(n, \lambda)$ , as shown on Fig. 4(b)). The difference between the two zones is the relative performance (throughput) of Sequencer and Token. This yields that Token is better than Sequencer when the CPU is a limiting factor. In fact, Sequencer is limited by the sequencer process which becomes a clear bottleneck. Conversely, Token spreads the load evenly among all processes, and so none becomes a bottleneck. Once again, both classical metrics (time and message complexity) fail to capture this aspect.

### 6 Conclusion

The paper proposes two metrics to predict the latency and the throughput of distributed algorithms. Unlike other existing metrics, the two complementary metrics that we present here take account of both network and CPU contention. This allows for more precise predictions and a finer grained analysis of algorithms than what time complexity and message complexity permit. In addition, our metrics make it possible to find out whether the bottleneck is the network or the CPU of one specific process.

The problem of resource contention is commonly recognized as having a major impact on the performance of distributed algorithms. Because other metrics do not take account of contention to the same extent as ours, our metrics fill a gap that exists between simple complexity measures and more complex performance evaluation techniques.

The system model for the metrics presented here can be

 $<sup>^{7}\</sup>text{The}$  full description of the analysis in model  $\mathcal{M}_{\mathrm{br}}(n,\lambda)$  is given in [30].

<sup>&</sup>lt;sup>8</sup>Realistic values for the parameters  $\lambda$  and n are:  $\lambda \ge 0$  and  $n \ge 2$ .

extended in a variety of ways. Modeling a separate network processor beside the CPUs would bring it closer to the architecture of current networks. Also, the simple bus-based network could be replaced by more complex topologies<sup>9</sup>. Careful experimentation is needed to decide which extensions result in more realistic models, without making the computation of the metrics unnecessarily difficult.

#### Acknowledgments

We would like to thank Jean-Yves Le Boudec for his numerous comments and advice on early versions of this work. We would also like to thank the anonymous reviewers for their comments, especially their ideas on how to develop our metrics further.

# References

- [1] M. K. Aguilera, W. Chen, and S. Toueg. Failure detection and consensus in the crash-recovery model. In *Proc. 12th Int'l Symp. on Distributed Computing (DISC)*, pages 231– 245, Sept. 1998.
- [2] M. Ajtai, J. Aspnes, C. Dwork, and O. Waarts. A theory of competitive analysis for distributed algorithms. In S. Goldwasser, editor, *Proc. 35th Annual Symp. on Foundations of Computer Science*, pages 401–411, Nov. 1994.
- [3] J. Aspnes and O. Waarts. A modular measure of competitiveness for distributed algorithms (abstract). In Proc. 14th ACM Symp. on Principles of Distributed Computing (PODC), page 252, Aug. 1995.
- [4] J. Aspnes and O. Waarts. Modular competitiveness for distributed algorithms. In *Proc. 28th ACM Symp. on Theory of Computing (STOC)*, pages 237–246, May 1996.
- [5] B. Awerbuch, S. Kutten, and D. Peleg. Competitive distributed job scheduling. In *Proc. 24th ACM Symp. on Theory* of *Computing (STOC)*, pages 571–580, May 1992.
- [6] Y. Bartal, A. Fiat, and Y. Rabani. Competitive algorithms for distributed data management. In *Proc. 24th ACM Symp.* on *Theory of Computing (STOC)*, pages 39–50, May 1992.
- [7] K. Birman, A. Schiper, and P. Stephenson. Lightweight causal and atomic group multicast. *ACM Trans. Comput. Syst.*, 9(3):272–314, Aug. 1991.
- [8] K. P. Birman and T. A. Joseph. Reliable communication in presence of failures. ACM Trans. Comput. Syst., 5(1):47–76, Feb. 1987.
- [9] D. E. Culler, R. M. Karp, D. Patterson, A. Sahay, E. E. Santos, K. E. Schauser, R. Subramonian, and T. von Eicken. LogP: A practical model of parallel computation. *Commun. ACM*, 39(11):78–85, Nov. 1996.
- [10] S. E. Deering. RFC 1112: Host extensions for IP multicasting, Aug. 1989.
- [11] X. Défago. Agreement-Related Problems: From Semi-Passive Replication to Totally Ordered Broadcast. PhD thesis, École Polytechnique Fédérale de Lausanne, Switzerland, Aug. 2000. Number 2229.
- [12] C. Dwork, M. Herlihy, and O. Waarts. Contention in shared memory algorithms. J. ACM, 44(6):779–805, Nov. 1997.

- [13] R. Guerraoui and A. Schiper. Total order multicast to multiple groups. In Proc. 17th Int'l Conf. on Distributed Computing Systems (ICDCS), pages 578–585, May 1997.
- [14] V. Hadzilacos and S. Toueg. Fault-tolerant broadcasts and related problems. In S. Mullender, editor, *Distributed Systems*, chapter 5, pages 97–146. Second edition, 1993.
- [15] A. Heddaya and K. Park. Congestion control for asynchronous parallel computing on workstation networks. *Parallel Computing*, 23(13):1855–1875, Dec. 1997.
- [16] J.-H. Huang, C.-C. Yang, and N.-C. Fang. A novel congestion control mechanism for multicast real-time connections. *Computer Communications*, 22:56–72, 1999.
- [17] M. F. Kaashoek and A. S. Tanenbaum. Fault tolerance using group communication. ACM Operating Systems Review, 25(2):71–74, Apr. 1991.
- [18] E. V. Krishnamurthy. Complexity issues in parallel and distributed computing. In A. Y. H. Zomaya, editor, *Parallel & Distributed Computing Handbook*, pages 89–126. 1996.
- [19] L. I. Kronsjö. PRAM models. In A. Y. H. Zomaya, editor, Parallel & Distributed Computing Handbook, pages 163– 191. McGraw-Hill, 1996.
- [20] L. Lamport. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM*, 21(7):558–565, July 1978.
- [21] C.-C. Lim, L.-J. Yao, and W. Zhao. A comparative study of three token ring protocols for real-time communications. In *Proc. 11th Int'l Conf. on Distributed Computing Systems* (*ICDCS*), pages 308–317, May 1991.
- [22] N. A. Lynch. Distributed Algorithms. Morgan Kaufmann, 1996.
- [23] B. Rajagopalan and P. McKinley. A token-based protocol for reliable, ordered multicast communication. In *Proc. 8th Symp. on Reliable Distributed Systems (SRDS)*, pages 84– 93, Oct. 1989.
- [24] M. Raynal. Networks and Distributed Computation. MIT Press, 1988.
- [25] L. Rodrigues, R. Guerraoui, and A. Schiper. Scalable atomic multicast. In Proc. 7th IEEE Int'l Conf. on Computer Communications and Networks (IC3N'98), pages 840–847, Oct. 1998.
- [26] A. Schiper. Early consensus in an asynchronous system with a weak failure detector. *Distributed Computing*, 10(3):149– 157, 1997.
- [27] N. Sergent. Soft Real-Time Analysis of Asynchronous Agreement Algorithms Using Petri Nets. PhD thesis, École Polytechnique Fédérale de Lausanne, Switzerland, 1998.
- [28] D. D. Sleator and R. E. Tarjan. Amortised efficiency of list update and paging rules. *Commun. ACM*, 28(2):202–208, Feb. 1985.
- [29] K. Tindell, A. Burns, and A. J. Wellings. Analysis of hard real-time communications. *Real-Time Systems*, 9(2):147– 171, Sept. 1995.
- [30] P. Urbán, X. Défago, and A. Schiper. Contention-aware metrics: Analysis of distributed algorithms. Technical Report DSC/2000/012, École Polytechnique Fédérale de Lausanne, Switzerland, Feb. 2000.
- [31] L. G. Valiant. A bridging model for parallel architectures. *Commun. ACM*, 33(8):103–111, Aug. 1990.

<sup>&</sup>lt;sup>9</sup>We primarily see the bus-based network of the model as the simplest mechanism we could think of to introduce network contention.