

On the Performance of Synchronized Programs in Distributed Networks with Random Processing Times and Transmission Delays

Sergio Rajsbaum and Moshe Sidi, *Senior Member, IEEE*

Abstract— A synchronizer is a compiler that transforms a program designed to run in a synchronous network into a program that runs in an asynchronous network. The behavior of a simple synchronizer, which also represents a basic mechanism for distributed computing and for the analysis of marked graphs, was studied by Even and Rajsbaum under the assumption that message transmission delays and processing times are constant. In this paper, we study the behavior of the simple synchronizer when processing times and transmission delays are random. Our main performance measure is the rate of a network, i.e., the average number of computational steps executed by a processor in the network per unit time. We analyze the effect of the topology and the probability distributions of the random variables on the behavior of the network. For random variables with exponential distribution, we provide tight (i.e., attainable) bounds and study the effect of a bottleneck processor on the rate.

Index Terms— Distributed networks, synchronizer, performance analysis, stochastic behavior, marked graphs

I. INTRODUCTION

CONSIDER a network of processors that communicate by sending messages along communication links. The network is *synchronous* if there is a global clock whose beats are heard by all the processors simultaneously, and the time interval between clock beats is long enough for all messages to reach their destinations and for local *computational steps* to be completed before the clock beats again. The network is *asynchronous* if there is no global clock, and the transmission times of messages are unpredictable.

In general, a program designed for a synchronous network will not run correctly in an asynchronous network. Instead of designing a new program for the asynchronous network, it is possible to use a *synchronizer* [1], i.e., a compiler that converts a program designed for a synchronous network, to run correctly in an asynchronous network. Synchronizers are a useful tool because programs for synchronous networks are easier to design, debug, and test than are programs for asynchronous networks. Furthermore, an important use of synchronizers is the design of more efficient asynchronous algorithms [2]. The

problem of designing efficient synchronizers has been studied in the past (e.g., [1], [3], [28]).

The (worst case) time complexity of a distributed algorithm is usually computed under the assumption that processing times and message transmission delays are equal to some constant that represents an upper bound on these durations. The goal of this paper is to study the effect of *random* processing times and transmission delays on the performance of synchronous programs running in an asynchronous network under the control of a simple synchronizer. We compare the results with the deterministic case [18], [19], in which processing times, as well as message delays, are constant (or bounded).

The operation of the synchronizer is as follows. Each processor waits for a message to arrive on each of its incoming links before performing the next computational step. When a computational step is completed (after a random time), it sends one message on each of its outgoing links. The implementation of this synchronizer may require, for instance, that every message be followed by an end-of-message marker, even if the message is empty. These end-of-message markers model the flow of information that must exist between every pair of processors connected by a link in each computational step [1]. This is how a processor knows it has to wait for a message that was sent to it, or if no message was sent.

We use this synchronizer in our analysis because it is very simple, yet it captures the essence of the synchronizer methodology; i.e., it ensures that a processor does not initiate a new phase of computation before knowing that all the messages sent to it during the previous phase have already arrived. Moreover, the synchronizer is equivalent to a marked graph (e.g., [15]) in which the initial marking has one token per edge. In [31] and [32], the relationship between synchronizers and marked graphs is studied, and it is shown how the simple synchronizer can model the behavior of any marked graph, of the synchronizers of [1], and of distributed schedulers in [6], [23]. Thus, our work is closely related to problems in stochastic Petri nets, where, because of the huge size of the state space, the solution techniques often rely on simulation (e.g., [22], [24], [25]).

Many distributed protocols are based on this simple synchronizer, e.g., the snapshot algorithm [16], clock synchronization algorithms (e.g. [12], [26]), the synchronizers of [1], the distributed schedulers in [6], [23], and the *optimistic* synchronizer [20]. The synchronizer is similar to synchronizer

Manuscript received October 31, 1992; revised June 22, 1993.

S. Rajsbaum is with the Instituto de Matematicas, UNAM, Ciudad Universitaria, D.F. 04510 Mexico; e-mail: rajsbaum@redvax1.dgsc.unam.mx.

M. Sidi is with the Department of Electrical Engineering, Technion—Israel Institute of Technology, Haifa, Israel 32000; e-mail: moshe@ee.technion.ac.il.

IEEE Log Number 9401219.

α in [1], but can be used also in directed networks, as opposed to other synchronizers suggested in [1] that require all links to be bidirectional. In [18] and [19], the benefits of using the synchronizer as an initialization procedure are described.

Main Results: This paper is devoted to the performance analysis of strongly connected directed networks controlled by the simple synchronizer, in which transmission delays, as well as the time it takes a processor to complete a computational step are random variables. Our main performance measure is the *rate* of computation R_v , i.e., the average number of computational steps executed by a processor in the network per unit time. To facilitate the presentation, we first assume that the transmission delays are negligible; only at the end of the paper, we describe how to extend the results for networks with non-negligible delays.

In Section III, we study the case in which the random variables have general probability distributions. We consider two approaches. First, in Section III-A, we analyze the effect of the topology on the rate. We use stochastic comparison techniques to compare the rate of networks with different topologies. We give examples of networks with different topologies, but with the same rate. Then, in Section III-B, we analyze networks with the same topology, but different processing times. By defining a partial order on the set of distributions, we show that deterministic (i.e., constant) processing times maximize the rate of computation. For this case, it is shown in [18] that if the processing times are equal to λ^{-1} , the rate of the network is λ , regardless of the number of processors in the network or its topology. In the next section, we show that in case the processing times are random and unbounded, the rate may be degraded by a logarithmic factor in the number of processors. This occurs in the case of exponentially distributed processing times. However, in this section, we show that the exponential is the worst among a large and natural class of distributions; i.e., it yields the minimum rate within a class of distributions.

In Section IV, we concentrate on the case of processing times that are exponentially distributed random variables with mean λ^{-1} . We prove that the rate is between $\lambda/4 \log(\Delta + 1)$ and $\lambda/\log(\delta + 1)$, where Δ (δ) is the maximum (minimum) vertex in-degree or out-degree. Hence, for regular-degree (either in or out-degree) networks, the rate is $\Theta(\lambda/\log(\delta + 1))$. We compute the exact rate and the stationary probabilities for the extreme cases of a directed cycle and a complete graph. Finally, we study the effect of having one processor that runs slower than the rest of the processors, and we show that in some sense, the directed cycle network is more sensitive to such a bottleneck processor than a complete network.

In the last section, we show that it is easy to extend the results to networks with non-negligible transmission delays. We consider the exponential distribution case, and show that adding transmission delays to a regular degree network may reduce its rate by at most a constant factor, provided that they are not larger (w.r.t. the partial order) than the processing times. In networks with processing times exponentially distributed with mean 1, and larger delays with mean λ^{-1} , we compare the results with those of [19], where it was shown that for the corresponding deterministic case, the rate is λ .

In the probabilistic case of a regular-degree network, the rate is at least $\Theta(\lambda/\log \delta)$. Thus, in both cases (small and large delays), the rate of a bounded degree network is reduced only by a constant factor.

Previous Work: There exist several results related to our results in Section III-B in the literature on stochastic Petri nets. For instance, dominance results for rather general stochastic Petri nets have been obtained in [5], and, more recently, in [8], by using subadditive ergodic theory (e.g., [21]). It should be noted, however, that the proofs we provide for the simple synchronizer are different and much simpler, and do not require heavy mathematical tools. Other stochastic ordering studies exist. Papers on acyclic networks and fork-join queues are [29] and [9]–[11], respectively. For closed queueing networks, the effect of increasing the service rate of a subset of stations for systems, such that the distribution of the number of works in each station has a product form solution, is studied in [35].

A model similar to our model in Section IV is considered in [13], where it is claimed that the rate is $\theta(1/\log \delta_{\text{out}})$, for regular networks with out-degree equal to δ_{out} , identically exponentially distributed transmission delays with mean 1, and negligible processing times. In [12], only a lower bound of $\Theta(1/\log \delta_{\text{in}})$ on the rate is given, for regular networks with in-degree equal to δ_{in} , with negligible transmission delays, and identically exponentially distributed processing times. Recently, it has been shown in [7] that subadditive ergodic theory can be used to derive more general lower bounds on the rate. A bottleneck problem related to ours has been considered in [4], where an asymptotic analysis of cyclic queues as the number of costumers grows is presented. Asymptotic performance of stochastic marked graphs as the number of tokens grows is studied in [25]. The class of networks with exponentially distributed processing times belongs to the more general model of stochastic Petri nets (see [22] for a survey), where it is usually assumed that the state space (of exponential size, in our case) is given.

II. THE MODEL

The *network* is modeled by a (finite) directed, strongly connected graph $G(V, E)$, where $V = \{1, 2, \dots, n\}$ is the set of vertices of the graph and $E \subseteq V \times V$ is the set of directed edges. A vertex of the graph corresponds to a processor that is running its own program, and a directed edge $u \rightarrow v$ corresponds to a communication link from processor u to processor v . In this case, we shall say that u is an *in-neighbor* of v , and that v is an *out-neighbor* of u in the network. The processors communicate by sending messages along the communication links. To facilitate the presentation, we assume that the message transmission delays are negligible. At the end of this paper, we briefly discuss the case of non-negligible transmission delays.

Initially, all processors are in a *quiescent state*, in which they send no messages and perform no computations. Once a processor leaves the quiescent state, it never reenters it and is considered *awake*. When awakened, each processor operates in phases as described in the sequel. Assume that at an arbitrary

time, $t(v)$, processor v leaves the quiescent state and enters its first *processing state*, PS_0 . This may be caused by a message from another processor, or by a signal from the outside world, not considered in our model. Then processor v remains in PS_0 for $\tau_0(v)$ units of time and then transits to its first *waiting state*, WS_0 . From this time on, let PS_k and $WS_k, k \geq 0$, denote the processing state and the waiting state, respectively, for the k th phase. Observe that we are concerned with the rate of computation of the network; the nature of the computation is of no concern to us here. Thus, we take the liberty of denoting with the same symbol the k th processing state of all the processors.

The transition rules between states are as follows. If a processor v transits from state PS_k to WS_k , it sends one message on each of its outgoing edges. These messages are denoted by M_k . Note that this labeling is not needed for the implementation of the protocol; it is used only for its analysis. When v sends the M_k messages, we say that v has completed its k th *processing step*.

If a processor v is in state WS_k , and has received a message (M_k) on each of its incoming edges, it removes one message from each of its incoming edges, transits to state PS_{k+1} , remains there for $\tau_{k+1}(v)$ units of time, and then transits to state WS_{k+1} . Otherwise, if at least on one incoming edge, M_k , has not yet arrived, processor v remains in state WS_k until it receives a message from each of its in-neighbors, and then operates as described above.

The *processing times*, $\tau_k(v)$, correspond to the time it takes for processor v to complete the k th computation step. The processing times $\tau_k(v), k \geq 0$, and $v \in V$, are positive, real-valued random variables defined over some probability space.

For $k \geq 0$, let $t_k^G(v)$ (or $t_k(v)$, whenever G is understood) be the k th completion time, i.e., the time at which processor v sends messages M_k in network G . Let the *in-set* of a vertex v in G , $IN^G(v)$ (or simply $IN(v)$), be the set of vertices in G that have an edge to v , including v itself, that is, $IN(v) = \{u : u \rightarrow v \in E\} \cup \{v\}$. With this notation, the operation of processor $v \in V$ is as follows. Once v has sent a message M_k at time $t_k(v)$, it waits until all processors with an edge to it send message M_k , and then starts its $(k + 1)$ st computation step; that is, after the maximum of $t_k(u), u \in IN(v)$, it starts the $(k + 1)$ st computation step, which takes $\tau_{k+1}(v)$ units of time, and then sends out M_{k+1} . For this reason, we shall assume in the rest of the paper that for each vertex v , the edge $v \rightarrow v$ is in E . The evolution of the network can be described by the following recursions:

$$t_0(v) = t(v) + \tau_0(v)$$

$$t_{k+1}(v) = \max_{u \in IN(v)} \{t_k(u)\} + \tau_{k+1}(v), \quad k \geq 0. \quad (1)$$

It is interesting to note that the completion times $t_k(v)$ have a simple graph theoretic interpretation. For a vertex v , let $S_k(v)$ be the set of all directed paths of length k ending in v . For $k = 0$, the only path of length 0 ending in v consists of v itself. For a path $P_k = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k (= v)$, let $T(P_k) \triangleq t(v_0) + \sum_{i=0}^k \tau_i(v_i)$, and let $T(S_k(v)) \triangleq \{T(P) :$

$P \in S_k(v)\}$. Thus, $T(S_k(v))$ is a set of random variables; each one is the sum of $k + 1$ random variables. Note that these random variables are not independent, even if the $\tau_i(v)$'s are independent. The explicit computation of $t_k(v)$ is as follows.

Theorem 2.1: For every $v \in V, k \geq 0, t_k(v) = \max T(S_k(v))$.

Proof: By induction on k . For $k = 0$, note that the only path of length 0 to v is v itself; i.e., $S_0(v) = \{v\}$, and $T(S_0(v)) = \{t(v_0) + \tau_0(v)\}$. Hence, we get the following equation:

$$t_0(v) = \max T(S_0(v)) = t(v) + \tau_0(v).$$

Assume that the theorem holds for $k \geq 0$. From the recursion above, we have the following condition:

$$t_{k+1}(v) = \max_{u \in IN(v)} \{t_k(u)\} + \tau_{k+1}(v).$$

By the inductive hypothesis, we have the following:

$$t_{k+1}(v) = \max_{u \in IN(v)} \{\max T(S_k(u)) + \tau_{k+1}(v)\},$$

which gives the following desired result:

$$t_{k+1}(v) = \max T(S_{k+1}(v)). \quad \square$$

The Performance Measures: The most important performance measures investigated in this paper are the completion times $t_k(v), k \geq 0, v \in V$. A related performance measure of interest is the counting process $N_t^G(v)$ (or simply $N_t(v)$) associated with processor v , which is defined by the following expression:

$$N_t(v) \triangleq \sup \{k : t_k(v) \leq t\};$$

that is, $N_t(v)$ is the number of computation steps (minus 1) completed by v up to time t , or the highest index of an M_k message that has been sent by v up to time t . Similarly, $N_t \triangleq \sum_{v=1}^n N_t(v)$ denotes the total number of processing steps (minus n) executed in the network up to time t . The following claim indicates that no processor can advance (in terms of executed processing steps) too far ahead of any other processor.

Claim 2.2: Let d be the diameter of a directed, strongly connected graph G . Then, for all $u, v \in V$, and $t \geq 0, |N_t(u) - N_t(v)| \leq d$.

Proof: Denote by l the length of a simple path from u to v . A simple inductive argument on l shows that the fact that the last message sent by u up to time t is $M_{N_t(u)}$ implies that $N_t(v) \leq N_t(u) + l$. Thus, $N_t(v) - N_t(u) \leq l \leq d$. The same argument for a simple path from v to u proves that $N_t(u) - N_t(v) \leq d$. \square

Another important performance measure is the *computation rate*, $R^G(v)$ (or simply $R(v)$) of processor v in network G , defined by the following expression:

$$R(v) \triangleq \lim_{t \rightarrow \infty} \frac{N_t(v)}{t},$$

whenever the limit exists. Similarly, the computation rate of the network is defined by the following expression:

$$R \triangleq \lim_{t \rightarrow \infty} \frac{N_t}{t}.$$

Claim 2.2 implies that for every $u, v \in V$, $R(u) = R(v)$; therefore, $R = n \cdot R(v)$.

III. GENERAL PROBABILITY DISTRIBUTIONS

In this section, we compare the performance of different networks, with general distributions of the processing times $\tau_k(v)$. We first show that adding edges to a network with an arbitrary topology slows down the operation of each of the processors in the network. We show how the theory of graph embedding can be used to compare the rates of different networks. As an example, we present graphs, which have the same rate (up to a constant factor) for general distributions, although they have different topologies. Finally, we compare networks with the same (arbitrary) topology, but different distributions of the processing times. Specifically, we show that determinism maximizes the rate, and exponential distributions minimize the rate, among a large class of distributions.

A. Topology Of The Network

Monotonicity: Here we show that adding edges to a network with an arbitrary topology slows down the operation of each of the processors in the network. The basic methodology used is the *sample path* comparison; that is, we compare the evolution of message transmissions in different networks for every instance, or realization, of the random variables $\tau_k(v)$. This yields a stochastic ordering between various networks [33], [34].

Theorem 3.1: Let $G(V, E)$ be a graph, and let $E' \subseteq V \times V$ be a set of directed edges. Let $H(V, E \cup E')$ be the graph obtained from G by adding edges E' . Assume that processor v , $1 \leq v \leq n$ awakens in both G and H at the same time $t(v)$. For every realization of the random variables $\tau_k(v)$, $k \geq 0$, $1 \leq v \leq n$, the following inequalities hold:

$$t_k^G(v) \leq t_k^H(v),$$

for all $k \geq 0$, $1 \leq v \leq n$.

Proof: The proof is by induction on k . The basis of the induction is trivial, because the following is true:

$$t_0^G(v) = t(v) + \tau_0(v) = t_0^H(v).$$

The induction hypothesis is $t_k^G(v) \leq t_k^H(v)$. We need to show that $t_{k+1}^G(v) \leq t_{k+1}^H(v)$. From (1), we have the following condition:

$$\begin{aligned} t_{k+1}^G(v) &= \max_{u \in \text{IN}^G(v)} \{t_k^G(u)\} + \tau_{k+1}(v), \\ t_{k+1}^H(v) &= \max_{u \in \text{IN}^H(v)} \{t_k^H(u)\} + \tau_{k+1}(v). \end{aligned} \quad (2)$$

Since $\text{IN}^G(v) \subseteq \text{IN}^H(v)$, it follows that:

$$\max_{u \in \text{IN}^G(v)} \{t_k^G(u)\} \leq \max_{u \in \text{IN}^H(v)} \{t_k^H(u)\},$$

and therefore it follows from (2) that $t_{k+1}^G(v) \leq t_{k+1}^H(v)$, for all v . \square

The previous theorem implies immediately.

Corollary 3.2: Under the conditions of Theorem 3.1, we have the condition that $N_t^G(v) \geq N_t^H(v)$ and $R^G(v) \geq R^H(v)$ (when the limits exist) for all $v \in V$. Also $N_t^G \geq N_t^H$.

Remark 1: Notice that no assumption was made about the random variables $\tau_k(v)$. In particular, they need not be independent.

Remark 2: The sample path proof above implies that the random variable N_t^G is stochastically larger than the random variable N_t^H , denoted $N_t^G \geq_d N_t^H$; i.e., $\Pr\{N_t^G \geq \alpha\} \geq \Pr\{N_t^H \geq \alpha\}$ for all α .

Remark 3: The above implies that if one starts with a simple, directed cycle (a strongly connected graph with the least number of edges) and successively adds edges, a complete graph is obtained, without ever increasing the rate.

Embedding: The theory of graph embedding has been used to model the notion of one network simulating another on a general computational task (see for example [30]). Here we show how the notion of graph embedding can be helpful in comparing the behavior and the rates of different networks controlled by the synchronizer.

An *embedding* of graph G in graph H is specified by a one-to-one *assignment* $\alpha: V_G \rightarrow V_H$ of the nodes of G to the nodes of H , and a *routing* $\rho: E_G \rightarrow \text{Paths}(H)$ of each edge of G along a distinct path in H . The *dilation* of the embedding is the maximum amount by which the routing ρ "stretches" any edge of G :

$$\text{dilation}(\alpha, \rho) = \max_{u \rightarrow v \in E_G} \text{length}(\rho(u \rightarrow v)).$$

The dilation is a measure of the delay incurred by the simulation according to the embedding. The following theorem is a generalization of Theorem 3.1.

Theorem 3.3: Let (α, ρ) be an embedding with dilation D of a graph $G(V_G, E_G)$ in a graph $H(V_H, E_H)$. Assume that $t(v) = t(\alpha(v))$ for all $v \in V_G$, and that $\tau_k(v)$ and $\tau_{kD}(\alpha(v))$ for all $k \geq 0$, $v \in V_G$, have the same distribution. For every realization of the random variables $\tau_k^G(v) = \tau_{kD}^H(\alpha(v))$, $k \geq 0$, $v \in V_G$, the following inequalities hold:

$$t_k^G(v) \leq t_{kD}^H(\alpha(v)), \quad k \geq 0, v \in V_G.$$

Proof: For each path of length $k \geq 0$ in G , the following is true:

$$P_k^G = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k (= v),$$

one can use ρ to construct a path in H of length less than or equal to $k \cdot D$ from $\alpha(v_0)$ to $\alpha(v_k)$, passing through $\alpha(v_1), \alpha(v_2), \dots, \alpha(v_{k-1})$:

$$\rho(v_0 \rightarrow v_1) \rightarrow \rho(v_1 \rightarrow v_2) \rightarrow \dots \rightarrow \rho(v_{k-1} \rightarrow v_k).$$

Moreover, there is such a path of length exactly $k \cdot D$, because one can revisit vertices (each vertex has a self-loop) each time between a pair of vertices $\alpha(v_i)$ and $\alpha(v_{i+1})$, there are less than D edges in the path $\rho(v_i, v_{i+1})$. Thus, there is the following path in H :

$$P_{kD}^H = u_0 \rightarrow u_1 \rightarrow \dots \rightarrow u_{kD},$$

where $u_{iD} = \alpha(v_i)$, $0 \leq i \leq k$.

We are assuming a realization $\tau_k(u) = \tau_{kD}(\alpha(u))$, for every $u \in V_G$. It follows that for every path P_k^G , there exists a path P_{kD}^H , such that the following is true:

$$T(P_k^G) \leq T(P_{kD}^H),$$

and thus the following is also true:

$$\max T(S_k^G(v)) \leq \max T(S_{kD}^H(\alpha(v))).$$

By Theorem 2.1, $t_k^G(v) \leq t_{kD}^H(\alpha(v))$. □

Corollary 3.4: Under the conditions of Theorem 3.3, we have the condition that $D \cdot N_t^G(v) \geq N_t^H(\alpha(v))$ and $D \cdot R^G(v) \geq R^H(\alpha(v))$ (when the limits exist) for all $v \in V_G$.

Remarks 1-3 hold in this case, too.

A simple corollary of Theorem 3.3 is that if G is a subgraph of H , $N_t^G(v) \geq_d N_t^H(\alpha(v))$. This is because if G is a subgraph of H , then there is an embedding from G in H with dilation 1. In addition, if the number of vertices in G and H are equal, and the dilation of the embedding is D , then G is a D -spanner of H (e.g., [28]), and we have the following.

Corollary 3.5: If H has a D -spanner G , then $R^G/D \leq R^H \leq R^G$.

A motivation for the theory of embedding is simulation. Namely, one expects that if there is an embedding (α, ρ) from G in H with dilation D , then the architecture H can simulate T steps of the architecture G on a general computation in order of $D \cdot T$ steps, by routing messages according to ρ . In our approach, we compare the performance of G and of H under the synchronizer, without using ρ . The embedding is used only for the purpose of proving statements about the performance of the networks. Consider, for example, the following two results of the theory of embedding [30].

Proposition 3.6: For all $n \geq 1$, one can embed the order n shuffle-exchange graph in the order n deBruijn graph with dilation 2. One can embed the order n deBruijn graph in the order n shuffle-exchange graph with dilation 2.

Proposition 3.7: For all $n \geq 1$, one can embed the order n cube-connected-cycles graph in the order n butterfly graph with dilation 2. One can embed the order n butterfly graph in the order n Cube-Connected-Cycles graph with dilation 2.

By Theorem 3.3, the average rate of the graphs of Proposition 3.6 (3.7) are equal up to a constant factor of 2, provided that the processing times of corresponding processors have the same distributions (regardless of what these distributions are).

B. Probability Distributions

Deterministic Processing Times: Now we compare networks, say, $G(V, E)$ and $H(V, E)$, having the same (arbitrary) topology, but operate with different distributions of the random variables $\tau_k(v)$. To that end, we assume that the processing times $\tau_k^G(v)$, $k \geq 0$, $v \in V$ are independent and have finite mean $E[\tau_k^G(v)] = \lambda_v^{-1}$.

We say that λ_v is the *potential rate* of v , because this would be the rate of v if it would not have to wait for messages from its in-neighbors. The processing times in H are distributed as in G , except for a subset $V' \subseteq V$ of processors, for which the processing times are assumed to be deterministic; i.e., $\tau_k^H(v) = \lambda_v^{-1}$, $v \in V'$, for $k \geq 0$. We let

$\tau_k^H(v) = \tau_k^G(v) = \tau_k(v)$, $k \geq 0$, $v \notin V'$, be any specific realization of the random variables in G . Again, it is assumed that the processors are awakened at the same time in both networks.

Theorem 3.8: Under the above conditions, we have the following condition:

$$t_k^H(v) \leq E[t_k^G(v)],$$

for all processors v , and $k \geq 0$. The expectation is taken over the respective distributions of processing times of processors of G in V' .

Proof: The proof is by induction on k . For the basis, $k = 0$, we make the following observations:

$$E[t_0^G(v)] = t_0^G(v) = t(v) + \tau_0(v) = t_0^H(v),$$

for $v \notin V'$, and

$$E[t_0^G(v)] = t(v) + \lambda_v^{-1} = t_0^H(v),$$

for $v \in V'$.

The induction hypothesis is $t_k^H(v) \leq E[t_k^G(v)]$, and we need to show that $t_{k+1}^H(v) \leq E[t_{k+1}^G(v)]$, for all $v \in V$.

From (1), we have the following condition:

$$t_{k+1}^G(v) = \max_{u \in \text{IN}(v)} \{t_k^G(u)\} + \tau_{k+1}^G(v),$$

for $v \in V$. Jensen's inequality implies the following:

$$E[t_{k+1}^G(v)] \geq \max_{u \in \text{IN}(v)} \{E[t_k^G(u)]\} + E[\tau_{k+1}^G(v)].$$

By the induction hypothesis, we have the following equation:

$$E[t_{k+1}^G(v)] \geq \max_{u \in \text{IN}(v)} \{t_k^H(u)\} + E[\tau_{k+1}^G(v)] = t_{k+1}^H(v),$$

because $E[\tau_{k+1}^G(v)] = \tau_{k+1}(v)$ for $v \notin V'$, and $E[\tau_{k+1}^G(v)] = \lambda_v^{-1} = \tau_{k+1}^H(v)$, for $v \in V'$. □

Remark 4: Theorem 3.8 holds also if the processing times $\tau_k^H(v)$ of processors v of H in V' , are deterministic, but not necessarily the same for every k .

When all processing times in the network H are deterministic, the computation of the network rate is no longer a stochastic problem, but a combinatorial one. Thus, a conclusion of Theorem 3.8 is that in this case, the computation rate of H , obtained via combinatorial techniques ([18] and [19]), yields an upper bound on the average rate of G . Furthermore, if the times $t_k^H(v)$ are computed, they give a lower bound on $E[t_k^G(v)]$ for every $k \geq 0$.

More Variable Processing Times: More generally, we study the effect of substituting a random variable in the network (e.g., the processing time of a given processor, for a given computational step) with a given distribution, for a random variable with another distribution on the rate of the network, and define an ordering among probability distributions.

Recall that a function h is *convex* if, for all $0 < t < 1$, x_1, x_2 , $h(tx_1 + (1-t)x_2) \leq th(x_1) + (1-t)h(x_2)$. A random variable X with distribution F_X is said to be *more variable* than a random variable Y with distribution F_Y , denoted $X \geq_c Y$ or $F_X \geq_c F_Y$, if $E[h(X)] \geq E[h(Y)]$ for all increasing

convex functions h . The partial order \leq_c is called *convex* order (e.g., [33], [34]). Intuitively, X will be more variable than Y if F_X gives more weight to the extreme values than F_Y . For instance, if $E[X] = E[Y]$, then $\text{Var}(X) \geq \text{Var}(Y)$, because $h(x) = x^2$ is a convex function.

Here we compare networks, say, $G(V, E)$ and $H(V, E)$ having the same arbitrary topology, but some of the processing times in G are more variable than the corresponding processing times in H , i.e., for some k 's and some v 's, $\tau_k^G(v) \geq_c \tau_k^H(v)$, while all other processing times have the same distributions in both graphs. When $t^G(v) = t^H(v)$ and all processing times in G (H) are independent of each other, the following holds.

Theorem 3.9: Under the above conditions, the following holds for all processors v , and $k \geq 0$:

$$t_k^H(v) \leq_c t_k^G(v).$$

Proof: From Theorem 2.1, we have the following condition:

$$t_k(v) = \max\{T(P_k) : P_k \in S_k(v)\},$$

where $P_k = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k (= v)$ is a directed path of length k ending in v , and $T(P_k) = t(v_0) + \sum_{i=0}^k \tau_i(v_i)$.

From the fact that the τ 's are positive and \max and \sum are convex increasing functions, it follows that $t_k(v)$ is a convex increasing function of its arguments $\{\tau_i(u) : 0 \leq i \leq k, u \in P_k, P_k \in S_k(v)\}$. Now we can use Proposition 8.5.4 in [33].

Proposition 8.5.4: If X_1, X_2, \dots, X_n are independent r.v., and Y_1, Y_2, \dots, Y_n are independent r.v., and $X_i \geq_c Y_i, i = 1, 2, \dots, n$, then $g(X_1, X_2, \dots, X_n) \geq_c g(Y_1, Y_2, \dots, Y_n)$ for all increasing convex function g that are convex in each of its arguments.

The proof of the theorem now follows, because, by assumption, the τ 's in G are independent, the τ 's in H are independent, and $\tau_k^H(v) \leq_c \tau_k^G(v), k \geq 0, v \in V$. Note that the random variables $T(P_k)$ are not independent.

Corollary 3.10: Under the above conditions $N_t^G(v) \leq N_t^H(v), R^G(v) \leq R^H(v)$ and $R^G \leq R^H$.

In the next section, we show that if the processing times are independent and have the same exponential distribution with mean λ^{-1} , then the rate of any network is at least $\lambda|V|/\log|V|$. We conclude this subsection by characterizing a set of distributions for which the same lower bound holds.

Assume that the expected time until a processor finishes a processing step, given that it has already been working on that step for α time units, is less than or equal to the original expected processing time for that step. Namely, we assume that the distributions of the processing times $\tau_k(v)$, for all $v \in V, k \geq 0$, are *new better than used in expectation* (NBUE) (e.g., [33], [34]), so that if τ is a processing time, then we have the following equation:

$$E[\tau - a | \tau > a] \leq E[\tau], \quad \forall a \geq 0.$$

Let $G_d(V, E)$ be a network with deterministic processing times; let $G_e(V, E)$ be a network with corresponding processing times with the same mean, but independent, exponentially distributed; and let $G(V, E)$ be a network with corresponding processing times with the same mean and independent, but

with any NBUE distribution. The following theorem follows from the fact that the deterministic distribution is the minimum, whereas the exponential distribution is the maximum with respect to the ordering \leq_c , among all NBUE distributions [33], [34].

Theorem 3.11: For every $v \in V, k \geq 0$, it holds that $t_k^{G_d}(v) \leq_c t_k^G(v) \leq_c t_k^{G_e}(v)$.

Some examples of distributions that are less variable than the exponential (with appropriate parameters) are the Gamma, Weibull, and Uniform. We should conclude this section by pointing out that the interested reader can find similar results for rather general stochastic Petri nets in [5] and [8].

IV. EXPONENTIAL DISTRIBUTIONS

In this section, we assume that the processing times $\tau_k(v), k \geq 0, v \in V$ are independent and exponentially distributed with mean λ^{-1} . We first consider general topologies and derive upper and lower bounds on the expected values of $t_k(v)$, and thus obtain upper and lower bounds on the rate of the network. These bounds depend on the in-degrees and out-degrees of processors in the network, but not on the number of processors itself. Then, exploring the Markov chain of the underlying process, we derive the exact rates of two extreme topologies: the directed ring and the fully connected (complete) network. For these two topologies, we study also the effect of having a single slower processor within the network.

A. Upper and Lower Bounds

Denote by $d_{\text{out}}(v)(d_{\text{in}}(v))$ the number of edges going out of (into) v in G , and let the following be true:

$$\Delta_{\text{out}} = \max_{v \in V} d_{\text{out}}(v), \quad \Delta_{\text{in}} = \max_{v \in V} d_{\text{in}}(v);$$

$$\delta_{\text{out}} = \min_{v \in V} d_{\text{out}}(v), \quad \delta_{\text{in}} = \min_{v \in V} d_{\text{in}}(v).$$

Lemma 4.1 (Lower Bound):

- 1) For every $k \geq 0$, there exists a processor $v \in V$ for which the following condition exists:

$$E[t_k(v)] \geq \max_{u \in V} t(u) + \lambda^{-1}[1 + k \cdot \log \delta_{\text{out}}].$$

- 2) For every $k \geq 0$, and every $v \in V$, the following holds:

$$E[t_k(v)] \geq \min_{u \in V} t(u) + \lambda^{-1}[1 + k \cdot \log \delta_{\text{in}}].$$

Proof: We present a detailed proof for part (1) only; the proof of part (2) is discussed at the end of this paper. We start by proving that for every $k \geq 0$, there exists a (not necessarily simple) path $v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k$, such that the following expression is true:

$$E[t_{i+1}(v_{i+1})] - E[t_i(v_i)] \geq \lambda^{-1} \log \delta_{\text{out}}, \quad 0 \leq i < k.$$

We assume the statement holds for $k \geq 0$, and prove it for $k + 1$. The proof of the basis is identical. Let v_{k+1} be the processor for which the processing time during the $(k + 1)$ th

computational step is maximum, among the out-neighbors of v_k , i.e., as follows:

$$\tau_{k+1}(v_{k+1}) = \max_{v_k \rightarrow v} \tau_{k+1}(v).$$

Because v_{k+1} will not start the $(k + 1)$ st computational step before v_k finishes the k th computational step, we have the condition that $t_{k+1}(v_{k+1}) - t_k(v_k) \geq \tau_{k+1}(v_{k+1})$. The quantity $\tau_{k+1}(v_{k+1})$ is equal to the maximum of at least δ_{out} independent and identically distributed exponential random variables with mean λ^{-1} . It is well known (e.g., [13], [17]) that the mean of the maximum of c such random variables is at least $\lambda^{-1} \log c$. It follows that:

$$E[t_{k+1}(v_{k+1})] - E[t_k(v_k)] \geq \lambda^{-1} \log \delta_{\text{out}}.$$

We can choose v_0 to be the one with latest waking time $t(v_0)$, and thus $E[t_0(v_0)] = t(v_0) + \lambda^{-1}$. Therefore, for every $k \geq 0$, there exists a processor v such that the following is true:

$$E[t_k(v)] \geq \max_{u \in V} t(u) + \lambda^{-1} [1 + k \cdot \log \delta_{\text{out}}],$$

completing the proof of (1). The proof of part (2) evolves along the same lines, except that we start from v_k and move backward along a path. \square

Remark 5: From its proof, one can see that Lemma 4.1 holds for any distribution F of the processing times, for which the expected value m_c of the maximum of c independent r.v. with distribution F exists. In this case, it implies that $R_v \leq 1/m_c$, with $c = \delta_{\text{out}}$ or $c = \delta_{\text{in}}$.

Remark 6: Lemma 4.1 implies that for the exponential case, the slowdown of the rate is at least logarithmic in the maximum degree of G . By Remark 5, there are distributions (not NBUE by Theorem 3.11), for which the slowdown is larger; an example is $F(x) = 1 - 1/x^2$, $x \geq 1$, for which the slowdown is at least the square root of the maximum degree of G [17, p. 58].

Lemma 4.2 (Upper Bound):

1) For every $k \geq 1$, for every processor v , we get the following:

$$E[t_{k-1}(v)] \leq \max_{u \in V} t(u) + \frac{4}{\lambda} (1 + k \cdot \log \Delta_{\text{in}}).$$

2) For every $k \geq 1$, for every processor v , we get the following:

$$E[t_{k-1}(v)] \leq \max_{u \in V} t(u) + \log |V| + \frac{4}{\lambda} (1 + k \cdot \log \Delta_{\text{out}}).$$

Proof: Again, we restrict ourselves to the proof of part (1). Recall that Theorem 2.1 states that for every $v \in V$, $k \geq 0$, $t_k(v) = \max T(S_k(v))$. Also, for a path $P_k = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k$, $T(P_k) = t(v_0) + \sum_{i=0}^k \tau_i(v_i)$, but for the moment, let $t(v) = 0$ for every v .

By Proposition D.2 of the Appendix, we have the following condition:

$$\Pr \left(T(P_{k-1}) \geq \frac{ck}{\lambda} \log \Delta_{\text{in}} \right) \leq e^{-\frac{ck}{4} \log \Delta_{\text{in}}},$$

for every $c > 4$, because $\log 2 / \log \Delta_{\text{in}} \leq 1$. It follows that:

$$\begin{aligned} \Pr \left(t_{k-1}(v) \geq \frac{ck}{\lambda} \cdot \log \Delta_{\text{in}} \right) &\leq \Delta_{\text{in}}^k \cdot e^{-\frac{ck}{4} \log \Delta_{\text{in}}} \\ &= e^{-k(\frac{c}{4}-1) \log \Delta_{\text{in}}}, \end{aligned}$$

for every $c \geq 4$, and

$$\begin{aligned} E[t_{k-1}(v)] &\leq \int_0^{\frac{4k}{\lambda} \log \Delta_{\text{in}}} 1 dt + \int_{\frac{4k}{\lambda} \log \Delta_{\text{in}}}^{\infty} e^{-k(\frac{c}{4}-1) \log \Delta_{\text{in}}} \frac{4}{\lambda} \log \Delta_{\text{in}} dc \\ &= \frac{4k}{\lambda} \log \Delta_{\text{in}} + \frac{4}{\lambda}. \end{aligned} \quad \square$$

Combining Lemma 4.1 and Lemma 4.2, we obtain Theorem 4.3.

Theorem 4.3:

$$\frac{\lambda}{4 \log \min(\Delta_{\text{out}}, \Delta_{\text{in}})} \leq R(v) \leq \frac{\lambda}{\log \max(\delta_{\text{out}}, \delta_{\text{in}})}.$$

B. Exact Computations

Theorem 4.3 implies the following bounds for the rate of a directed cycle $C_n(\Delta = \delta = 2)$ and of a complete graph $K_n(\Delta = \delta = n)$, where n is the number of processors:

$$0.36\lambda \leq R^{C_n}(v) \leq \lambda,$$

$$\frac{\lambda}{4 \log n} \leq R^{K_n}(v) \leq \frac{\lambda}{\log n}.$$

In this section, we shall compute the exact values for the rates of C_n and K_n . To that end, we consider the Markov chain associated with the network. This Markov chain is denoted by $X(t) = (X_1(t), X_2(t), \dots, X_m(t))$, where $X_i(t)$ is the number of messages stored in the buffer of edge i at time t , and m is the number of edges in the network. Note that a processor with a positive number of messages on each of its incoming edges is in a processing state. When such a processor completes its processing (after an exponential time), one message is deleted from each of its incoming edges and one message is put on each of its outgoing edges. We denote by s_0 the state in which $X_i(0) = 1$, $1 \leq i \leq m$. Thus, the network can be represented as a marked graph (e.g., [15]).

The number of states in the Markov chain is finite, say, N , because a transition of the chain does not change the total number of messages in a circuit in the network. Moreover, if the network is strongly connected, then the Markov chain is irreducible. Therefore, the limiting probabilities P_i , $1 \leq i \leq N$, of the states s_i of the chain exist; they are all positive; and their sum is equal to 1 (e.g., [14], [33]). However, as we shall see, N can be exponential in n ; therefore, it is infeasible to compute the rate by directly solving the Markov chain. Here we show how to solve the Markov chain for two network classes without having to produce the entire chain. We hope this combinatorial approach could be applied to other networks as well.

Let G_X denote the transition diagram (directed graph) of the Markov chain X . Consider a breadth-first search (BFS) tree of G_X , rooted at s_0 . The level $L(v)$ of a vertex v will be

equal to the distance from s_0 to v . Thus, $L(s_0) = 0$. Denote by L_i , $i \geq 0$, the set of vertices at level i , and by L the number of levels of G_X .

A Simple Directed Cycle: We study the performance of a simple directed cycle of n processors $C_n = p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n \rightarrow p_1$. It is not difficult to observe that the Markov chain associated with C_n corresponds to that of a closed queuing network. We return to this approach later. Here we choose to use a combinatorial approach.

Theorem 4.4:

- 1) All the states associated with C_n have the same limiting probability.
- 2) For any graph G that is not a simple directed cycle, (1) does not hold.

Proof: 1) The proof follows from two observations. First, by symmetry, all the states in one level have the same probability. Second, the in-degree of any state in the transition diagram is equal to its out-degree. Then a simple inductive argument can be used to prove part (1).

2) If G is not a cycle, then it has a node v , s.t. $d_{in}(v) \geq 1$. Let v_1, v_2 be two nodes with edges to v . Consider the state s , reached from s_0 , by the processing completion (or, in marked graphs terminology, firing) of vertex v . The outdegree of $s = n - 1$, because apart from v , all vertices are still enabled. But the indegree of s is at most $n - 2$, because by the firing of v_1 or of v_2 , it is not possible to reach s , because there are no messages on the edges from v_1 and v_2 to v , in s . Therefore, we have proved that $d_{in}(s) \neq d_{out}(s)$.

Consider the balance equation that holds at state s : $P_1 + P_2 + \dots + P_k = n - 1 \cdot P_s$, where P_i , $1 \leq i \leq k$ are the limiting probabilities of the states that have an edge to s , $k = d_{in}(s)$, P_s is the limiting probability of s , and $n - 1 = d_{out}(s)$. We have just proved that $k \neq n - 1$. It follows that it is not possible that all the probabilities of the last equation are equal. \square

The next theorem states that each processor of C_n works at least at half of its potential rate λ , regardless of the value of n .

Theorem 4.5: The rate $R(v)$ of a processor in C_n is as follows:

$$R(v) = \lambda \cdot \left(1 - \frac{n-1}{2n-1}\right) \sim \frac{\lambda}{2},$$

$$N = \frac{(2n-1)!}{n!(n-1)!},$$

and the limiting probability of each state is $1/N$, where N is the number of states in the associated chain.

Proof: If M is the number of states in which at least one message is in an edge, going into a processor, say, v , then the running rate will be M/N times the expected firing rate. This is because v will be enabled when it has more than 0 messages in its input edge, and because all states have the same probability (Theorem 4.4), the percentage of the time that is enabled is simply M/N .

The number of ways of putting n objects in k places is as follows:

$$P(n, k) = \frac{(n+k-1)!}{n!(k-1)!}.$$

It is not difficult to see that $N = P(n, n)$ and $M = N - P(n, n-1)$. Thus, we have the following equation:

$$\frac{M}{N} = 1 - \frac{n-1}{2n-1},$$

which gives the desired results. \square

A Complete Graph: Let K_n be a complete graph with n processors. Recall that N is the number of states in the associated Markov chain, and let s_0 be the state in which each edge has one token. A state is at level l , $0 \leq l \leq n-1$, if it can be reached from s_0 by the firing of l processors. The limiting probability of a state at level l is denoted by $P(l)$.

Theorem 4.6: The rate of a processor in K_n is as follows:

$$R(v) = \lambda / \sum_{i=1}^n \frac{1}{i} \sim \frac{\lambda}{\log n}.$$

Proof: A simpler proof can be derived, as in the proof of Theorem 4.10; here we give a combinatorial proof that also yields the number and the limiting probabilities of the states of the associated Markov chain.

We consider a Markov chain T , similar to the Markov chain associated with network K_n . The root of T , s_0 , is the state with a message in each edge. A state s will have one son for each one of the enabled processors at state s ; a son of s corresponds to the state arrived at from s by the firing (completion of a processing step) of one of the enabled processors in state s . Note that in chain T , there are several vertices corresponding to the same state of the chain associated with K_n .

In T , the number of states in level l is $n!/(n-l)!$, because each time a processor fires, it cannot fire again until the rest of the processors have fired. Thus, the number N^T , of states in T is as follows:

$$N^T = \sum_{i=1}^n \frac{n!}{i!}.$$

The number of states in which a given processor is enabled at level l , $en(l)$ (edges from level l to level $l+1$), is as follows:

$$en(l) = \frac{1}{n} \frac{n!}{(n-l-1)!},$$

because at level l , there are $n!/(n-l-1)!$ enabled processors, and, by symmetry, each processor is enabled the same number of times at each level.

Let us denote by P_l^T the limiting probability of a state of T in level l . One can show that $P_l^T = (n-l-1)!/K$, where the following is true:

$$K = \sum_{l=0}^{n-1} \frac{n!}{(n-l)!} P_l^T = n! \sum_{i=1}^n \frac{1}{i}.$$

It follows that the percentage of time that a processor is enabled is as follows:

$$ut = \sum_{l=0}^{n-1} ut(l) = \frac{1}{\sum_{i=1}^n \frac{1}{i}},$$

where $ut(l) = en(l)P_l^T$, and its rate is $\lambda \cdot ut$. \square

Corollary 4.7: For a network K_n , we have the following conditions:

$$N = 2^n - 1,$$

$$P_l = \frac{l!(n-l-1)!}{n! \sum_{i=1}^n \frac{1}{i}}.$$

Proof: As noted before, it may be that two states of T correspond to the same state, say, s , of K_n . In fact, if a state of T is reached from s_0 by firing a sequence of processors of length k , then all $k!$ permutations of the processors in this sequence constitute a valid firing sequence, which leads to the same state s . Thus, the limiting probability of a state s at level l is as follows:

$$P_l = l!P_l^T = \frac{l!(n-l-1)!}{n! \sum_{i=1}^n \frac{1}{i}}.$$

The number of different states at level l is $n!/l!(n-l)!$, and the total number of different states is as follows:

$$\sum_{l=0}^{n-1} n!/l!(n-l)! = 2^n - 1.$$

□

Corollary 4.8: Asymptotically, the rate of any network of n processors is between $\lambda n/2$ and $\lambda n/\log n$.

Observe that the best possible rate of a processor is $2/3$ of the potential rate, in the case of a cycle of two processors. Adding more processors can only lower this rate, but not below $1/2$. Yet the rate of the network grows linearly with n . In the case of a complete graph, the rate of a processor reduces as n grows, but also here the total number of computational steps executed per unit time ($n/\log n$) grows with n .

C. Bottlenecks

Suppose that the potential rate of all processors of a graph is λ , except for one, which has a lower rate μ . We shall now show that such a bottleneck has a stronger effect in a network that is a directed cycle, than in one that is a complete graph.

Consider the case of a simple directed cycle with n vertices CB_n , where $n-1$ processors have rate λ and one processor has rate μ . Using standard techniques of queuing theory, we prove the following.

Theorem 4.9: The following is the rate of a processor in CB_n :

$$\mu \left(1 - \frac{\binom{2n-2}{n} \rho^n}{\sum_{i=0}^n \binom{n+i-2}{i} \rho^i} \right); \quad \rho = \frac{\lambda}{\mu}.$$

Proof: Let X_i , $1 \leq i \leq n$ be the number of messages in the buffer of the incoming edge to processor i . The total number of messages in the cycle is equal to n . Since this is a closed queueing system, we have the condition that the limiting probability of the system being in state (x_1, x_2, \dots, x_n) is

given by the following product form [33]:

$$\Pr(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = K \left(\frac{1}{\mu} \right)^{x_n} \prod_{j=1}^{n-1} \left(\frac{1}{\lambda} \right)^{x_j}$$

if $\sum_{j=1}^n x_j = n$,

and is equal to 0 otherwise, where K is a normalization constant that guarantees that the sum of all the above probabilities is equal to 1. Thus, the following is the probability of having l messages on the incoming edge to processor n :

$$\Pr(X_n = l) = \Pr(X_1 + X_2 + \dots + X_{n-1} = n-l, X_n = l)$$

$$= \sum_{S'} K \left(\frac{1}{\mu} \right)^l \left(\frac{1}{\lambda} \right)^{n-l}$$

$$= |S'| K \left(\frac{1}{\lambda} \right)^n \rho^l; \quad 0 \leq l \leq n,$$

where $\rho = \lambda/\mu$ and $S' = \{(x_1, x_2, \dots, x_{n-1}) : x_j \geq 0, 1 \leq j \leq n-1, \sum_{j=1}^{n-1} x_j = n-l\}$. Hence,

$$\Pr(X_n = l) = K' \binom{2n-l-2}{n-l} \rho^l, \quad 0 \leq l \leq n,$$

where K' is the normalization factor determined by the condition $\sum_{l=0}^n \Pr(X_n = l) = 1$. Finally, we get the following:

$$\Pr(X_n = l) = \frac{\binom{2n-l-2}{n-l} \rho^l}{\sum_{j=0}^n \binom{2n-j-2}{n-j} \rho^j}, \quad 0 \leq l \leq n.$$

Now observe that the rate is simply $\mu n[1 - \Pr(X_n = 0)]$, because the rate of the processor is μ while there are messages in the incoming edge to processor n . □

Several conclusions can be derived from Theorem 4.9. First, observe that the rate of the cycle cannot exceed μn ; thus, the slow processor bounds the rate of the network. Moreover, for a fixed n and a very slow processor ($\mu \rightarrow 0$ or $\rho \rightarrow \infty$), the rate of the network is $\mu n \left[1 - \binom{2n-2}{n} \rho^{-n} \right]$; namely, as ρ increases, the rate approaches its upper bound μn .

Next we consider the case where the graph is a complete graph KB_n . We continue to assume that the rate of $n-1$ processors is λ and the n th processor is slower, operating at rate μ . We shall show that for fixed μ and λ , as the number of processors n grows to infinity, the influence of the slow processor diminishes, and in the limit, the rate of the network is the same as that of a network with all processors running with the same rate λ .

Theorem 4.10: The rate of a processor in KB_n is at least $\frac{\lambda}{\rho + \log n}$, $\rho = \lambda/\mu$.

Proof: Suppose that the network is in state s_0 at a given time, and that after some time T_1 , it returns to that state. Then, after some time T_2 , it returns again, and so on. Then, $\{T_i, i = 1, 2, \dots\}$ is a sequence of non-negative independent random variables with a common distribution F , and expected value $E[T_i]$.

Denote by $N(t)$ the number of events (returning to s_0) by time t . The counting process $\{N(t), t \geq 0\}$ is a renewal process. Therefore, with probability 1, we have the following equation:

$$\frac{N(t)}{t} \rightarrow \frac{1}{E[T_i]} \quad \text{as } t \rightarrow \infty.$$

(See, e.g., [33].) Moreover, since each time the process returns to s_0 , each processor of the network has completed exactly one computational step, it follows that the rate of the network is $1/E[T_i]$. We proceed to bound $E[T_i]$.

The expected time of T_i that it takes to return to s_0 is of the following form:

$$\begin{aligned} \alpha_j = & \frac{1}{(n-1)\lambda + \mu} + \frac{1}{(n-2)\lambda + \mu} + \dots \\ & + \frac{1}{(n-j)\lambda + \mu} + \frac{1}{(n-j)\lambda} + \frac{1}{(n-j-1)\lambda} + \dots \\ & + \frac{1}{\lambda}, \end{aligned}$$

for some $1 \leq j \leq n$, depending on when the slow processor completes a computational step. If the system leaves s_0 because the slow processor completed a computational step, $E[T_i]$ is α_1 . In general, if the j th ($1 \leq j \leq n$) processor to complete a computational step, after leaving s_0 , is the slow processor, then $E[T_i]$ is α_j .

The probability of $E[T_j]$ being equal to α_j is not necessarily the same for every j , but for the case $\lambda \geq \mu$, it holds that $\alpha_j < \alpha_{j+1}$. Thus, the following equation:

$$\alpha_n = \sum_{k=0}^{n-1} \frac{1}{\lambda \cdot k + \mu}$$

gives an upper bound on the time $E[T]$ that it takes to return to s_0 , and $1/\alpha_n$, is a lower bound on the rate of a processor in the network.

We have the following condition:

$$\alpha_n = \sum_{i=0}^{n-1} 1/(\lambda i + \mu) \leq \frac{1}{\mu} + \frac{1}{\lambda} \sum_{i=1}^{n-1} \frac{1}{i} \leq \frac{1}{\mu} + \frac{1}{\lambda} \log n.$$

We see that $E[T_i] = O(\frac{1}{\mu} + \frac{1}{\lambda} \log n)$, and thus that the rate is at least $\frac{1}{\frac{1}{\mu} + \frac{1}{\lambda} \log n}$. \square

For fixed μ , R_v is $\Theta(\lambda/\log n)$; but observe that the rate of the network cannot exceed μ . However, when the number of processors n increases to infinity, the rate of the network decreases in proportion to $1/\log n$, as if the slower processor were not in the network.

V. NON-NEGLIGIBLE TRANSMISSION DELAYS

We briefly discuss the case of non-negligible transmission delays. In this model, the processing times are random, as before, but the transmission delays are also random. Denote the transmission delay of message M_k , $k \geq 0$, along edge $u \rightarrow v$, by $\tau_k(u, v)$, and let $\tau_k(u, u) = 0$, for all $u \in V$. It follows that the behavior of the system is described by the following recursions:

$$\begin{aligned} t_0(v) &= t(v) + \tau_0(v) \\ t_{k+1}(v) &= \max_{u \in \text{IN}(v)} \{t_k(u) + \tau_k(u, v)\} + \tau_{k+1}(v) \quad k \geq 0. \end{aligned}$$

Note that this system is not equal to the one of [13], in which the processing times are negligible, and the delays are non-negligible, with a self-loop in each processor (to model its processing delay).

Let $P_k = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k (= v)$ be a path of length k . It is easy to see how to modify the definition of $T(P_k)$:

$$T(P_k) \triangleq t(v_0) + \sum_{i=0}^{k-1} [\tau_i(v_i) + \tau_i(v_i, v_{i+1})] + \tau_k(v_k).$$

Thus, a theorem similar to Theorem 2.1 holds, and the corresponding results for general distributions follow.

Consider the case in which the processing times, as well as the transmission delays, are exponentially distributed, with the same mean, say, 1. It is easy to see that Lemma 4.1 still holds, and that Lemma 4.2 holds up to a factor of 2. Namely, by Theorem 3.9, a regular network with non-negligible delays runs at the same rate as the same network, up to a constant factor, provided that the delays are less than or equal (in the convex order) to the processing times. In [18], we show that for a network with negligible delays and deterministic processing times equal to 1, the rate of any network is equal to 1. Thus, in this case, random processing times degrade the rate by at most a logarithmic factor in the maximum degree of a processor.

Now consider the case in which all processing times have a mean of 1, but the delays have mean λ^{-1} greater than 1, both exponentially distributed. The rate in the deterministic case is equal to λ [19]; thus, by Theorem 3.3, in our case, the rate is at most λ . One can prove (by using Proposition D.2), that also in the case of non-negligible delays, the rate is degraded by at most a logarithmic factor in the maximum degree of a processor, with respect to the (optimal) deterministic case, for any NBUE distribution.

As for the exact computations for networks with average processing times and delays exponentially distributed with mean 1, the rate of a simple cycle can be computed by using the same tools of queuing theory that we used in the case of negligible delays. Computing the rate of a complete network K_n is not as straightforward; the structure of the Markov process is more complicated, but by the arguments made above, we have the condition that the rate is between $1/8 \log n$ and $1/\log n$. However, using the ideas of embedding, let us show that the rate of K_n is at least $1/4 \log n$. Let K'_n be a complete network with negligible delays. Construct G , from K'_n by inserting one vertex in each of its edges. By Theorem 3.1 (or by Theorem 3.9), the rate of any processor in G is at

least the following:

$$\frac{1}{\log(n + n(n - 1))} = \frac{1}{2 \log n}.$$

One can show that the rate of any processor v in K_n is greater than or equal to half the rate of the corresponding processor in G using the fact that there is an embedding of K_n into G of dilation 2. Therefore, the rate of v in K_n is at least $1/4 \log n$.

VI. CONCLUSION

In this paper, we have studied the behavior of synchronizers in networks with random transmission delays and processing times. We attempted to present a self-contained general study of the synchronizer performance from the view point of distributed algorithms, rather than providing a deep mathematical study of the underlying stochastic process. In particular, we were interested in comparing the behavior of synchronizers with random delays, as opposed to the usual approach of analyzing distributed algorithms with bounded delays. Our main conclusion is that if the delays belong to the natural class of NBUE distributions, the rate of the network is degraded only by a small, local (vertex degree) factor.

We presented several properties of the behavior of the synchronizer for general probability distributions, and described techniques useful to compare the rate of the synchronizer running in networks with different topologies. For exponential distributions, we showed that the expected duration of a round of computation depends on the logarithm of a vertex degree, and hence that the rate of computation does not diminish with the number of processors in the network. We presented techniques to prove upper and lower bounds on the rate, and to obtain exact computations. We hope that the combinatorial approach of these techniques, which was applied to rings, complete networks, and regular degree networks, will be used in the future to obtain results for other topologies as well.

APPENDIX

The following proposition (similar to pp. 672 in [13]) is used to prove the lower bounds on the rate of a network.

Proposition 7.1 (D.2): Let $\{X_i\}$ be a sequence of independent exponential random variables with mean λ^{-1} . For every positive integer k and any $c \geq 4 \log 2$, we have the following condition:

$$\Pr \left(\sum_{i=1}^k X_i \geq \frac{ck}{\lambda} \right) \leq e^{-\frac{c}{4}}.$$

Proof: Fix $\beta \in (0, \lambda)$, and let γ be a positive scalar. A direct calculation yields the following equation:

$$\mathbb{E} \left[e^{\beta(X_i - \gamma)} \right] = \int_0^\infty e^{\beta(x - \gamma)} \lambda e^{-\lambda x} dx = e^{-\beta\gamma} \frac{\lambda}{\lambda - \beta}.$$

In particular, we can choose γ sufficiently large so that $\mathbb{E} \left[e^{\beta(X_i - \gamma)} \right] \leq 1$. If $\beta = \lambda/2$, then $\gamma \geq 2\lambda^{-1} \log 2$ satisfies the last equation. Using the independence of the random

variables X_i , we obtain the following:

$$\mathbb{E} \left[e^{\beta \sum_{i=1}^k (X_i - \gamma)} \right] = \prod_{i=1}^k \mathbb{E} \left[e^{\beta(X_i - \gamma)} \right] \leq 1.$$

Using the Markov inequality, we obtain the following:

$$e^{\beta km} \Pr \left(e^{\beta \sum_{i=1}^k (X_i - \gamma)} \geq e^{\beta km} \right) \leq 1.$$

This in turn implies the following equation:

$$\Pr \left(\sum_{i=1}^k X_i \geq mk + \gamma k \right) \leq e^{-\beta mk},$$

$$\Pr \left(\sum_{i=1}^k X_i \geq k(m + \gamma) \right) \leq e^{-\beta mk}.$$

Let $c = 2m$ and $c_0 = 2\gamma$. Then, if $c \geq c_0$, we have $m \geq \gamma$ and the following equation:

$$\Pr \left(\sum_{i=1}^k X_i \geq kc \right) = \Pr \left(\sum_{i=1}^k X_i \geq k2m \right)$$

$$\leq \Pr \left(\sum_{i=1}^k X_i \geq k(m + \gamma) \right)$$

$$\leq e^{-\beta mk} = e^{-\alpha ck},$$

where $\alpha = \beta/2$. For our choice of β , we have $\alpha = \lambda/4$, and $c_0 = 4\lambda^{-1} \log 2$. □

ACKNOWLEDGMENT

We would like to thank G. Singh and G. Sideman for helpful comments.

REFERENCES

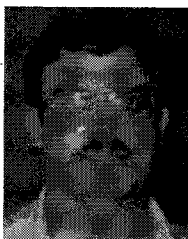
- [1] B. Awerbuch, "Complexity of network synchronization," *J. ACM*, vol. 32, no. 4, pp. 804-823, Oct. 1985.
- [2] ———, "Reducing complexities of distributed max-flow and breadth-first-search algorithms by means of network synchronization," *Networks*, vol. 15, pp. 425-437, 1985.
- [3] B. Awerbuch and D. Peleg, "Network synchronization with polylogarithmic overhead," *Proc. 31st IEEE FOCS*, 1990.
- [4] O. J. Boxma, "Sojourn times in cyclic queues: The influence of the slowest server," in G. Iazeolla, P. J. Courtois, and O. J. Boxma, Eds., *Computer Performance and Reliability*. New York: Elsevier (North-Holland), 1988.
- [5] F. Baccelli, "Ergodic theory of stochastic Petri networks," Rapport INRIA 1037, 1989.
- [6] V. C. Barbosa and E. Gafni, "Concurrency in heavily loaded neighborhood-constrained systems," *ACM Trans. Programming Languages Syst.*, vol. 11, pp. 562-584, Oct. 1989.
- [7] F. Baccelli, P. Konstantopoulos, "Estimates of cycle times in stochastic Petri nets," in I. Karatzas, Ed., *Proc. Rutgers Conf. Stochastic Analysis*. New York: Springer-Verlag, 1991.
- [8] F. Baccelli and Z. Liu, "Comparison properties of stochastic decision-free Petri nets," INRIA Res. Rapport 1433, May 1991.
- [9] F. Baccelli and A. M. Makowski, "Queueing models for systems with synchronization constraints," *Proc. IEEE*, vol. 77, pp. 138-161, Jan. 1989.
- [10] F. Baccelli, A. M. Makowski, and A. Schwartz, "The fork-join queue and related systems with synchronization constraints: Stochastic ordering and computable bounds," *Advances in Applied Probability*, vol. 21, pp. 629-660, Sept. 1989.

- [11] F. Baccelli, W. A. Massey, and A. Towsley, "Acyclic fork-join queueing networks," *J. ACM*, vol. 36, pp. 615-642, July 1989.
- [12] P. Berman and J. Simon, "Investigations of fault-tolerant networks of computers," *Proc. 20th ACM STOC*, 1988.
- [13] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation*. Englewood Cliffs, NJ: Prentice-Hall, 1989.
- [14] K. L. Chung, *Markov Chains with Stationary Transition Probabilities*, 2nd ed. New York: Springer-Verlag, 1967.
- [15] F. Commoner, A. W. Holt, S. Even, and A. Pnueli, "Marked directed graphs," *J. Comput. Syst. Sci.*, vol. 5, no. 5, Oct. 1971.
- [16] K. M. Chandy and L. Lamport, "Distributed snapshots: Determining global states of distributed systems," *ACM Trans. Comput. Syst.*, vol. 3, no. 1, Feb. 1985.
- [17] H. A. David, *Order Statistics*. New York: Wiley, 1970.
- [18] S. Even and S. Rajsbaum, "Lack of global clock does not slow down the computation in distributed networks," Tech. Rep. 522, Dept. of Comput. Sci., Technion, Haifa, Israel, Oct. 1988, to appear in *Mathemat. Syst. Theory* (first part of this paper appears under the title, "Unison in distributed networks" in R. M. Capocelli, Ed., *Sequences: Combinatorica, Compression, Security, and Transmission*. New York: Springer-Verlag, 1990.)
- [19] S. Even and S. Rajsbaum, "The use of a synchronizer yields maximum rate in distributed networks," *Proc. 22nd ACM STOC*, 1990.
- [20] J. Garofalakis, S. Rajsbaum, P. Spirakis, and B. Tampakas, "Tentative and definite distributed computations: An optimistic approach to network synchronization," in *Proc. 6th Int. Workshop on Distributed Algorithms*, Haifa, Israel, Nov. 1992, *Lecture Notes in Computer Science 647*. New York: Springer-Verlag, pp. 110-119 (to appear in *Theoretical Comput. Sci.*).
- [21] J. F. C. Kingman, "Subadditive ergodic theory," *Ann. Prob.*, vol. 1, no. 6, pp. 883-909, 1973.
- [22] M. A. Marsan, "Stochastic Petri nets: An elementary introduction," in *Advances in Petri Nets 1989, Lecture Notes in Computer Science 424*. New York: Springer-Verlag, 1989, pp. 1-29.
- [23] J. Malka, S. Moran, and S. Zaks, "Analysis of a distributed scheduler for communication networks," Tech. Rep. 495, Dept. of Comput. Sci., Technion, Haifa, Israel, Feb. 1988 (also in *Lecture Notes on Computer Science 319*. New York: Springer Verlag, 1988, pp. 351-360).
- [24] M. K. Molloy, "Performance analysis using stochastic Petri nets," *IEEE Trans. Comput.*, vol. 31, no. 9, pp. 913-917, Sept. 1982.
- [25] M. K. Molloy, "Fast bounds for stochastic Petri nets," *International Workshop on Timed Petri Nets*, Torino, Italy, July 1985, pp. 244-249.
- [26] Y. Ofek and I. Gopal, "Generating a global clock in a distributed system," IBM Res. Rep., 1987.
- [27] D. Peleg and A. A. Schaffer, "Graph spanners," *J. Graph Theory*, vol. 13, pp. 99-116, 1989.
- [28] D. Peleg and J. D. Ullman, "An optimal synchronizer for the hypercube," *SIAM J. Computing*, vol. 18, pp. 740-747, Aug. 1989.
- [29] N. Pekergin and J.-M. Vincent, "Stochastic bounds on execution times of task graphs," Rep. EHEI, 1989.
- [30] A. L. Rosenberg, "Shuffle-oriented interconnection networks," COINS Tech. Rep. 88-84, Univ. of Massachusetts, 1988.
- [31] S. Rajsbaum, "Upper and lower bounds for stochastic marked graphs," *Inform. Processing Lett.*, vol. 49, pp. 291-295, 1994.
- [32] Y. Malka and S. Rajsbaum, "Analysis of distributed algorithms based on recurrence relations," in *Proc. 5th Int. Workshop on Distributed Algorithms, Lecture Notes in Computer Science 579*. New York: Springer-Verlag, 1992, pp. 242-253.
- [33] S. M. Ross, *Stochastic Processes*. New York: Wiley, 1983.
- [34] D. Stoyan, *Comparison Methods for Queues and Other Stochastic Models*, (English translation). New York: Wiley, 1984.
- [35] J. G. Shanthikumar and D. D. Yao, "The effect of increasing service rates in a closed queueing network," *J. Appl. Prob.*, vol. 23, pp. 474-483, 1986.



S. Rajsbaum received the degree in computer engineering from the National Autonomous University of Mexico (UNAM) in 1986, and the Ph.D. degree in computer science from the Technion—Israel Institute of Technology, Haifa, Israel, in 1991.

Since 1991, he has been an Associate Professor at the Institute of Mathematics, UNAM. He is currently a Visiting Scientist at the Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA. His general research interests are distributed computation, theoretical computer science, and graph theory. In particular, he has worked on real-time computing, algorithm design and analysis, and on the study of graph tournaments.



M. Sidi (S'77-M'82-SM'87) received the B.Sc., M.Sc., and D.Sc. degrees in electrical engineering from the Technion—Israel Institute of Technology, Haifa, Israel, in 1975, 1979, and 1982, respectively.

In 1982, he joined the faculty of the Department of Electrical Engineering at the Technion. During the academic year 1983-1984, he was a postdoctoral associate at the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA. During 1986-1987, he was a Visiting Scientist at the IBM T.J. Watson Research Center, Yorktown Heights, NY, USA. His current research interests are in queueing systems and computer communication networks.

Dr. Sidi received the New England Academic Award in 1989. He coauthored the book, *Multiple Access Protocols: Performance and Analysis* (Springer-Verlag 1990). Currently, he serves as the Area Editor for communication networks in *IEEE TRANSACTIONS ON COMMUNICATIONS*, as an Editor of *IEEE/ACM TRANSACTIONS ON NETWORKING*, and as the Associate Editor for communication networks and computer networks of *IEEE TRANSACTIONS ON INFORMATION THEORY*.