

Therefore, we have shown the following.

**Theorem 9.7:** If agreement is reached, then all processes eventually terminate with the agreed-upon value. No process may terminate unless all terminate with the same correct value.

*Proof:* By Theorem 9.6 and Theorem 9.5.  $\square$

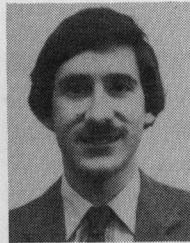
#### ACKNOWLEDGMENT

The author would like to express his appreciation to S. Toueg for the insights and encouragement provided during the course of this work.

#### REFERENCES

- [1] D. Dolev, C. Dwork, and L. Stockmeyer, "On the minimal synchronism needed for distributed consensus," in *Proc. 24th Symp. Foundations of Comput. Sci.*, Tucson, AZ, Nov. 1983, pp. 393-402.
- [2] D. Dolev and H. R. Strong, "Polynomial algorithms for multiple processor agreement," in *Proc. 14th ACM Symp. Theory of Comput.*, 1982, pp. 401-407.
- [3] M. Fisher, N. Lynch, and M. Patterson, "Impossibility of distributed consensus with one faulty process," in *Proc. 2nd ACM Symp. Principles of Database Syst.*, 1983.
- [4] L. Lamport, R. Shostak, and M. Pease, "The Byzantine generals problem," *ACM Trans. Program. Lang. Syst.*, vol. 4, pp. 382-401, July 1982.
- [5] N. Lynch, M. Fischer, and R. Fowler, "A simple and efficient Byzantine generals algorithm," in *Proc. 2nd IEEE Symp. Reliability in Distributed Software and Database Syst.*, Pittsburgh, PA, 1982, pp. 46-52.
- [6] M. Pease, R. Shostak, and L. Lamport, "Reaching agreement in

- the presence of faults," *J. ACM*, vol. 27, no. 2, pp. 228-234, 1980.
- [7] K. Perry, "Early stopping protocols for fault-tolerant distributed agreement," Ph.D. dissertation, Dep. Comput. Sci., Cornell Univ., Ithaca, NY, Jan. 1985.
- [8] M. Rabin, "Randomized Byzantine generals," in *Proc. 24th Symp. Foundations of Comput. Sci.*, Tucson, AZ, Nov. 1983, pp. 403-409.
- [9] R. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Commun. ACM*, vol. 21, pp. 120-126, Feb. 1978.
- [10] A. Shamir, "How to share a secret," *Commun. ACM*, vol. 22, pp. 612-613, Nov. 1979.
- [11] S. Toueg, K. Perry, and T. K. Srikanth, "Simple and efficient Byzantine generals algorithms with early stopping," Dep. Comput. Sci., Cornell Univ., Ithaca, NY, Tech. Rep. TR 84-621, July 1984.
- [12] R. Turpin, and B. Coan, "Extending binary Byzantine agreement to multivalued Byzantine agreement," *Inform. Processing Lett.*, vol. 18, pp. 73-76, Feb. 1984.



Kenneth J. Perry (S'77-M'85) received the B.S.E. degree in electrical engineering and computer science from Princeton University, Princeton, NJ, in 1979, and the M.S. and Ph.D. degrees in computer science from Cornell University, Ithaca, NY, in 1983 and 1985, respectively.

He recently completed a dissertation in the area of fault-tolerant distributed computing. His current research interests include distributed systems, methodology, and artificial intelligence.

Dr. Perry is a member of the Association for Computing Machinery.

# Streets of Byzantium: Network Architectures for Fast Reliable Broadcasts

ÖZALP BABAOĞLU, MEMBER, IEEE, AND ROGÉRIO DRUMMOND

**Abstract**—A site broadcasting its local value to all other sites in a fault-prone environment is a fundamental paradigm in constructing reliable distributed systems. Time complexity lower bounds and network connectivity requirements for reliable broadcast protocols in point-to-point communication networks are well known. In this paper, we consider the reliable broadcast problem in distributed systems with broadcast networks (for example, Ethernets) as the basic communication architecture. We show how properties of such network architectures can be used to effectively restrict the externally visible behavior of faulty

processors. We use these techniques to derive simple protocols that implement reliable broadcast in only two rounds, independent of the failure upper bounds.

**Index Terms**—Byzantine agreement, distributed computing, Ethernet, fault-tolerance, network partitions.

## I. INTRODUCTION

A distributed computing system consists of a set of autonomous processors (sites) that can communicate through a network. To cooperate, these processors must periodically coordinate their activities. One way to do this is by the use of broadcasts—an operation that allows a processor to communicate a local value to all other processors in the system. Since fault-tolerance is one of the primary motivations for con-

Manuscript received May 15, 1984; revised February 7, 1985. This work was supported in part by the National Science Foundation under Grant MCS 82-10356 and in part by the CAPES Agency of the Government of Brazil through a fellowship to R. Drummond.

Ö. Babaoğlu is with the Department of Computer Science, Cornell University, Ithaca, NY 14853.

R. Drummond is with the Department of Computer Science, Cornell University, Ithaca, NY, on leave from the Department of Computer Science, State University at São Paulo, Campinas, Brazil.

structuring a distributed system, the semantics of a broadcast operation have to accommodate failures. In the systems we consider, both processors and communication components can fail and are said to be *faulty* after their first failure. We define the *reliable broadcast* (RB) problem in an unreliable environment as follows.

*Reliable Broadcast*: a distinguished processor, called the *transmitter*, broadcasts a local value (drawn from a domain of at least two elements) such that the following two requirements are satisfied.

1) *Unanimity*: All nonfaulty processors decide on the same value.

2) *Nontriviality*: If the transmitter is nonfaulty, then all nonfaulty processors decide on the value that was broadcast by the transmitter.

The abstract formulation of this problem was first studied by Pease *et al.* [18] under the name *interactive consistency* (later called *Byzantine agreement* [14]), and subsequently by numerous other researchers (see [6] for a brief survey).

The reliable broadcast has been identified as the appropriate primitive in solutions to a large number of problems [10]. These include distributed consensus [6], distributed synchronization [20], replicated data update (e.g., distributing updates to routing tables in a packet switching network) [9], and transaction management in database systems [17].

The properties of the environment in which RB is studied have profound implications on the nature of the solution. In fact, Fischer *et al.* [8] have shown that the problem has no deterministic solution in systems where failures cannot be detected and message delivery times cannot be bounded. In this paper, we will be concerned only with deterministic protocols in synchronous systems. We define such systems in the next section.

Most of the previous formulations of the RB problem assume a reliable point-to-point network as the communication architecture. In such a system, let  $t$  denote the maximum number of processors that can be faulty during the execution of a protocol solving RB. Lower bounds for the total number of processors participating in the protocol, and time and message complexity of the protocol have been obtained as functions of  $t$ . It is customary to measure time complexity in units called *rounds* consisting of message exchange and computation steps, such that during a round every processor can communicate with all others in the system. The messages sent by a processor in a given round cannot depend on the messages it receives during the same round. Of particular importance is the result that at least  $t + 1$  rounds are required for any solution to the RB, even under the most restrictive failure modes and unconstrained amounts of information exchange [7], [4], [11].

In this paper, we dispel the belief that RB is an inherently expensive primitive by breaking the  $t + 1$  round time barrier. We show that by selecting communication architectures that can restrict the external behavior of faulty processors, RB can be implemented in two rounds. We argue that our assumptions about the behavior of the communication network are reasonable in the environments where RB is applicable as a primitive in the first place—highly integrated distributed systems built around local area networks.

## II. MOTIVATION

Arbitrary processor failures in a distributed system using a point-to-point communications network are troublesome for several reasons. First, a faulty processor can send conflicting information along disjoint routes to other processors. Second, unless the communication graph is fully connected, faulty processors on the route of a message can tamper with it. Finally, faulty processors can generate spurious messages and pretend to be forwarding them on behalf of other processors.

The visible external behavior that faulty processors may exhibit can be effectively restricted through the use of certain mechanisms. For example, if the system supports an authentication mechanism such as digital signatures [3], faulty processors are prevented from altering messages and from spontaneously generating spurious messages. This restriction on the behavior of a faulty processor results in algorithms that are simpler and that can tolerate more faults than their counterparts without authentication [4]. In what follows, we propose communication architectures as mechanisms to further restrict the behavior of faulty processors.

Our motivation is based on the following observation. If we can design a communication network architecture such that 1) there exists a single route between any given processor and all other processors, and 2) messages never need to be forwarded (all processors are effectively immediate neighbors of each other), then the behavior of a faulty processor can be limited. In particular, designs based on such architectures need not worry about message tampering and confounding through conflicting messages. *Broadcast networks*, such as ALOHA [1], DCS [5], and Ethernet [16], come very close to realizing these properties.

We pursue this approach through the following organization. The next section is a description of the global system properties and classification of the processor failures that can occur. In Section IV, the major communication network components are identified and their failure behaviors discussed. In Section V, we introduce the class of *R-redundant* broadcast networks. Sections VI and VII contain our main results—protocols to solve RB along with proofs of correctness in *R-redundant* broadcast networks. Network partition failures are examined in Section VIII. Finally, we extend our results to accommodate arbitrary failure behavior of the communication components. A comparison of our results to those of point-to-point communication networks concludes the paper. We note that other researchers have also proposed using broadcast networks to implement reliable broadcasts, however, under much stricter failure assumptions [2], [15].

## III. FAILURE MODES AND SYSTEM PROPERTIES

In a distributed system, the only visible behavior of a processor is the sequence of messages it sends. Consequently, the various failure modes of processors (with respect to a given protocol) can be classified as follows.

*Crash Fault*: A faulty processor stops sending messages after some arbitrary time, never to resume again.

*Omission Fault*: A faulty processor omits sending some of the messages prescribed by its protocol. The messages it does send are always correct.

*Malicious Fault:* A faulty processor exhibits arbitrary behavior. Such a processor can withhold messages that it is required to send, send messages not prescribed by its protocol, and even collude with other (faulty) processors to disrupt the system.

A *nonfaulty* processor sends exactly those messages that are prescribed by its protocol—no more and no less. Clearly, a protocol that can tolerate malicious faults can also tolerate either omission or crash faults. Similarly, a protocol that can tolerate omission faults will tolerate crash faults.

We assume that the distributed system under consideration possesses the following properties.

*Authentication Assumption:* Messages can be signed using digital signatures [3].

*Network Assumption:* In the absence of failures, the communication system guarantees bounded delivery times for messages to all processors.

*Clock Assumption:* Nonfaulty processors maintain local times that do not differ by more than a known constant, and are within a known constant around (unobservable) absolute time.

In systems composed of homogeneous processors that are in close physical proximity (as they are in local area networks), these assumptions seem quite reasonable. Halpern *et al.* [12] have developed an efficient algorithm to achieve the clock assumption in the presence of faulty processors as long as the first two assumptions hold. Recently, Srikanth and Toueg [21] have shown that clock synchronization can be achieved without authentication as long as only one third of the total number of processors can be faulty. In a future paper, we will show how clock synchronization can be realized efficiently and without authentication by exploiting the properties of the networks we are about to propose. Without loss of generality, we will assume that clocks are perfectly synchronized so that we can think of the distributed computation proceeding in lock-step rounds.

Algorithms for approximating the authentication assumption are well known. Depending on the degree of maliciousness that processors may exhibit, messages can be augmented with checksums, cyclic redundancy checks, or digital signatures based on public-key cryptography techniques [19].

#### IV. BROADCAST NETWORKS

Broadcast networks have become a popular topology for local area networks. Some of the reasons for this popularity are simplicity, high reliability, high performance, and very low cost. Perhaps the best-known example of this class of networks is the Ethernet. Fig. 1 shows a typical distributed system configuration based on a broadcast network.

The central component of a broadcast network is the *channel*. Typically, this is a passive transmission medium such as a coaxial cable, fiber optic cable, or simply space (in the case of a radio broadcast network). Processors are connected to the channel through *links*. Links contain active components such as signal transceivers, network interface units, and buffer memory. Access to the channel by the links is controlled through a multiple-access transmission protocol [22]. We note that collision-free transmission protocols exist for multiple-access channels that can guarantee the network assump-

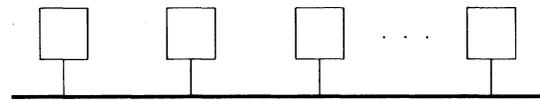


Fig. 1. Broadcast network.

tion [22]. Even with transmission protocols where collisions may occur (such as in Ethernet), the message delivery time can be bounded with high probability by keeping the channel lightly loaded.

The fundamental communication operation implemented by a broadcast network is **broadcast**( $v$ ), where  $v$  is the message being communicated. Through the authentication assumption, a processor receiving a message is able to reliably identify the broadcaster under any possible processor failure mode. We abstract the operation of the channel and the links as follows.

*Nonfaulty Channel Semantics:* In the absence of failures, the channel arbitrarily selects one of the links to listen to and simply echos the incoming message (if any) to all of the other links.

*Faulty Channel Semantics:* A fault in the channel causes some of the incoming messages to be lost (i.e., echoed to none of the other links).

*Nonfaulty Link Semantics:* In the absence of failures, the link simply transfers messages between the channel and its associated processor.

*Faulty Link Semantics:* A fault in the link causes some arbitrary incoming messages, either from the processor or the channel, to be lost.

Broadcast networks that exhibit these behaviors can be readily implemented. We further require that neither a faulty channel nor a faulty link can spontaneously generate (valid) messages. This requirement is usually achieved by standard transmission protocols.

As defined, our channel semantics ensure that either all (in case the channel is nonfaulty) or none (in case the channel is faulty) of the processors connected to a given channel through nonfaulty links receive the broadcast message. There are rare channel failures that result in the partitioning of the links (and thus processors) such that some of them receive a broadcast message and others do not. We consider these types of failures in Section VIII.

To simplify our presentation, we will assume that message authentication and validation (using the appropriate mechanisms) are performed at a level below the RB protocol. Messages that are invalid are simply discarded and are never presented to the RB protocol level. If no message is delivered to the RB protocol level after a time interval equal to the network delivery time bound, we say that a null message (denoted as  $\phi$ ) is received.

Given that processors communicate through a network with the above properties, we note that there is only one externally visible system behavior for any possible failure scenario.

*Broadcast Network Failure Property:* In response to a broadcast, those processors that receive a (nonnull) message receive the same message.

This behavior is ensured even if processor failures are malicious, since it is physically impossible for a processor to send conflicting messages to different processors in the same broad-

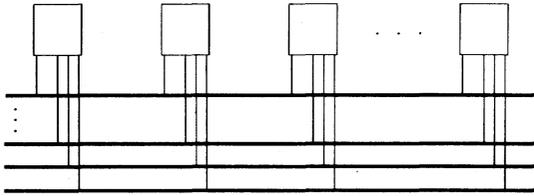


Fig. 2.  $R$ -redundant broadcast network.

cast. In other words, we have succeeded in restricting the visible affects of processor failures.

We now discuss the power of this property when applied to solve the RB problem in the presence of processor, link, and channel failures.

### V. REDUNDANT BROADCAST NETWORKS

Even with the broadcast network failure property, the  $\text{broadcast}(v)$  operation alone does not achieve the unanimity requirement for RB since a single broadcast network cannot guarantee 1-connectivity<sup>1</sup> among processors in the presence of communication faults. To increase connectivity, we propose an  $R$ -redundant broadcast network generated by replicating each communication component  $R$  fold. As shown in Fig. 2, the construction requires  $R$  independent broadcast networks, where each of the processors is connected to each of the  $R$  channels through an independent link. In such a network, we account separately the various components that exhibit faults. During some execution of a RB protocol, let

- $\pi$  be the actual number of faulty processors,
- $\lambda$  be the actual number of faulty links, and
- $\gamma$  be the actual number of faulty channels.

We say that a processor is *properly connected* to a channel if its link to that channel is nonfaulty. By construction, it is trivial to show the following property of  $R$ -redundant broadcast networks.

*Property 1:* In an  $R$ -redundant broadcast network, the condition  $R > \lambda + \gamma$  is sufficient to guarantee that any pair of processors will remain at least 1-connected for the duration of the protocol.

In the subsequent discussions, we will assume that this property holds without explicitly repeating it.

Since every processor in an  $R$ -redundant broadcast network is connected to  $R$  channels, we augment the broadcast primitive syntax as follows:  $\text{broadcast}(v, i)$ ; where  $v$  is the message to be communicated, and  $i \in \{1, 2, \dots, R\}$  selects a channel. We will use the following shorthand for notational convenience:

```

broadcast( $v, *$ )  $\equiv$  for each  $i$  in  $\{1, 2, \dots, R\}$ 
do
    broadcast( $v, i$ )
od;
    
```

Note that we do not require this operation to be atomic nor do we make any assumptions about the relative order or tim-

<sup>1</sup>Since no message forwarding ever takes place in the networks we consider, we are only concerned with *edge connectivity*. We say that a network has  $k$ -connectivity if there exist at least  $k$  disjoint paths consisting of links and channels between every pair of processors.

Round 1: For  $p_1$  (the transmitter with local value  $v \in \mathcal{V}$ )

```

broadcast( $v, *$ );
decide on  $v$ ;
    
```

Round 2: For all  $p_i, i \neq 1$

```

if  $|C_i^1| \neq 0$  then
     $m_i$  is the (unique) message received in round 1
    foreach  $k$  in  $\{1, 2, \dots, R\} - C_i^1$  do
        broadcast( $m_i, k$ );
    od;
    decide on  $m_i$ ;
fi;
    
```

At the conclusion of round 2: For all  $p_i$

```

if not yet decided then
    if  $|C_i^2| \neq 0$  then
         $m_i$  is the (unique) message received in round 2
        decide on  $m_i$ ;
    else
        decide on  $\delta$ ;
    fi;
fi;
    
```

Fig. 3. Protocol  $P1$ —reliable broadcast for omission fault processors.

ing of the broadcasts over the various channels. Furthermore, no assumptions are made about the behavior of faulty processors with respect to messages broadcast over different channels. In particular, a faulty processor can broadcast conflicting messages over different channels. All that we rely on is the broadcast network failure property to hold for each of the broadcast networks.

### VI. RELIABLE BROADCAST PROTOCOL FOR OMISSION FAULT PROCESSORS

While crash faults are more restrictive failures than omission faults, it is interesting to note that restricting processors and network components to exhibit crash faults does not render the RB problem any easier. In fact, Hadzilacos [11] has shown that there exists a transformation procedure for converting any solution to the RB for crash fault components into an equivalent solution (in the sense of having the same time and information complexities) for omission fault components. Therefore, we will initially consider processors that fail through omission faults.

Let  $p_1, p_2, \dots, p_N$  denote processors and assume  $p_1$  is the transmitter. Let  $\delta$  denote a default value known *a priori* to all of the processors such that all of the nonfaulty ones will decide on  $\delta$  in case the transmitter is deemed faulty. Let  $\mathcal{V}$  denote the domain of all possible valid message values not including  $\phi$ . We say that an RB protocol is *symmetric* if the computation carried out by each processor is independent of its name, with the exception of the transmitter which is distinguished from the other processors. We do not consider the message authentication steps, which can be processor name dependent, to be part of the RB protocol. Each processor executes the protocol of Fig. 3. Note that all variables are subscripted with the processor name only to denote that they are local variables. Consequently, the protocol is symmetric according to our definition. Recall that the values broadcast in a given round cannot depend on messages received in that round. In the protocol, the notation  $C_i^r$  denotes the set of channels from which processor  $p_i$  received the (nonnull) message  $m_i \in \mathcal{V}$  at round  $r$ . Note that message  $m_i$  is unique since

processor failures are restricted to omission faults. The following theorem establishes the correctness of the protocol.

*Theorem 1:* If  $N \geq \lambda + \pi$ , then protocol *P1* implements RB for processors that may fail according to omission faults.

*Proof:* We proceed by case analysis.

If the transmitter is nonfaulty, all processors receive the broadcast in round 1 (and know that it is from the transmitter) and decide on  $v$  at round 2. Thus, nontriviality and unanimity are assured. The following assure unanimity in cases of transmitter failure.

If the transmitter is faulty and broadcasts no messages over nonfaulty channels to which it is properly connected in the first round, by the end of the second round only  $\phi$  messages have been received and all nonfaulty processors decide on  $\delta$ .

Next, consider the case where the transmitter is faulty and some nonfaulty processor  $p_i$  receives the initial broadcast. Since processor failures are restricted to omission faults, the nonnull messages  $m_i$  received by  $p_i$  must all be equal to  $v$ . The protocol requires  $p_i$  to echo  $m_i$  through all channels not in  $C_i^1$ . Therefore, by the end of the second round,  $v$  messages will have reached all nonfaulty channels to which  $p_i$  is properly connected. By the 1-connectivity property, it is guaranteed that all processors receive  $v$  in one of the two rounds. As receiving a nonnull message in either of the two rounds implies a decision, all processors decide on  $v$ .

Finally, we consider what happens if the transmitter is faulty and only other faulty processors receive the initial broadcast. This can only occur when all nonfaulty processors have faulty links to the nonfaulty channels  $p_1$  broadcast over in the first round. The minimum number of link faults necessary by the end of round 1 for this scenario is  $N - \pi$ , and comes about when  $p_1$  broadcasts over a single nonfaulty channel to which it is properly connected and all of the nonfaulty processors have faulty links to this channel. Since  $N \geq \lambda + \pi$  and  $\lambda = N - \pi$ , there can be no further links. If one of the faulty processors echos  $v$  during round 2 over a channel other than the one  $p_1$  broadcast over, then all nonfaulty processors are guaranteed to receive it and decide on it. Otherwise, all nonfaulty processors receive only  $\phi$  messages by the end of round 2 and decide on  $\delta$ .  $\square$

From the protocol description, it is obvious that *P1* results in  $N \cdot R$  messages to be exchanged, and requires two rounds to implement RB. The next result shows the optimality of *P1* with respect to time complexity within the class of redundant broadcast networks.

*Theorem 2:* For the class of  $R$ -redundant broadcast networks, there exists no symmetric one-round protocol to implement RB in the presence of even nonfaulty channels, a single faulty processor, and a single faulty link.

*Proof:* Without loss of generality, assume that the transmitter is the faulty processor and fails according to crash faults. Under any (one-round) protocol, there must be a step where  $p_1$  broadcasts its local value. Assume that  $p_1$  fails after having broadcast  $v$  over  $k$  channels ( $0 < k \leq R$ ). As before,  $C_i^1$  denotes the set of channels over which processor  $p_i$  received nonnull messages from the transmitter by the end of round 1. Note that during the single round, any two processors can also exchange information through broadcasts. However, this infor-

mation cannot depend on what each one received from the transmitter since this occurs in the same (single) round. Furthermore, since we require that the protocol be deterministic and symmetric, the information content of this communication between the processors is zero. Note that each channel could have a different name with respect to different processors. This fact, along with the requirement that the protocol be symmetric, prevents the decision to be based on the actual elements of the set  $C_i^1$ . This leaves the number of nonnull messages received by each processor from the transmitter by the end of the round as the only relevant information for deciding. Therefore, for any symmetric one-round protocol, there must exist a common function  $\psi$  such that each processor  $p_i$  decides on  $\psi(|C_i^1|)$  at the end of the round. From the nontriviality requirement of RB, we must have  $\psi(R) = v$ . Furthermore,  $\psi(0) = \delta$  is simply due to a complete lack of global information. Consequently, there must exist some threshold value  $z > 0$  such that  $\psi(z - 1) = \delta$  and  $\psi(z) = v$ . Now, in the above failure scenario, let  $k = z$ , whereby all proper processors decide on  $v$ , except for the processor with the single faulty link to one of the  $k$  channels which has to decide on  $\delta$ . We conclude that the unanimity requirement for RB cannot be assured under any symmetric one-round protocol.  $\square$

We next consider the RB problem when processor faults are unconstrained.

## VII. RELIABLE BROADCAST PROTOCOL FOR MALICIOUS FAULT PROCESSORS

By far the most important property of a broadcast network is that it ensures unanimity among the receivers of a broadcast in the absence of network failures. It is impossible for a (maliciously) faulty processor to send conflicting values to different processors in a single broadcast over the same channel. Our motivation for the construction of  $R$ -redundant broadcast networks was to increase connectivity so that network failures could be tolerated in an effective manner. However, in so doing, we have left the system vulnerable to malicious failures, whereby conflicting messages can be sent by a faulty processor over each of the  $R$  components of the redundant network during a single **broadcast**( $v, *$ ) operation.

In this section, we present an efficient protocol that implements RB in  $R$ -redundant broadcast networks in the presence of malicious fault processors. Our assumptions about the behavior of the communication system continue to apply. In terms of fault-tolerant system design, malicious fault processors represent a "worst case" design strategy. Even in this case, our protocol requires only two rounds to implement RB. The only additional cost with respect to protocol *P1* is an increase in the total number of processors that must participate in the protocol for a given level of fault tolerance.

Before discussing the protocol, we introduce some notation and definitions. Let  $W = \{\phi, \perp\} \cup V$  denote the domain of message values including the null message and the special error symbol  $\perp$ . The set of all possible bags (multisets) of size  $k$ , constructed from the elements of a set  $S$ , is denoted as  $S^k$ . In the protocol, every nonfaulty processor broadcasts at most one message over a channel in a single round. Let  $B_i^r(j) \in W^R$  be the *in-bag* of  $R$  messages (one from each channel) processor

Round 1: For  $p_1$  (the transmitter with local value  $v \in \mathbf{V}$ )  
**broadcast**( $v, *$ );

Round 2: For all  $p_i$   
 $m_i \leftarrow \sigma(B_i^1(1));$   
**broadcast**( $m_i, *$ );

At the conclusion of round 2: For all  $p_i$   
 $M^i \leftarrow \{j=1,2,\dots,N: \sigma(B_i^2(j))\};$   
**decide on**  $\psi^z(M^i)$ ;

Fig. 4. Protocol  $P2^z$  with threshold  $z$ -reliable broadcast for malicious fault processors.

$p_i$  receives from processor  $p_j$  during round  $r$ . If a processor  $p_i$  receives more than one message over the same channel from some (faulty) processor  $p_j$  during some round  $r$ , it sets the in-bag  $B_i^r(j) = \{\perp\}^R$  immediately, regardless of what messages it received (or will receive) from  $p_j$  over other channels in the same round.

We say that a bag  $B$  is *consistent in*  $v$  if and only if  $B = \{\phi, v\}^R$  for some unique  $v \in W$  and  $v \neq \perp$ . In other words, the bag contains copies of either a single nonerror message or two messages, one of which has to be the null message. The *filter function*

$$\sigma: W^R \rightarrow W$$

maps bags of size  $R$  to the valid message domain is defined as follows:

$$\sigma(B) = \begin{cases} v & \text{if } B \text{ is consistent in } v \\ \delta & \text{otherwise.} \end{cases}$$

Finally, we define the *decision function* with threshold  $z$

$$\psi^z: W^N \rightarrow V$$

that maps bag  $M \in W^N$  of  $N$  messages such that

$$\psi^z(M) = \begin{cases} v & \text{if } \rho(M) = v \text{ and } \theta(M, v) \geq z \\ \delta & \text{otherwise} \end{cases}$$

where  $\rho(M)$  selects the most frequent non- $\phi$  element in the bag  $M$ , and  $\theta(M, x)$  counts the number of times element  $x$  appears in  $M$ .

The protocol for malicious fault processors is shown in Fig. 4. Let us consider the result of a processor failure during the **broadcast**( $v, *$ ) operation in this system. Note that if the failure is such that a processor broadcasts conflicting messages over faulty channels, or through faulty links to nonfaulty channels, the failure is hidden from all other processors and the transmitter appears to be nonfaulty. It is only when the conflicting messages are visible to other processors (because they reach nonfaulty channels) that unanimity can be compromised. With this in mind, we say that a sending processor is *overtly malicious* if the bag of message values, as observed on the channels due to a broadcast, is not consistent. A receiving processor  $p_i$  is called a *witness* to this type of failure of processor  $p_j$  in round  $r$  if and only if  $B_i^r(j)$  is not consistent.

Similarly, a processor exhibits *overt omission* failure if it is not overtly malicious, but broadcasts a null message over at least one nonfaulty channel to which it is properly connected.

In an  $R$ -redundant broadcast network, processor  $p_i$  *witnessing* an overt omission failure of processor  $p_j$  in round  $r$  requires that  $B_i^r(j) = \{\phi\}^R$ .

As we have defined it, the filter function  $\sigma$  is designed to detect exactly these failures. It results in  $\phi$  for the witnesses of overt omission failures, and in  $\delta$  for the witnesses of overt commission failures. We will next present results that bound the number of processors that can be witnesses to the failure of the transmitter, and we show that under no circumstance can other faulty processors collude to prevent unanimity. This is accomplished by ensuring that a sufficiently large number of nonfaulty processors are biased toward the same decision by the end of the first round.

*Lemma 1:* In a  $R$ -redundant broadcast network, the number of witnesses to an overt omission failure of the transmitter by the end of the first round in protocol  $P2$  is either  $N - 1$  or in the range  $[0, \lambda]$ .

*Proof:* Let  $x$  be the number of nonfaulty channels to which  $p_1$  is properly connected and sends nonnull messages. Clearly,  $0 \leq x < R$  since  $x = R$  would imply that the transmitter is nonfaulty. Since 1-connectivity is ensured, when  $x = 0$  all  $N - 1$  remaining processors witness the failure. If  $x > 0$ , then for a processor to witness the failure it must have faulty links to all  $x$  channels. As there are at most  $\lambda$  actual faulty links, the number of such witness processors is  $\lfloor \lambda/x \rfloor$ . This results in the range  $[0, \lambda]$ .  $\square$

*Lemma 2:* In an  $R$ -redundant broadcast network, the number of witnesses to an overtly malicious transmitter by the end of the first round in protocol  $P2$  is in the range  $[N - \lambda - 1, N - 1]$ .

*Proof:* Let  $p_1$  fail by broadcasting two different (nonnull) messages over two nonfaulty channels to which it is properly connected. For a processor not to witness this failure, it must not be properly connected to one of these channels. Since there are  $\lambda$  possible faulty links, there can be at most  $\lambda$  such processors. Consequently, the number of witnesses (not counting the faulty transmitter) is at least  $N - \lambda - 1$  and at most  $N - 1$ .  $\square$

Recall that  $t$  is the upper bound for the number of actual faulty processors  $\pi$ .

*Theorem 3:* If  $N > t + \pi + 2\lambda$ , then protocol  $P2$  with threshold  $z = t + 1$  implements RB for processors that may fail according to malicious faults.

*Proof:* We proceed by case analysis.

If the transmitter is nonfaulty, all of the in-bags will be consistent in  $v$  by the end of the first round. For all nonfaulty processors, the filter function yields  $v$ , which is echoed in round 2. The  $N - \pi$  such echos will be filtered by the end of the second round and added as  $v$  values to the  $M^i$  bag for each processor. For all nonfaulty  $p_i$ , the bags  $M^i$  will contain at least  $N - \pi$  occurrences of  $v$  and most  $\pi$  occurrences of some other value. Even if all of the  $\pi$  faulty processors collude and try to force a decision on this other value, it will fail since  $\pi \leq t$  and the threshold of the decision function will not be exceeded. The condition  $N > t + \pi + 2\lambda$  can be written as  $N - \pi \geq t + 2\lambda + 1$ , and since  $\lambda \geq 0$ , we have  $N - \pi \geq t + 1$  so that  $v$  will be selected by all of the decision functions of the nonfaulty processors. Thus, we have shown nontriviality for pro-

tolocol  $P2$ . We next argue that it also ensures unanimity if the transmitter is faulty.

Assume that the transmitter fails by overt omission. By Lemma 1, there will be  $k \in \{0, 1, \dots, \lambda, N-1\}$  witnesses. If  $k = N-1$ , then at least  $N - \pi$  processors will broadcast  $\phi$  in the second round. By an argument identical to that of the previous paragraph, all nonfaulty processors decide on  $\delta$ . If  $k \leq \lambda$ , then  $N - \pi - k$  nonfaulty processors do not witness the failure, and for each, the filter function yields the value  $v$  since each has an in-bag consistent in  $v$ . In the second round,  $N - \pi - k$  consistent  $v$  messages are broadcast such that for each nonfaulty processor  $p_i$ , the bag  $M^i$  will contain  $N - \pi - k$  occurrences of  $v$ . Since  $k \leq \lambda$ , the inequality  $N - k - \pi \geq t + 1$  is assured by the constraint  $N > t + \pi + 2\lambda$ , and the decision function yields  $v$  for all nonfaulty processors.

Finally, assume that the transmitter is overtly malicious. By Lemma 2,  $N - \lambda - 1 \leq k \leq N - 1$  processors witness the failure. In the second round, at least  $k - (\pi - 1)$  processors will echo  $\delta$  messages. By the end of the second round,  $M^i$  for each nonfaulty  $p_i$  will contain at least  $k - \pi + 1$  occurrences of  $\delta$ . Since for all possible values of  $k$  we have  $k - \pi + 1 \geq t + 1$ , the decision function yields  $\delta$ .  $\square$

### VIII. COPING WITH CHANNEL PARTITIONS

Consider the failure mode of an individual broadcast network that partitions the set of links into two disjoint groups—those that receive a message and those that do not. For example, in an Ethernet, the coaxial cable physically breaking, the failure of the repeater unit between two channel segments and an interface unit causing collisions shorter in duration than the end-to-end propagation delay of the channel are rare instances of partition failures. We now relax the faulty channel semantics to allow these partition failures.

*Faulty Channel Semantics:* A faulty channel fails to deliver a broadcast message to a subset of the links connected to it. Those links that do receive a message receive the same message.

With respect to a given broadcast, the group of links receiving the message is called the *in-partition*, and the group not receiving it is called the *out-partition*. Note that our definition allows the composition of these partitions to change dynamically from one message to another and from one broadcaster to another. Let us observe the behavior of the communication system from a processor's point of view for a given broadcast. The case where the processor's link is nonfaulty but in an out-partition, and the case where its link has failed and is in an in-partition, are indistinguishable—they both result in the processor receiving a  $\phi$  message. This observation suggests a duality between link and channel failures. We can either assume that channels are failure free and simulate a partition failure by pretending that all of the links in the out-partition have failed regardless of their actual state, or we can assume that links are failure free and simulate their failure by constructing the appropriate out-partitions for channel failures. We will pursue the first alternative.

Now, links can be "faulty" either because they have failed or because they belong to the out-partition of a faulty channel. With this new interpretation, all of our previous results continue to hold for the more general faulty channel seman-

tics. All that is required is that we be able to account properly for the total number of faulty links. Next, we characterize each channel according to the configuration of link failures to that channel. For channel  $i$ , let  $\lambda_i^m$  denote the number of processors that do not receive a given broadcast message  $m$  over that channel. If  $\lambda_i^m = 0$  then channel  $i$  has not failed. If  $\lambda_i^m = N$  then channel  $i$  is said to have *failed totally*. For any other value, channel  $i$  is said to have *failed partially*. Throughout the execution of a protocol, a channel is nonfaulty if it has not failed for any broadcast message. If all of the channel failures are total, then it is said to exhibit a *total fault*.<sup>2</sup> Any other channel behavior is called a *partition fault*. We let  $\gamma'$  count the number of faulty channels (either total or partition).

We will first examine the implications of the new semantics for faulty channels on the replication factor of our  $R$ -redundant broadcast networks. Recall that in such networks,  $R > \lambda + \gamma$  is a sufficient condition for 1-connectivity. Clearly, the network can sustain more than  $R - \gamma$  faulty links and still preserve 1-connectivity if there are multiple faulty links to the same channel. Since we do not allow any message forwarding by processors, 1-connectivity can be assured if there is at least one nonfaulty channel with all nonfaulty links. According to our new notation, 1-connectivity is guaranteed if  $R > \gamma'$ .

For each channel, we are interested in the maximum number of processors that do not receive some message after a broadcast. Let  $\lambda_i$  be the maximum of  $\lambda_i^m$  over all messages  $m$  broadcast over channel  $i$  during the execution of the protocol, such that  $\lambda_i^m \neq N$ . If the channel exhibits a total fault, then we define  $\lambda_i$  to be zero for that channel. Note that this definition results in  $\lambda_i = 0$  for both channels that fail totally and that are nonfaulty. Under both circumstances, all processors receive the same message (either  $\phi$  or what was broadcast) and the  $N$  "faulty" links need not be accounted (except for connectivity reasons). Let  $\lambda'$  be the sum of  $\lambda_i$  over all channels. With this new interpretation, all of our previous results can be extended to channel partition failures simply by replacing all occurrences of  $\lambda$  with  $\lambda'$ . In essence,  $\lambda'$  is an upper bound on the total number of link failures that can occur in each round. Note that the number of links that fail at least once during the execution of the protocol can be much larger than  $\lambda'$ .

For certain network configurations, we can tolerate even more link faults. These configurations are characterized as permitting a bound on the  $\lambda_i$  for each channel. For example, consider a broadcast network consisting of several Ethernet segments connected through repeaters. If repeater failures are the only sources of channel partitions, and the number of processors connected to each segment is known, we can easily express the minimum number of processors that will receive every message broadcast over that channel. Since  $\lambda_i$  denotes the maximum number of processors that do not receive some message over channel  $i$ ,  $N - \lambda_i$  denotes the minimum number of processors that receive every message over the same channel. If  $N - \lambda_i > \pi$  can be ensured for each channel, we know that at least one nonfaulty processor will receive each message broadcast over each channel. For omission fault processors,

<sup>2</sup>Note that this behavior corresponds to the earlier semantics for faulty channels whereby some messages are lost completely.

this condition is sufficient to guarantee that protocol  $P1$  continues to implement RB in two rounds.

In summary, if the faulty links that simulate partitions are well structured (as in the case with multiple Ethernet segments and repeaters), our protocols can sustain many more than what is allowed by the pessimistic accounting through  $\lambda'$ . As an example, Theorem 1 requires that the total number of faulty links not exceed  $N - \pi$  for protocol  $P1$  to work. However, if  $\lambda_i < N - \pi$  for each channel, the system can tolerate almost  $R - 1$  times as many faulty links.

## IX. MALICIOUS COMMUNICATION COMPONENT FAILURES

The failure mode for links and channels we have examined so far can be characterized as omission faults. In this section, we examine the consequences of network component failures that can undetectably alter the contents of messages.

Given that the communication media can alter messages, receiving a message over a single channel does not permit the processor to infer the actual message broadcast over that channel. Let  $G > \gamma'$  be an upper bound on the number of faulty channels. If we ensure that  $R > 2G$  and each processor communicates through **broadcast**( $v$ , \*), we can overcome the behavior of faulty communication components. If a nonfaulty processor broadcasts the message  $v$ , then all processors will receive at least  $G + 1$  copies of  $v$  since, at most,  $G$  copies could have been altered by faulty paths. Since there are more non-faulty paths between processors than faulty ones, receivers can always determine the actual broadcast message by simply selecting the majority element from the bag of received messages. Having guaranteed effective communication through this mechanism, the subsequent discussion will assume channels that are failure free. We can do this without loss of generality since malicious channel failure behavior can be simulated through the failure of links.

We modify slightly our previous notion of consistency for bags. A bag  $B$  is said to be *consistent in  $v$*  if and only if  $v$  is a valid message and it is the majority element. A bag that is not consistent in any value is called *inconsistent*.

Let  $B^*$  be the bag of messages as observed on the channels due to the execution of a broadcast by some processor over the  $R$  channels. We will now prove that protocol  $P2$  achieves RB under malicious processor and link faults.

*Lemma 3:* If  $B^*$  is consistent in  $v$ , then

- 1) at most  $\lambda'$  processors may receive an inconsistent message,
- 2) at most  $\lambda'$  processors may receive a consistent message in some value  $y$  different from  $v$ , and
- 3) at least  $N - \lambda'$  processors receive a message consistent in  $v$ .

*Proof:* For the first two cases, a faulty link to a processor could turn a bag consistent in  $v$ , as observed over the channels, into an inconsistent bag, or a bag consistent in some other value  $y$  as seen by the processor. In the worst case, there can be  $\lambda'$  such processors since there are, at most,  $\lambda'$  link failures in each round. As a consequence, at least  $N - \lambda'$  processors have nonfaulty links and receive the bag of messages consistent in  $v$ .  $\square$

*Lemma 4:* If  $B^*$  is inconsistent then, at most,  $\lambda'$  processors receive a bag of messages consistent in some value.

*Proof:* This trivially follows from the observation that a

change to a single element of an inconsistent bag could make it consistent.  $\square$

Given these two Lemmata, we next show how protocol  $P2$ , with the new definition of consistency, can be adapted to support arbitrary communication failures. Let  $L$  denote the upper bound for  $\lambda'$  for all executions of the protocol.

*Theorem 4:* If  $N > t + L + \pi + \lambda'$ , then protocol  $P2$  with threshold  $z = t + L + 1$  implements RB when processors and communication components may fail according to malicious faults.

*Proof:* If the transmitter is nonfaulty, the bag of messages observed over the channels resulting from its initial broadcast is consistent in  $v$ . By Lemma 3, at least  $N - \lambda'$  processors receive bags consistent in  $v$ . Among them, at least  $N - \lambda' - \pi$  processors are nonfaulty and echo  $v$  in the second round. Since by the statement of the theorem  $N - \lambda' - \pi$  is greater than or equal to  $z$  and  $\lambda' + \pi$  is less than  $N - \lambda' - \pi$ , nontriviality is ensured.

If the transmitter is faulty, the bag of messages observed over the channels resulting from its initial broadcast can be either consistent in  $v$  or inconsistent. If it is consistent, we have a decision on  $v$  as above. Otherwise, by Lemma 4, at least  $N - \lambda' - \pi$  nonfaulty processors receive inconsistent bags. In the second round these processors echo  $\delta$ . By the same argument, as in the previous paragraph, a decision on  $\delta$  is ensured.  $\square$

Note that the modified protocol  $P2$  with threshold  $z = L + 1$  implements RB for omission fault processors and arbitrary communication failures, provided  $N > L + \pi + \lambda'$ .

## X. DISCUSSION AND CONCLUSIONS

In this paper, we have shown that the reliable broadcast operation in distributed systems need not be expensive (in time) if the correct model of communication is realized through the network architecture. Furthermore, additional degrees of fault-tolerance can be "bought" simply by increasing the number of appropriate system components but without slowing the system down.

From a practical point of view, we must be careful in comparing two different protocols based on their time complexity as measured in rounds. The total communication time associated with a single round of protocols  $P1$  and  $P2$  can be structured to require exactly  $N$  single message transmission (end-to-end delay) times. If the system is composed of a very large number of processors, this communication step could account for an inordinate amount of time. However, in such cases, only the minimum number of processors required by the protocol for a given degree of fault-tolerance need actively participate in the protocol execution. For example, in protocol  $P1$ , only  $\lambda + \pi$  processors are required in implement RB for a given value of  $\lambda + \pi$ . If  $N \gg \lambda + \pi$ , then we can simply arrange that only  $\lambda + \pi$  of the total available processors actively participate in the protocol execution (i.e., by broadcasting messages) while the others participate as passive listeners. All of the processors execute the same protocol, but the passive ones simply refrain from performing any of the broadcasts that are prescribed. In this manner, the decision they reach will be consistent with that of the active processors. Obviously, the same observation applies to the other protocols.

The parallelism within the communication step of a round

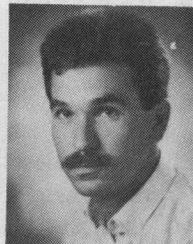
for a broadcast network architecture is limited by the number of physical media that exist within the channel implementation. If there is a single medium that has to be time multiplexed (as it usually is the case), the communication step duration grows proportionally to the number of processors in the system. On the other hand, a point-to-point communication architecture allows for parallelism of message transmissions over all of the available links. In particular, if the communication topology forms a fully connected graph, all processors can communicate with each other in a single end-to-end delay time since all message transmissions can proceed in parallel. However, we conjecture that for point-to-point communication networks with  $N$  processors and only  $O(N)$  links (i.e., less than fully connected), the communication step duration of a round will require  $O(N)$  end-to-end delay times just as in broadcast networks.

#### ACKNOWLEDGMENT

We are grateful to F. Schneider and G. Neiger for discussing this material with us and commenting on earlier versions of this manuscript. We would also like to thank F. Cristian for his helpful comments.

#### REFERENCES

- [1] N. Abramson, "The ALOHA system," in *Computer-Communication Networks*, N. Abramson and F. Kuo, Eds. Englewood Cliffs, NJ: Prentice-Hall, 1973.
- [2] J. M. Chang and N. F. Maxemchuk, "Reliable broadcast protocols," *ACM Trans. Comput. Syst.*, vol. 2, pp. 251-273, Aug. 1984.
- [3] W. Diffie and M. Hellman, "New directions in cryptography," *IEEE Trans. Inform. Theory*, vol. IT-22, pp. 644-654, 1976.
- [4] D. Dolev and H. R. Strong, "Authenticated algorithms for Byzantine agreement," *SIAM J. Comput.*, vol. 12, pp. 656-666, Nov. 1983.
- [5] D. Farber, J. Feldman, F. Heinrich, M. Hopwood, K. Larson, D. Lomms, and L. Rowe, "The distributed computing system," in *Proc. CompCon 73*, IEEE Computer Society, Feb. 1973, pp. 31-34.
- [6] M. J. Fisher, "The consensus problem in unreliable distributed systems (A brief survey)," Dep. Comput. Sci., Yale Univ., New Haven, CT, Tech. Rep. YALEU/DCS/RR-273, June 1983.
- [7] M. J. Fisher and N. A. Lynch, "A lower bound for the time to assure interactive consistency," *Inform. Proc. Lett.*, vol. 14, pp. 183-186, Apr. 1982.
- [8] M. J. Fisher, N. A. Lynch, and M. S. Paterson, "Impossibility of distributed consensus with one faulty process," in *Proc. 2nd ACM Symp. Principles of Database Syst.*, Mar. 1983.
- [9] H. Garcia-Monlina, F. Pitelli, and S. Davidson, "Applications of byzantine agreement in database systems," Dep. Elec. Eng. Comput. Sci., Princeton Univ., Princeton, NJ, Tech. Rep. TR-316, June 1984.
- [10] N. H. Gehani, "Broadcasting satellite processes (BSP)," *IEEE Trans. Software Eng.*, vol. SE-10, pp. 343-351, July 1984.
- [11] V. Hadzilacos, "Issues of fault tolerance in concurrent computations," Ph.D. dissertation, Aiken Computation Lab., Harvard Univ., Cambridge, MA, Tech. Rep. TR-11-84, June 1984.
- [12] J. Halpern, B. Simons, R. Strong, and D. Dolev, "Fault-tolerant clock synchronization," in *Proc. 3rd ACM Symp. Principles Distributed Comput.*, Vancouver, B. C., Canada, Aug. 1984, pp. 89-102.
- [13] L. Lamport, "Using time instead of timeout for fault-tolerant distributed systems," *ACM Trans. Program. Lang. Syst.*, vol. 6, pp. 254-280, Apr. 1984.
- [14] L. Lamport, R. Shostak, and M. Pease, "The Byzantine Generals problem," *ACM Trans. Program. Lang. Syst.*, vol. 4, pp. 382-401, July 1982.
- [15] G. LeLann, "A distributed system for transaction processing," *IEEE Computer*, vol. 14, pp. 43-48, Feb. 1981.
- [16] R. Metcalfe and D. R. Boggs, "Ethernet: Distributed packet switching for local computer networks," *Commun. ACM*, vol. 19, pp. 396-403, July 1976.
- [17] C. Mohan, H. R. Strong, and S. Finkelstein, "Method for distributed transaction commit and recovery using Byzantine agreement within clusters of processors," in *Proc. 2nd ACM Symp. Principles Distributed Comput.*, Montreal, P.Q., Canada, Aug. 1983, pp. 89-103.
- [18] M. Pease, R. Shostak, and L. Lamport, "Reaching agreement in the presence of faults," *J. ACM*, vol. 27, pp. 228-234, Apr. 1980.
- [19] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Commun. ACM*, vol. 21, pp. 120-126, Feb. 1978.
- [20] F. B. Schneider, "Synchronization in distributed programs," *Trans. Program. Lang. Syst.*, vol. 4, pp. 125-148, Apr. 1982.
- [21] T. K. Srikanth and S. Toueg, "Clock synchronization," Dep. Comput. Sci., Cornell Univ., Ithaca, NY, Tech. Rep. TR 84-656, Dec. 1984.
- [22] A. Tanenbaum, *Computer Networks*. Englewood Cliffs, NJ: Prentice-Hall, 1981.



Özalp Babaoğlu (S'78-M'81) was born in Ankara, Turkey. He received the B.Sc. degree in electrical engineering from George Washington University, Washington, DC, in 1976, and the M.Sc. and Ph.D. degrees in computer science from the University of California, Berkeley, in 1977 and 1981, respectively.

He spent the Summer of 1978 at IBM San Jose Research Laboratory, San Jose, CA, as an Academic Associate, and the first half of 1980 as a Foreign Scholar at the Numerical Analysis Institute of the Italian National Research Council in Pavia, Italy. In 1981 he joined the Department of Computer Science at Cornell University, Ithaca, NY, where he is currently an Assistant Professor. As a graduate student at Berkeley, he was the principal designer and co-implementor of the virtual memory extensions to the UNIX operating system for the VAX computer known as 3.0bsd. His current research interests include fault-tolerance, distributed computing, performance evaluation, and modeling.

Dr. Babaoğlu is a member of the Association for Computing Machinery and the IEEE Computer Society. He was a co-recipient of the 1982 Sakrison Memorial Award for his contribution to the Berkeley UNIX operating system.



Rogério Drummond was born in Campinas, Brazil in 1955. He received the B.Sc. and M.Sc. degrees in computer science from the State University at São Paulo, Campinas, Brazil, in 1978 and 1980, respectively.

He joined the Department of Computer Science at the State University of São Paulo in 1980. In the Spring of 1981 he took a leave of absence from his post at the University of São Paulo to join the graduate program in computer science at Cornell University, Ithaca, NY. He is currently completing the requirements for the Ph.D. degree in computer science. His research interests are in the area of distributed computing, where he is studying the impact of communication architectures on the complexity of distributed algorithms.