

Reaching Approximate Agreement in the Presence of Faults

Danny Dolev¹
Nancy A. Lynch²
Shlomit S. Pinter³
Eugene W. Stark⁴
William E. Weihl⁵

February, 1985

ABSTRACT

This paper considers a variant on the Byzantine Generals problem, in which processes start with arbitrary real values rather than Boolean values or values from some bounded range, and in which approximate, rather than exact, agreement is the desired goal. Algorithms are presented to reach approximate agreement in asynchronous, as well as synchronous systems. The asynchronous agreement algorithm is an interesting contrast to a result of Fischer, Lynch, and Paterson, who show that exact agreement is not attainable in an asynchronous system with as few as one faulty process. The algorithms work by successive approximation, with a provable convergence rate that depends on the ratio between the number of faults and the number of processes. Lower bounds on the convergence rate for algorithms of this form are proven, and the algorithms presented are shown to be optimal.

Keywords: Agreement protocols, fault-tolerance, Byzantine Generals, approximate agreement, distributed systems

©1985 Massachusetts Institute of Technology, Cambridge, MA. 02139

¹Hebrew University, Jerusalem, Israel

²M.I.T., Cambridge, MA. This work was supported in part by the NSF under Grant No. DCR-8302391, U.S. Army Research Office Contract #DAAG29-84-K-0058, and the Defense Advanced Research Projects Agency (DARPA) #N00014-83-K-0125.

³Technion, Haifa, Israel

⁴M.I.T., Cambridge, MA. This work was supported in part by the NSF under Grant No. DCR-8302391, U.S. Army Research Office Contract #DAAG29-84-K-0058, and the Defense Advanced Research Projects Agency (DARPA) #N00014-83-K-0125. Presently at SUNY, Stony Brook.

⁵M.I.T., Cambridge, MA. This work was supported in part by a graduate fellowship from the Fannie and John Hertz Foundation, and the Defense Advanced Research Projects Agency (DARPA) #N00014-83-K-0125.



1. Introduction

In designing fault-tolerant distributed systems, one often encounters questions of agreement among processes. In the Byzantine Generals problem [PSL 80, LSP 82], the objective is for nonfaulty processes to agree on a value, in spite of the presence of a small number of "Byzantine" types of faults — completely arbitrary, even possibly malicious, behavior. Several variations on the problem can be considered — the model can be synchronous or asynchronous, and either exact or approximate agreement can be demanded. In this paper, we consider a variant on the traditional Byzantine Generals problem, in which processes start with arbitrary real values, and where approximate, rather than exact, agreement is the desired goal. Approximate agreement can be used, for example, for clock synchronization and for stabilization of input from sensors.

We assume a model in which processes can send messages containing arbitrary real values, and can store arbitrary real values as well. We assume that each process starts with an arbitrary real value. For any preassigned $\epsilon \geq 0$ (as small as desired), an *approximate agreement algorithm* must satisfy the following two conditions:

- Agreement: All nonfaulty processes eventually halt with output values that are within ϵ of each other.
- Validity: The value output by each nonfaulty process must be in the range of initial values of the nonfaulty processes.

Thus, in particular, if all nonfaulty processes should happen to start with the same initial value, the final values are all required to be the same as the common initial value. This is consistent with the usual requirements for Byzantine agreement algorithms. However, should the nonfaulty processes start with different values, we do not require that the nonfaulty processes agree on a unique final value.

We consider both synchronous and asynchronous versions of the problem. Systems in which there is a finite bounded delay on the operations of the processes and on their intercommunication are said to be synchronous. In such systems, unannounced process deaths, as well as long delays, are considered to be faults. For synchronous systems, we give a simple and rather efficient algorithm for achieving approximate agreement. This algorithm works by successive approximation, with a provable convergence rate that depends on the ratio between the number of faults and the total number of processes. The algorithm is guaranteed to converge in the case where the total number of processes is more than three times the number of possible faults. Termination is achieved using a technique that ensures that all nonfaulty processes halt, but allows different processes to terminate at different times.

For asynchronous systems, in which a very slow process cannot be distinguished from a dead process, no exact agreement can be achieved [FLP 83], even if no malicious failures occur [DDS 83]. An interesting contrast to the results in [FLP 83, DDS 83] is our second algorithm, which enables processes in an asynchronous system to get as close to agreement as one chooses. Our algorithm for the asynchronous case also works by successive approximation. In this case, however, the total number of processes required by the algorithm is more than five times the number of possible faults. As in the synchronous case, we achieve termination using a technique that ensures that all nonfaulty processes halt, but permits different processes to terminate at different times. As we discuss later, it is possible to achieve simultaneous termination in the synchronous case; the technique used for simultaneous termination, however, does not extend to the asynchronous case.

Our algorithms to obtain approximate agreement are of a very simple form. Namely, at each round until termination is reached, each process sends its latest value to all processes (including itself). On receipt of a vector V of values, the process computes a certain function $f(V)$ as its next value. The function f is a kind of averaging function. Here we use functions that are appropriate for handling t faults. We will show that these functions have particularly nice approximation behavior. In particular, we will show that, for algorithms of a particular form, no approximation function can provide uniformly faster convergence than the functions used in this paper. An earlier paper [DLPSW 83] presented similar algorithms, but used approximation functions that provided slower convergence than achieved by the functions used in this paper.

The remainder of this paper is organized as follows: In Section 2, we prove some combinatorial properties of the approximation functions upon which our algorithms depend. Then, in Section 3, we introduce the synchronous model and present the synchronous approximate agreement algorithm, and in Section 4, we present the asynchronous model and algorithm. Next, in Section 5, we present lower bounds on the convergence rate for algorithms of the form presented in sections 3 and 4, and show that the approximation functions used in our algorithms are optimal. In Section 6, we discuss the resilience properties of our algorithms. Finally, in Section 7, we conclude with a short summary and some open questions.

2. Properties of the Approximation Functions

In this section, we will state and prove the relevant properties of the approximation functions. First, we require some preliminary definitions and properties of multisets.

2.1. Preliminary Definitions

Let \mathbb{N} be the natural numbers, including 0, and let \mathbb{R} be the real numbers. We view a finite multiset U of reals as a function $U: \mathbb{R} \rightarrow \mathbb{N}$ that is nonzero on at most finitely many $r \in \mathbb{R}$. Intuitively, the function U assigns a finite multiplicity to each value $r \in \mathbb{R}$. The *cardinality* of a multiset U is given by $\sum_{r \in \mathbb{R}} U(r)$, and is denoted by $|U|$. We say that a multiset is *empty* if its cardinality is zero; otherwise it is *nonempty*. The *difference* $U-V$ of multisets U and V is the multiset W defined by

$$\begin{aligned} W(r) &= U(r)-V(r) \text{ if } U(r)-V(r) \geq 0 \\ &= 0 \text{ otherwise.} \end{aligned}$$

The *intersection* $U \cap V$ of multisets U and V is the multiset W defined by $W(r) = \min(U(r), V(r))$.

In the sequel, the term "multiset" will always refer to finite multisets of real numbers as above. In addition, g^k will denote the k -fold iteration of g ; thus $g^1 = g$, $g^2 = g \circ g$, etc.

The *minimum* $\min(U)$ of a nonempty multiset U is defined by $\min(U) = \min\{r \in \mathbb{R} : U(r) \neq 0\}$. The *maximum* $\max(U)$ is defined similarly. If U is nonempty, let $\rho(U)$ (the *range* of U) be the interval $[\min(U), \max(U)]$, and let $\delta(U)$ (the *diameter* of U) be $\max(U) - \min(U)$. The *mean* $\text{mean}(U)$ of a nonempty multiset U is defined by

$$\text{mean}(U) = \frac{\sum_{r \in \mathbb{R}} r \cdot U(r)}{|U|}.$$

If U is a nonempty multiset, we define the multiset $s(U)$ (intuitively, the multiset obtained by removing one occurrence of the smallest value in U) to be the multiset W defined by

$$\begin{aligned} W(r) &= U(r)-1 \text{ if } r = \min(U), \\ &= U(r) \text{ otherwise.} \end{aligned}$$

The multiset $l(U)$ (remove one occurrence of the largest value in U) is defined similarly. Assume t is a fixed nonnegative integer. If $|U| \geq 2t$, then define $\text{reduce}(U) = s^t(l^t(U))$, the result of removing the t highest and t lowest elements of U .

The first lemma shows that the number of common elements in two nonempty multisets is reduced by at most 1 when the smallest (or the largest) element is removed from each.

Lemma 1: Suppose that V and W are nonempty multisets. Then

$$1. |V \cap W| - |s(V) \cap s(W)| \leq 1.$$

$$2. |V \cap W| - |l(V) \cap l(W)| \leq 1.$$

Proof: We prove the first inequality; the argument for the second is symmetric. Let $M = V \cap W$, and let $N = s(V) \cap s(W)$. Let $v = \min(V)$ and $w = \min(W)$. Now, $N(r) = M(r)$ if $r \neq v$ and $r \neq w$, so we can write $|M| - |N| = \sum_{r \in \{v, w\}} M(r) - \sum_{r \in \{v, w\}} N(r)$. There are two cases,

depending upon whether $v = w$ or $v \neq w$.

Suppose $v = w$. Then $|M| - |N| = M(v) - N(v) = \min(V(v), W(v)) - \min(V(v)-1, W(v)-1) \leq 1$.

Suppose $v \neq w$. Then

$$\begin{aligned} |M| - |N| &= (M(v) + M(w)) - (N(v) + N(w)) \\ &= (\min(V(v), W(v)) + \min(V(w), W(w))) \\ &\quad - (\min(V(v)-1, W(v)) + \min(V(w), W(w)-1)). \end{aligned}$$

Assume without loss of generality that $v < w$. Then $W(v) = 0$, and hence

$$\begin{aligned} |M| - |N| &= \min(V(w), W(w)) - \min(V(w), W(w)-1) \\ &\leq 1. \quad \square \end{aligned}$$

The next lemma extends the results of the previous lemma to removing the t largest and t smallest elements.

Lemma 2: Suppose that V and W are multisets such that $|V| \geq 2t$ and $|W| \geq 2t$. Then

$$|V \cap W| - |\text{reduce}(V) \cap \text{reduce}(W)| \leq 2t.$$

Proof: Follows from repeated application of Lemma 1. \square

The next lemma is fundamental to the correctness of the algorithms. It states that if at most kt values in a multiset V are not in another multiset U , then every value in $\text{reduce}^k(V)$ is in the range of U . For example, if the multiset of values held by nonfaulty processes at some point in the algorithm is U , and the multiset of values received by some process is V , then at most t of the values in V are not in U . The lemma then states that reducing V produces a multiset whose range is contained in the range of the values of the nonfaulty processes. This property is essential in showing that the validity condition is satisfied.

Lemma 3: Suppose that k is a nonnegative integer and that U and V are nonempty multisets such that $|V - U| \leq kt$ and $|V| > 2kt$. Then $\rho(\text{reduce}^k(V)) \subseteq \rho(U)$.

Proof: Suppose $\rho(\text{reduce}^k(V)) \not\subseteq \rho(U)$. Then either $\min(\text{reduce}^k(V)) < \min(U)$ or $\max(\text{reduce}^k(V)) > \max(U)$. Both cases lead to a contradiction. We argue the first; the second is symmetric.

If $\min(\text{reduce}^k(V)) < \min(U)$, then $\sum_{r < \min(U)} V(r) \geq kt + 1$. Hence, $|V - U| \geq kt + 1$, which contradicts a hypothesis. \square

2.2. The Approximation Functions

If U is a nonempty multiset with $|U| = m$, $u_0 \leq u_1 \leq \dots \leq u_{m-1}$ are the elements of U in nondecreasing order, and $i > 0$ and $j \geq 0$ are integers such that $ji + 1 \leq m \leq (j+1)i$ (i.e., $j = \text{floor}((m-1)/i)$), we define $\text{select}_i(U)$ to be the multiset consisting of the elements u_0, u_i, u_{2i}, \dots , and u_{ji} . Thus, $\text{select}_i(U)$ chooses the smallest element of U and every i^{th} element thereafter.

In this paper we will use two related approximation functions, one for the synchronous protocol (f_S) and one for the asynchronous protocol (f_A). Let V be a finite multiset of reals with $|V| > 2t$. The approximation functions are defined as follows:

$$f_S(V) = \text{mean}(\text{select}_i(\text{reduce}(V)))$$

$$f_A(V) = \text{mean}(\text{select}_{2t}(\text{reduce}(V)))$$

For notational convenience, we define the function f_i , for a positive integer i , by $f_i(V) = \text{mean}(\text{select}_i(\text{reduce}(V)))$. Thus, $f_S = f_i$ and $f_A = f_{2t}$.

The next two lemmas describe properties of the approximation functions. Lemma 4 is used in verifying the validity condition.

Lemma 4: Suppose that U and V are nonempty multisets such that $|V-U| \leq t$ and $|V| > 2t$, and i is a positive integer. Then $f_i(V) \in \rho(U)$.

Proof: Follows easily from Lemma 3 (with $k = 1$). \square

Lemma 5 will be used in the following manner: U is the multiset of values held by nonfaulty processes in a given round, and M and N are the reduced multisets of values received by two nonfaulty processes in that round. Nonfaulty processes use the appropriate approximation function to choose their values for the next round; the lemma tells us how quickly those values converge.

Lemma 5: Suppose U , M , and N are nonempty multisets, and $m > 0$, $i > 0$, and $j \geq 0$ are integers such that:

$$ji + 1 \leq m \leq (j+1)i,$$

$$|M| = |N| = m,$$

$$|M \cap N| \geq m - i,$$

$$\rho(M) \subseteq \rho(U) \text{ and } \rho(N) \subseteq \rho(U).$$

Then $|\text{mean}(\text{select}_i(M)) - \text{mean}(\text{select}_i(N))| \leq \delta(U)/(j+1)$.

Proof: From the definition of select_i , we know that $\text{select}_i(M)$ and $\text{select}_i(N)$ each contain exactly $j+1$ elements. Let $m_0 \leq m_1 \leq \dots \leq m_j$ be the elements of $\text{select}_i(M)$, and let $n_0 \leq n_1 \leq \dots \leq n_j$ be the elements of $\text{select}_i(N)$. Notice that there are at least $ik+1$ elements in M that are less than or equal to m_k , and at most ik elements in M that are strictly less than m_k ; similarly for N .

We begin by showing that $\max(m_k, n_k) \leq \min(m_{k+1}, n_{k+1})$ for $0 \leq k \leq j-1$. It suffices to show that $m_k \leq n_{k+1}$; a symmetric argument demonstrates that $n_k \leq m_{k+1}$.

We proceed by contradiction: Suppose that $m_k > n_{k+1}$. As noted above, there are at least $i(k+1)+1$ elements in N less than or equal to n_{k+1} . By our supposition, these elements are strictly less than m_k . However, there are at most ik elements in M strictly less than m_k . Therefore, there are at least $i(k+1)+1-ik (=i+1)$ elements in N that are not in M ; thus, $|N-M| \geq i+1$. This contradicts the hypothesis that $|M \cap N| \geq m-i$. Thus, $m_k \leq n_{k+1}$.

Now we will use the inequality shown above to obtain the desired result. Using the notation defined above,

$$\begin{aligned} & |\text{mean}(\text{select}_i(M)) - \text{mean}(\text{select}_i(N))| \\ &= 1/(j+1) |(\sum_{k=0}^j m_k) - (\sum_{k=0}^j n_k)| \\ &= 1/(j+1) |\sum_{k=0}^j (m_k - n_k)| \\ &\leq 1/(j+1) \sum_{k=0}^j |m_k - n_k| \quad (\text{by the triangle inequality}) \\ &= 1/(j+1) \sum_{k=0}^j (\max(m_k, n_k) - \min(m_k, n_k)). \end{aligned}$$

By the inequality demonstrated above, for $0 \leq k \leq j-1$, $(\max(m_k, n_k) - \min(m_k, n_k)) \leq (\min(m_{k+1}, n_{k+1}) - \min(m_k, n_k))$, so we get

$$\begin{aligned} & |\text{mean}(\text{select}_i(M)) - \text{mean}(\text{select}_i(N))| \\ &\leq 1/(j+1) [\max(m_j, n_j) - \min(m_j, n_j)] \\ &\quad + 1/(j+1) \sum_{k=0}^{j-1} (\min(m_{k+1}, n_{k+1}) - \min(m_k, n_k)). \end{aligned}$$

Collecting terms then shows that

$$\begin{aligned} & |\text{mean}(\text{select}_i(M)) - \text{mean}(\text{select}_i(N))| \\ &\leq 1/(j+1) (\max(m_j, n_j) - \min(m_0, n_0)). \end{aligned}$$

Since $\rho(M) \subseteq \rho(U)$ and $\rho(N) \subseteq \rho(U)$, we know that $\max(m_j, n_j) \leq \max(U)$ and $\min(m_0, n_0) \geq \min(U)$. This gives the desired result:

$$\begin{aligned} & |\text{mean}(\text{select}_i(M)) - \text{mean}(\text{select}_i(N))| \\ &\leq 1/(j+1) (\max(U) - \min(U)) \\ &= \delta(U)/(j+1). \quad \square \end{aligned}$$

3. The Synchronous Problem

A *synchronous approximation algorithm* P is a system of n processes, $n \geq 1$. Each process p has a set of states, including a subset of states called *initial states* and a subset called *halting states*. There is a *value mapping* which assigns a real number as the value of each state. For each real number r , there is exactly one initial state with value r . Each process acts deterministically according to a *transition function* and a *message generation function*. The transition function takes a non-halting process state and a vector of messages received from all processes (one message per process) and produces a new process state. The message generation function takes a non-halting state and produces a vector of messages to be sent to all processes (one per process).

We assume that the system acts synchronously, using a reliable communication medium. Each process is able to send messages to all processes (including itself), and the sender of each message is identifiable by the receiver.

A *configuration* consists of a state for each process. An *initial configuration* consists of an initial state for each process. Let T be any subset of the processes. A sequence of configurations (called *rounds*), C_0, C_1, C_2, \dots is a *T-computation* provided there exist messages sent by each process at each round such that: (a) C_0 is an initial configuration; (b) for every i , and every $p \in T$, the messages sent out by p after C_i are exactly those specified by p 's message generation function, applied to p 's state in C_i ; and (c) for every i , and every $p \in T$, p 's state in C_{i+1} is exactly the one specified by p 's transition function applied to p 's state in C_i and the messages sent to p after C_i . In a T -computation, processes in T are nonfaulty, while processes not in T may be faulty.

For the rest of the paper, assume a fixed small value ϵ , a fixed number of processes n , and a fixed maximum number of faults t .

A synchronous approximation algorithm is said to be *t-correct* provided that for every subset T of processes with $|T| \geq n-t$, and every T -computation, the following is true: Every $p \in T$ eventually enters a halting state, and the following two conditions hold for the values of those halting states:

- Agreement: If two processes in T enter halting states with values r and s , respectively, then $|r - s| \leq \epsilon$.
- Validity: If a process in T enters a halting state with value r , then there exist processes in T having x and y as initial values, such that $x \leq r \leq y$.

We will prove the following theorem.

Theorem 6: If $n \geq 3t + 1$, then there exists a t -correct synchronous approximation

algorithm with n processes.

Note that the following strategy would suffice to prove Theorem 1. The processes could run n executions of a general (unlimited value set) Byzantine Generals algorithm such as the one in [DS 82], in order to obtain common estimates for the initial values of all the processes. After this algorithm completes, all processes in T will have the same multiset, v , of values for all the processes. Then each process halts with value $f(v)$, where f is a predetermined averaging function that is the same for all processes. This algorithm actually achieves exact real-valued agreement, with the required validity condition. However, the solution presented below is much simpler and more elegant, and moreover extends directly to the asynchronous case, for which exact agreement is impossible. The algorithm has two additional advantages over using a Byzantine Generals algorithm: It is more resilient than typical Byzantine Generals algorithms, and it can, in some cases, terminate in fewer than $t + 1$ rounds.

We now present our synchronous approximation algorithm, S . First, we describe a nonterminating algorithm, S_0 , and then we discuss how termination is achieved. We assume that $n \geq 3t + 1$.

Synchronous Approximation Algorithm S_0 :

At each round, each nonfaulty process p performs the following steps:

1. Process p broadcasts its current value to all processes, including itself.
2. Process p collects all the values sent to it at that round into a multiset V . If p does not receive exactly one correct value from some particular other process (which means, in the synchronous model, that the other process is faulty), then p simply picks some arbitrary default value to represent that process in the multiset. The multiset V will therefore always contain exactly n values.
3. Process p applies the function f_S to the multiset V to obtain its new value.

The following result states how the diameter and range of the nonfaulty processes' values are affected by each round of algorithm S_0 .

Lemma 7: Let $j \geq 1$ and $t, n \geq 0$ be such that $(j + 2)t + 1 \leq n \leq (j + 3)t$. Let T be a set of processes, with $|T| \geq n - t$. Let h be a positive integer. Let U and U' be the multisets of values of processes in T , immediately before and after round h , respectively, in a particular T -computation of S_0 . Then

1. $\delta(U') \leq \delta(U)/(j + 1)$.
2. $\rho(U') \subseteq \rho(U)$.

Proof: Let p and q be arbitrary processes in T . Let V and W be the multisets of values

(including default values) received by p and q , respectively, at round h . Since there are at most t faulty processes, $|V - U| \leq t$ and $|W - U| \leq t$. Moreover, since V and W contain identical entries for all the processes in T , we know that $|V \cap W| \geq n-t$.

1. Let $m = n-2t$ and let $i = t$. Let $M = \text{reduce}(V)$ and $N = \text{reduce}(W)$. By Lemma 2, $|M \cap N| \geq |V \cap W| - 2t \geq n-3t = m-i$. By Lemma 3 (with $k = 1$), $\rho(M) \subseteq \rho(U)$ and $\rho(N) \subseteq \rho(U)$. Thus, U, M, N, m , and i satisfy the hypotheses of Lemma 5, implying that $|f_S(V) - f_S(W)| \leq \delta(U)/(j+1)$. Since p and q were chosen arbitrarily, the result follows.

2. Follows directly from Lemma 4.

□

Lemma 7 shows that, at each round, the diameter of the multiset of values held by nonfaulty processes decreases by a factor of $1/(j+1)$, where $(j+2)t+1 \leq n \leq (j+3)t$ and $j \geq 1$. Thus, the diameter of the multiset of values held by nonfaulty processes eventually decreases to ϵ or less. In addition, repeated application of part 2 of Lemma 7 shows that, at each round $h \geq 1$, the values held by nonfaulty processes immediately before round h are all in the range of the initial values of nonfaulty processes.

Algorithm S_0 is not a correct synchronous approximation algorithm, however, for as stated, it never terminates. We modify S_0 to obtain a terminating algorithm, S , as follows. At the first round, each nonfaulty process uses the range of all the values it has received at that round to compute a round number at which it is sure that the values of any two nonfaulty processes will be at most ϵ apart. Each process can do this because it knows the value of ϵ , the guaranteed rate of convergence and furthermore, it knows that the range of values it receives on the first round includes the initial values of all nonfaulty processes.

In general, different processes might compute different round numbers. Any process that reaches its computed round simply halts, and sends its value out with a special "halting" tag. When any process, say p , receives a value with a "halting" tag, it knows to use the enclosed value not only for the designated round, but also for all future rounds (until p itself decides to halt, based on p 's own computed round number). Although nonfaulty processes might compute different round numbers, it is clear that the smallest such estimate is correct. Thus, at the time the first nonfaulty process halts, the range is already sufficiently small. At subsequent rounds, the range of values of nonfaulty processes is never increased, although we can no longer guarantee that it decreases. The following lemma makes these ideas more precise.

Lemma 8: Assume that $n \geq 3t + 1$. Let T be a set of processes, with $|T| \geq n-t$. Let h be

a positive integer. Let U and U' be the multisets of values of processes in T , immediately before and after round h , respectively, in a particular T -computation of S . Then $\rho(U') \subseteq \rho(U)$.

Proof: Let p be an arbitrary process in T . Let v and v' be the values held by p immediately before and after round h , respectively. It suffices, since p is arbitrary, to show that $v' \in \rho(U)$. If p has terminated prior to the start of round h , then $v' = v \in \rho(U)$. If p has not halted prior to the start of round h , then let V be the multiset of values received by p in round h . Then V and U satisfy the hypotheses of Lemma 4, and since $v' = f_S(V)$, it follows that $v' \in \rho(U)$. \square

It is now easy to see that S is a correct synchronous approximation algorithm. It is clear that all processes terminate. Consider the agreement property. At the first round at which some process halts, it is already the case that all nonfaulty processes' values are within ϵ of each other. By Lemma 8, this diameter never increases at subsequent rounds, so the final values of all the nonfaulty processes are also within ϵ of each other. The validity property also follows from repeated application of Lemma 8. This completes the proof of Theorem 6.

It is possible to modify algorithm S in several interesting ways. First, observe that a process need not always wait for its computed round to arrive before halting: it can halt after it receives "halting" tags from at least $t + 1$ other processes.

Second, it is possible to have all nonfaulty processes halt simultaneously. This can be done using known solutions to the binary-valued Byzantine Generals problem as follows: After certain (a priori) selected rounds of the approximate agreement algorithm, all processes run a Byzantine Generals algorithm to decide whether the algorithm should continue. A process p will vote to halt only if the round number has reached p 's computed round number. Voting to halt implies that p knows that the range of values of nonfaulty processes is at most ϵ . If the decision is to halt, then some nonfaulty process must have voted to halt, and the range of values of the nonfaulty processes must be at most ϵ . (Note that, since exact agreement is impossible in an asynchronous model, this termination technique does not extend to the asynchronous case.)

An alternative way of guaranteeing that all nonfaulty processes halt at the same round uses an arbitrary solution to the Byzantine Firing Squad problem [BL 85, CDDS 85].

4. The Asynchronous Problem

In this section, we reformulate the problem in an asynchronous model adapted from the one in [FLP 83]. In an *asynchronous approximation algorithm*, we assume that processes have states as before, but now the operation of the processes is described by a transition function that in one step tries to receive a message, gets back either "null" or an actual message, and based on the message, changes state and sends out a finite number of other messages. Nonfaulty processes always follow the algorithm. Faulty processes, on the other hand, are constrained so that their steps at least follow the standard form — in each step, they try to receive a message as do nonfaulty processes. However, they can change state arbitrarily (not necessarily according to the given algorithm), and can send out any finite set of messages (not necessarily the ones specified by the algorithm). A *T-computation* of an asynchronous approximation algorithm is one in which the processes in T always follow the algorithm, all processes (faulty and nonfaulty) continue to take steps until they reach a halting state, and any process that fails to enter a halting state eventually receives all messages sent to it.

An asynchronous approximation algorithm is said to be *t-correct* provided for every subset T of processes with $|T| \geq n-t$, and every T -computation, every process in T eventually halts, and the same agreement and validity conditions hold as for the synchronous case.

It seems simplest here to insist on the standard form being followed by all processes. The requirement that faulty processes keep taking steps until they enter halting states is not a restriction, since they are free to enter halting states at any time they wish. Similarly, the requirement that faulty processes continue trying to receive messages is not a restriction, since they are free to do whatever they like with the messages received. Finally, the requirement that faulty processes only send finitely many messages at each step is needed so that faulty processes are unable to flood the message system, preventing messages from other processes from getting through.

We assume that processes take steps at completely arbitrary rates, so that there is no way (in finite time) to distinguish a faulty process from one that is simply slow in responding. Also, we assume that the message system takes arbitrary lengths of time to deliver messages, and delivers them in arbitrary order.

We will prove the following theorem:

Theorem 9: If $n \geq 5t + 1$, then there exists a t -correct asynchronous approximation algorithm with n processes.

We now describe the asynchronous approximation algorithm. As in the synchronous case, first we

describe a nonterminating algorithm, A_0 , in which processes compute better and better approximations, and we then modify A_0 to produce a terminating algorithm A . Assume that $n \geq 5t + 1$.

Asynchronous Approximation Algorithm A_0

At round h , each nonfaulty process p performs the following steps:

1. Process p labels its current value with the current round number h , and then broadcasts this labeled value to all processes, including itself.
2. Process p waits to receive exactly $n-t$ round h values, and collects these values into a multiset V . Since there can be at most t faulty processes, process p will eventually receive at least $n-t$ round h values. Note that, in contrast to the synchronous case, process p does *not* choose any default values.
3. Process p applies the function f_A to the multiset V to obtain its new value.

In analogy with Lemma 7, we have the following result, which states the convergence properties of the above algorithm.

Lemma 10: Let $j \geq 1$, and $t, n \geq 0$, be such that $(2j+3)t+1 \leq n \leq (2j+5)t$. Let T be a set of processes, with $|T| \geq n-t$. Let h be a positive integer. Let U and U' be the multisets of values of processes in T , immediately before and after round h , respectively, in a particular T -computation of A_0 . Then

1. $\delta(U') \leq \delta(U)/(j+1)$.

2. $\rho(U') \subseteq \rho(U)$.

Proof: Let p and q be arbitrary processes in T . Let V and W be the multisets of values received by p and q , respectively, at round h . Since there are at most t faulty processes, $|V - U| \leq t$ and $|W - U| \leq t$. Moreover, since V and W both contain identical entries for all the processes in T from which both p and q heard, we know that $|V \cap W| \geq n - 3t$.

1. Let $m = n - 3t$ and let $i = 2t$. Let $M = \text{reduce}(V)$ and $N = \text{reduce}(W)$. By Lemma 2, $|M \cap N| \geq |V \cap W| - 2t \geq n - 5t = m - i$. By Lemma 3 (with $k = 1$), $\rho(M) \subseteq \rho(U)$ and $\rho(N) \subseteq \rho(U)$. Thus, U, M, N, m , and i satisfy the hypotheses of Lemma 5, implying that $|f_A(V) - f_A(W)| \leq \delta(U)/(j+1)$. Since p and q were chosen arbitrarily, the result follows.

2. Follows directly from Lemma 4.

□

Lemma 10 shows that, at each round, the diameter of the multiset of values of nonfaulty processes

decreases by a factor of $1/(j+1)$, where $(2j+3)t+1 \leq n \leq (2j+5)t$ and $j \geq 1$. Thus, the diameter of the multiset of values held by nonfaulty processes eventually decreases to ϵ or less. In addition, repeated application of part 2 of Lemma 10 shows that, at each round $h \geq 1$, the values held by nonfaulty processes immediately before round h are all in the range of the initial values of nonfaulty processes.

The only remaining problem is termination. We cannot use the same technique that we used in the synchronous algorithm, because a process cannot wait until it hears from all other processes, and thus cannot obtain an estimate of the range of the initial values of the nonfaulty processes. We solve this problem by adding an initialization round at the beginning of the algorithm. In this initialization round (round 0), each nonfaulty process p performs the following steps:

Initialization Round for Asynchronous Approximation Algorithm A:

1. Process p labels its initial value with the round number 0, and then broadcasts this labeled value to all processes, including itself.
2. Process p waits to receive exactly $n-t$ round 0 values, and collects these values into a multiset V_p .
3. Process p chooses an arbitrary element of $\rho(\text{reduce}^2(V_p))$ (for definiteness, say $\text{mean}(\text{reduce}^2(V_p))$) as its initial value for use in round 1. Let x_p be this chosen value.

Suppose that p and q are arbitrary nonfaulty processes. Then since $|V_p| \geq 4t$ and $|V_p - V_q| \leq 2t$, it follows that V_p and V_q satisfy the hypotheses for the multisets V and U , respectively, in Lemma 3 (with $k=2$). An application of this result therefore shows that, for any nonfaulty processes p and q , it is the case that $x_p \in \rho(V_q)$. That is, the value x_p computed by process p as the result of the initialization round is contained in the range of all values received by process q in the initialization round. Since each nonfaulty process q knows: (1) that its range $\rho(V_q)$ contains all the round 1 values x_p for nonfaulty processes p ; (2) the value of ϵ ; and (3) the guaranteed rate of convergence, it can compute, before the beginning of round 1, a round number at which it is sure that the values of any two nonfaulty processes will be at most ϵ apart.

As in the synchronous case, different processes will calculate different round numbers at which they would like to halt. The same modification, of sending a value out with a special "halting" tag, works here as well. We obtain a lemma analogous to Lemma 8:

Lemma 11: Assume that $n \geq 5t + 1$. Let T be a set of processes, with $|T| \geq n-t$. Let h be a positive integer. Let U and U' be the multisets of values of processes in T , immediately before and after round h , respectively, in a particular T -computation of A .

Then $\rho(U') \subseteq \rho(U)$.

The remainder of the proof of Theorem 9 is analogous to that of Theorem 6.

5. Lower Bound Results

In this section, we assume that algorithms are of a standard form in which at each round, an old approximation is exchanged with other processes, and a new approximation is computed from the multiset of values received, by the application of an approximation function f . We assume that f is *cautious*, as defined below. (Our algorithms all fit this pattern.) The results show that, under these assumptions, the function f_S gives the best possible convergence factor for a synchronous algorithm for $n \geq 3t + 1$, and the function f_A gives the best possible convergence factor for an asynchronous algorithm for $n \geq 5t + 1$.

In [DLPSW 83], an earlier version of this work, we used different averaging functions in our algorithms. The discovery of the lower bounds in this section suggested that those functions did not give optimal rates of convergence, and led us to search for the improved averaging functions which appear in this paper.

In the remainder of this section, let n and t be fixed.

We say that an approximation function f , which takes a multiset M of real numbers to a real number $f(M)$, is *cautious* if $f(M) \in \rho(U)$ for all multisets U such that $|M - U| \leq t$. The cautious requirement seems reasonable for any approximation function that will tolerate up to t faults: regardless of the values received from the faulty processes, a cautious function will produce a value in the range of the values held by the nonfaulty processes. It is easy to see that f_S and f_A are cautious.

5.1. The Synchronous Problem

We will show the following theorem:

Theorem 12: Suppose $j \geq 0$, $t \geq 1$, and $(j + 2)t + 1 \leq n \leq (j + 3)t$. Suppose that f and g are cautious approximation functions. Then there exist multisets V , W , and U such that:

$$|V| = |W| = n,$$

$$|U| = n - t,$$

$$|V - U| = |W - U| = n - t, \text{ and}$$

$$|f(V) - g(W)| \geq \delta(U)/(j+1).$$

The implications of this result for the synchronous agreement algorithm are the following: Suppose we consider algorithms of a standard form, in which at each round, a process exchanges its current approximation with all other processes, and then applies a cautious approximation function to the multiset of values it receives to determine its new approximation. Theorem 12 then implies that there exist multisets V , W , and U , such that if correct processes p and q (using approximation functions f and g , respectively) receive multisets of values V and W , respectively, in some round of execution, and U is the multiset of values held by correct processes at the start of that round, then the new approximations held by p and q at the end of the round can be no closer than $\delta(U)/(j+1)$. Thus this result yields a fundamental limitation on the rate of convergence of algorithms of the standard form. The lower bound given by this result also matches the upper bound provided by the function f_S .

The proof of Theorem 12 requires the following lemma, which asserts the existence of a $j+1$ -link chain of multisets that span from a multiset M_1 upon which every cautious approximation function must yield 0, to a multiset M_{j+2} upon which every cautious approximation function must yield 1.

Lemma 13: Suppose $j \geq 0$, $t \geq 1$, and $(j+2)t+1 \leq n \leq (j+3)t$. Then there exist multisets M_1, M_2, \dots, M_{j+2} , and U_1, U_2, \dots, U_{j+1} , such that:

$$|M_i| = n \text{ for } 1 \leq i \leq j+2,$$

$$|U_i| = n-t \text{ for } 1 \leq i \leq j+1,$$

$$\delta(U_i) = 1 \text{ for } 1 \leq i \leq j+1,$$

$$|M_i - U_i| = |M_{i+1} - U_i| = n-t \text{ for } 1 \leq i \leq j+1, \text{ and}$$

$$f(M_1) = 0 \text{ and } f(M_{j+2}) = 1,$$

whenever f is a cautious approximation function.

Proof: Define M_i to have the value 0 with multiplicity $n-it$ and the value 1 with multiplicity it . Define U_i to have the value 0 with multiplicity $n-(i+1)t$ and the value 1 with multiplicity it . The cardinality and diameter constraints on these sets are easily checked. Suppose f is cautious. Then since M_1 has the value 0 with multiplicity $n-t$ ($> t$) and the value 1 with multiplicity t ($\leq t$), it follows that $f(M_1) = 0$. Also, since M_{j+2} has the value 0 with multiplicity $n-(j+2)t$ ($\leq t$) and the value 1 with multiplicity $(j+2)t$ ($> t$), it follows that $f(M_{j+2}) = 1$. \square

We can now present the proof of Theorem 12:

Proof: For $1 \leq i \leq j+2$, let the approximation function h_i be f if i is odd, and g if i is even.

By Lemma 13, there exists a chain M_1, M_2, \dots, M_{j+2} , and U_1, U_2, \dots, U_{j+1} , such that:

$$|M_i| = n \text{ for } 1 \leq i \leq j+2,$$

$$|U_i| = n-t \text{ for } 1 \leq i \leq j+1,$$

$$\delta(U_i) = 1 \text{ for } 1 \leq i \leq j+1,$$

$$|M_i - U_i| = |M_{i+1} - U_i| = n-t \text{ for } 1 \leq i \leq j+1, \text{ and}$$

$$h_1(M_1) = 0 \text{ and } h_{j+2}(M_{j+2}) = 1.$$

Suppose, to obtain a contradiction, that $|h_{i+1}(M_{i+1}) - h_i(M_i)| < \delta(U_i)/(j+1)$ for $1 \leq i \leq j+1$.

Then

$$\begin{aligned} 1 &= |h_{j+2}(M_{j+2}) - h_1(M_1)| \\ &= |h_{j+2}(M_{j+2}) - h_{j+1}(M_{j+1}) + h_{j+1}(M_{j+1}) - h_j(M_j) + \dots + h_2(M_2) - h_1(M_1)| \\ &\leq |h_{j+2}(M_{j+2}) - h_{j+1}(M_{j+1})| + |h_{j+1}(M_{j+1}) - h_j(M_j)| + \dots + |h_2(M_2) - h_1(M_1)| \\ &< (j+1)/(j+1) \\ &= 1. \end{aligned}$$

This is a contradiction, and we conclude that $|h_{i+1}(M_{i+1}) - h_i(M_i)| \geq \delta(U_i)/(j+1)$ for some i with $1 \leq i \leq j+1$. If i is odd, then $h_i = f$ and $h_{i+1} = g$, so letting $V = M_i$, $W = M_{i+1}$, and $U = U_i$ satisfies the requirements of the theorem. If i is even, then instead let $V = M_{i+1}$, $W = M_i$, and $U = U_i$. \square

5.2. The Asynchronous Problem

We will show the following theorem:

Theorem 14: Suppose $j \geq 0$, $t \geq 1$, and $(2j+3)t+1 \leq n \leq (2j+5)t$. Suppose that f and g are cautious approximation functions. Then there exist multisets V , W , and U such that:

$$|V| = |W| = n-t,$$

$$|U| = n-t,$$

$$|V-U| = |W-U| = n-t, \text{ and}$$

$$|f(V) - g(W)| \geq \delta(U)/(j+1).$$

The implications of this result for the asynchronous agreement algorithm are analogous to what Theorem 12 has to say about the synchronous algorithm: There exist multisets V , W , and U , such that if correct processes p and q (using approximation functions f and g , respectively) receive multisets of values V and W , respectively, in some round of execution, and U is the multiset of values held by correct processes at the start of that round, then the new approximations held by p and q at the end of the round can be no closer than $\delta(U)/(j+1)$. The lower bound given by this result also matches the upper bound provided by the function f_A .

As before, the theorem is proved with the aid of the following "chain lemma."

Lemma 15: Suppose $j \geq 0$, $t \geq 1$, and $(2j+3)t+1 \leq n \leq (2j+5)t$. Then there exist multisets M_1, M_2, \dots, M_{j+2} , and U_1, U_2, \dots, U_{j+1} , such that:

$$|M_i| = n-t \text{ for } 1 \leq i \leq j+2,$$

$$|U_i| = n-t \text{ for } 1 \leq i \leq j+1,$$

$$\delta(U_i) = 1 \text{ for } 1 \leq i \leq j+1,$$

$$|M_i - U_i| = |M_{i+1} - U_i| = n-t \text{ for } 1 \leq i \leq j+1, \text{ and}$$

$$f(M_1) = 0 \text{ and } f(M_{j+2}) = 1, \text{ whenever } f \text{ is a cautious approximation function.}$$

Proof: We split the proof into two cases: In case $(2j+4)t+1 \leq n \leq (2j+5)t$, then define M_i to have the value 0 with multiplicity $n-2it$ and the value 1 with multiplicity $(2i-1)t$, for each i with $1 \leq i \leq j+2$. Define U_i to have the value 0 with multiplicity $n-(2i+1)t$ and the value 1 with multiplicity $2it$, for each i with $1 \leq i \leq j+1$. In case $(2j+3)t+1 \leq n \leq (2j+4)t$, then we modify slightly the definition of M_{j+2} and U_{j+1} from the preceding case. That is, define M_{j+2} to have the value 0 with multiplicity $n-(2j+3)t$ and the value 1 with multiplicity $2(j+1)t$. Also, define U_{j+1} to have the value 0 with multiplicity $n-2(j+1)t$ and the value 1 with multiplicity $(2j+1)t$.

In both cases it is straightforward to check that the required properties hold. \square

The proof of Theorem 14 is entirely analogous to the proof of Theorem 12.

6. Resilience

The algorithms presented in this paper have some interesting resilience properties, stronger than those usually claimed for Byzantine agreement algorithms. So far, we have only claimed that the algorithms are resilient to t different processes exhibiting Byzantine faults during the entire course of the algorithm. However, we can claim more for situations where processes fail and recover repeatedly. Our algorithms actually support resilience to any t Byzantine faulty processes at a time (under suitable definitions of faultiness at a particular time); the total number of faulty processes can be much greater than t , since we can allow different processes to be faulty at different times.

We do not give a formal presentation of our resilience properties. Rather, we just give a brief sketch of the main ideas.

First, consider the synchronous case. A faulty process is able to recover easily and reintegrate itself into the algorithm. It can reenter the algorithm at any round, just by sending an arbitrary value, collecting values and averaging them as usual to get a new value. The process also needs to obtain an estimate of the number of rounds required before termination. It can obtain such an estimate in the reentry round, just as it could in the first round.

The asynchronous case is a little more complicated. A faulty process p needs to rejoin the algorithm at some particular (asynchronous) round; however, it must be careful to rejoin at some round that is not "out of date." That is, in the absence of additional failures of p , it must be guaranteed to receive all of its messages for that and subsequent rounds. Process p could not simply wait until it received $n-t$ messages for some particular round k , since those messages might have been delivered very late, and messages for round $k+1$ might have already been lost. However, it suffices for p to send out a "recovery" message, and await acknowledgements from $n-t$ processes carrying the number of their current round. Process p knows that the $t+1^{\text{st}}$ smallest of these round numbers, plus 1, is an allowable round number for it to use for reentry.

The recovering process is not able to use the same method of estimating a termination round as it did initially. Therefore, it seems necessary to modify the asynchronous algorithm to enable recovering processes to obtain termination estimates when needed. An easy modification that works is to have every process piggyback its estimate of the number of rounds to termination on every message it sends. Then a recovering process can obtain a new estimate just by taking the $t+1^{\text{st}}$ smallest of the estimates it receives at the reentry round.

7. Summary and Open Questions

We have defined a problem of approximate agreement on real numbers by processes in a distributed system. We integrated simple approximation functions into two simple-to-implement algorithms for achieving approximate agreement — one for a synchronous distributed system, and the other for an asynchronous system. In addition, we showed that both algorithms achieve the fastest possible convergence rate for algorithms of a particular form. The algorithm for an asynchronous system provides an interesting contrast to the results in [FLP 83, DDS 83], which show that exact agreement is impossible in an asynchronous system.

The ideas of this paper have been used in the design of algorithms for synchronizing clocks in distributed systems [LL 84].

For the synchronous case, it is not difficult to show that $3t + 1$ processes are necessary to solve the approximate agreement problem. The proof is an adaptation of the lower bound proof in [LSP 82], and appears in [FLM 85]. For the asynchronous case, our number of processes is not optimal. In fact, it appears possible to reduce the number of processes to as few as $3t + 1$. This reduction is obtained using a more complex algorithm, based on some of the interesting ideas of [B 84]. This algorithm has a slower rate of convergence than ours.

To obtain the lower bound results, we had to restrict our attention to algorithms of a standard form (ones that operate by broadcasting values and using received values to compute a new approximation), and to functions with a natural, but apparently restrictive property (the "cautious" property). It would be interesting to obtain answers to the following questions:

- Can the cautious property be weakened or removed entirely?
- Can algorithms not of the standard form considered here produce agreement faster?

We would also like to have a better understanding of the relationship between the number of processes and the rate of convergence for approximate agreement algorithms. For instance, the more complex asynchronous algorithm mentioned above uses fewer processes, but has a slower rate of convergence than ours. Is there a tradeoff?

We can state a variant of the approximate agreement problem which uses a fixed number, k , of rounds, and in which ϵ is not predetermined. Each process starts with a real value, as before. After k rounds, the processes must output their final values. The validity condition is the same as before. The object of the algorithm is to insure the best possible agreement, expressed as a ratio of the new

diameter of the nonfaulty processes' values to the original diameter. For given n , t and k , we would like to know the best ratio.

As before, if the algorithm is constrained to operate round-by-round, applying cautious functions at each round, we obtain lower bounds which are exactly the same as are achieved by our averaging functions. However, if the algorithm is unconstrained, the best bounds we have are not at all tight. Consider the synchronous case, for example. The best upper bound we have still arises from repeated application of our averaging function f_g , and is approximately $(t/n)^k$. We can obtain a lower bound by extending our chain argument of this paper to a k -dimensional hypercube (along the lines in [FL 82]). This extension gives a lower bound of approximately $(t/nk)^k$. This is still a considerable gap, which we would like to see closed.

References

- [B 84] G. Bracha, "An Asynchronous $L(n - 1)/3J$ -Resilient Consensus Protocol," Proceedings of 3rd ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing, pp. 154-162, August 1984.
- [BL 85] J. Burns and N. A. Lynch, "The Byzantine Firing Squad Problem," Submitted for publication.
- [CDDS 85] B. Coan, D. Dolev, C. Dwork and L. Stockmeyer, "The Distributed Firing Squad Problem," to appear in STOC 85.
- [DDS 83] D. Dolev, C. Dwork and L. Stockmeyer, "On the Minimal Synchronism Needed for Distributed Consensus," Proceedings of 24th Annual Symposium on Foundations of Computer Science, pp. 393-402, Nov. 1983.
- [DLPSW 83] D. Dolev, N. A. Lynch, S. Pinter, E. W. Stark, and W. E. Weihl, "Reaching Approximate Agreement in the Presence of Faults," Proceedings of 3rd Annual IEEE Symposium on Reliability in Distributed Software and Database Systems, pp. 145-154, October 1983.
- [DS 82] D. Dolev and H. R. Strong, "Polynomial Algorithms for Multiple Processor Agreement," Proceedings of the 14th ACM SIGACT Symposium on Theory of Computing, pp. 401-407, May 1982.
- [FL 82] N. A. Lynch and M. Fischer, "A Lower Bound for the Time to Assure Interactive Consistency," *Information Processing Letters* 14, 4, pp. 183-186, June 1982.
- [FLM 85] M. Fischer, N. A. Lynch, and M. Merritt, "Shifting Scenarios: Easy Impossibility Proofs for Distributed Consensus Problems," submitted for publication.
- [FLP 83] M. J. Fischer, N. A. Lynch, and M. S. Paterson, "Impossibility of Distributed Consensus with one faulty process," Proceedings of 2nd ACM Symposium on Principles of Database Systems, March 1983.

- [LL 84] N. A. Lynch and J. Lundelius, "A New Fault-Tolerant Algorithm for Clock Synchronization," Proceedings of 3rd ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing, pp. 75-88, August 1984.
- [LSP 82] L. Lamport, R. Shostak, and M. Pease, "The Byzantine Generals Problem," *ACM Trans. on Programming Languages and Systems*, Vol. 4, No. 2, pp. 382-401 (1982).
- [PSL 80] M. Pease, R. Shostak, and L. Lamport, "Reaching Agreement in the Presence of Faults," *JACM*, Vol. 27, No. 2, pp. 228-234 (1980).