

LABORATORY FOR
COMPUTER SCIENCE



MASSACHUSETTS
INSTITUTE OF
TECHNOLOGY

MIT/LCS/TM-270

CONSENSUS IN THE PRESENCE
OF PARTIAL SYNCHRONY

Cynthia Dwork, Nancy Lynch and Larry Stockmeyer

October 1985

545 TECHNOLOGY SQUARE, CAMBRIDGE, MASSACHUSETTS 02139

CONSENSUS IN THE PRESENCE OF PARTIAL SYNCHRONY

Cynthia Dwork
IBM Research Laboratory
San Jose, CA 95193

Nancy Lynch
MIT Laboratory for Computer Science
Cambridge, MA 02139

Larry Stockmeyer
IBM Research Laboratory
San Jose, CA 95193

October, 1985
Revision of October, 1984 version of this paper.

ABSTRACT

The concept of partial synchrony in a distributed system is introduced. Partial synchrony lies between the cases of a synchronous system and an asynchronous system. In a synchronous system, there is a known fixed upper bound Δ on the time required for a message to be sent from one processor to another and a known fixed upper bound Φ on the relative speeds of different processors. In an asynchronous system, no fixed upper bounds Δ and Φ exist. In one version of partial synchrony, fixed bounds Δ and Φ exist but they are not known *a priori*. The problem is to design protocols which work correctly in the partially synchronous system regardless of the actual values of the bounds Δ and Φ . In another version of partial synchrony, the bounds are known but they are only guaranteed to hold starting at some unknown time T , and protocols must be designed to work correctly regardless of when the time T occurs. Fault tolerant consensus protocols are given for various cases of partial synchrony and various fault models. Lower bounds are also given which show in many cases that our protocols are optimal with respect to the number of faults tolerated. Our consensus protocols for partially synchronous processors use new protocols for fault-tolerant "distributed clocks" which allow partially synchronous processors to reach some approximately common notion of time.

Keywords: Consensus, partial synchrony, Byzantine agreement, fault-tolerance and distributed algorithms

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

This work was supported in part by the Office of Naval Research under Contract N00014-82-K-0154, by the U. S. Army Research Office under Contracts DAAG29-79-C-0155 and DAAG29-84-K-0058, and by the National Science Foundation under Grant 832391-A01-DCR.

1. INTRODUCTION

1.1. Background

The role of synchronism in distributed computing has recently received considerable attention [FLP, DDS, ADG]. One method of comparing two models with differing amounts or types of synchronism is to examine a specific problem in both models. Because of its fundamental role in distributed computing, the problem chosen is often that of reaching agreement. (See [F] for a survey; see also [GLPT, Sc, G, DLPSW] for example.) One version of this problem considers a collection of N processors, p_1, \dots, p_N , which communicate by sending messages to one another. Initially each processor p_i has a value v_i drawn from some domain V of values, and the correct processors must all decide on the same value; moreover, if the initial values are all the same, say v , then v must be the common decision. In addition, the consensus protocol should operate correctly if some of the processors are faulty, e.g., crash (fail-stop faults), fail to send or receive messages when they should (omission faults), or send erroneous messages (Byzantine faults).

Fix a particular type of fault. Given assumptions about the synchronism of the message system and of the processors, one can characterize the model by its *resiliency*, the maximum number of faults which can be tolerated in any protocol in the given model. For example, it might be assumed that there is a fixed upper bound Δ on the time for messages to be delivered (*communication is synchronous*), and a fixed upper bound Φ on the rate at which one processor's clock can run faster than another's (*processors are synchronous*), and that these bounds are known *a priori* and can be "built into" the protocol. In this case, N -resilient consensus protocols exist for Byzantine failures with authentication [LSP, DS] and, therefore, also for fail-stop and omission failures; in other words, any number of faults can be tolerated. For Byzantine faults without authentication, t -resilient consensus is possible iff $N \geq 3t+1$ [LSP, L1].

Recent work has shown that the existence of both bounds Δ and Φ is necessary to achieve any resiliency, even under the weakest type of faults. Dolev, Dwork and Stockmeyer [DDS], building on earlier work of Fischer, Lynch and Paterson [FLP], prove that either if a fixed upper bound Δ on message delivery time does not exist (*communication is asynchronous*) or if a fixed upper bound Φ on relative processor speeds does not exist (*processors are asynchronous*), then there is no consensus protocol resilient to even one fail-stop fault.

In this paper, we define and study practically motivated models which lie between the completely synchronous and the completely asynchronous cases.

1.2. Partially Synchronous Communication

We first consider the case in which processors are completely synchronous (that is, $\Phi = 1$) and communication lies "between" synchronous and asynchronous. There are at least two natural ways in which communication might be partially synchronous.

One reasonable situation could be that an upper bound Δ on message delivery time exists but we do not know what it is *a priori*. On the one hand, the impossibility results of [FLP, DDS] do not apply since communication is, in fact, synchronous. On the other hand, participating processors in the known consensus protocols need to know Δ in order to know how long to wait during each round of message exchange. Of course, it is possible to pick some arbitrary Δ to use in designing the protocol, and say that whenever a message takes longer than this Δ , then either the sender or the receiver is considered to be faulty. This is not an acceptable solution to the problem since if we picked Δ too small, all the processors could soon be considered faulty, and by definition the decisions of faulty processors do not have to be consistent with the decision of any other processor. What we would like is a protocol that does not have Δ "built in". Such a protocol would operate correctly whenever it is executed in a system where some fixed upper bound Δ exists. It should also be mentioned that we do not assume any probability distribution on message transmission time which would allow Δ to be estimated by doing experiments.

Another situation could be that we know Δ , but the message system is sometimes unreliable, delivering messages late or not at all. As noted above, we do not want to consider a late or lost message as a fault. However, without any further constraint on the message system, this "unreliable" message system is at least as bad as a completely asynchronous one, and the impossibility results of [DDS] apply. Therefore, we impose an additional constraint: that for each execution, there is a *global stabilization time* (GST), unknown to the processors, such that the message system respects the upper bound Δ from time GST onward.

This constraint might at first seem too strong: in realistic situations, the upper bound cannot reasonably be expected to hold forever after GST, but perhaps only for a limited time. However, any good solution to the consensus problem in this model would have an upper bound L on the amount of time after GST required for consensus to be reached; in this case, it is not really necessary that the bound Δ hold forever after time GST, but only up to time GST + L . We find it technically convenient to avoid explicit mention of the interval length L in the model, but will instead present the appropriate upper bounds on time for each of our algorithms.

Instead of requiring that the consensus problem be solvable in the GST model, we might think of separating the correctness conditions into *safety* and *termination* properties. The safety conditions are that no two correct processors should ever reach disagreement, and that no correct processor should ever make a decision which is contrary to the specified validity conditions. The termination property is just that each correct processor should eventually make a decision. Then we might require an algorithm to satisfy the safety conditions no matter how asynchronously the message system behaves, i.e., even if Δ does not hold eventually. On the other hand, we might only require termination in case Δ holds eventually. It is easy to see that these safety and termination conditions are equivalent to our GST condition: if an algorithm solves the consensus problem when Δ holds from time GST onward, then that algorithm cannot possibly violate a safety property even if the message system is completely asynchronous. This is because safety violations must occur at some finite point in time, and there would be some continuation of the violating execution in which Δ eventually holds.

Thus, the condition that Δ holds from some time GST onward provides a second reasonable definition for partial communication synchrony. Once again, it is not clear how we could apply previously known consensus protocols to this model. For example, the same argument as for the case of the unknown bound shows that we cannot treat lost or delayed messages in the same way as processor faults.

For succinctness, we say that communication is *partially synchronous* if one of these two situations holds: Δ exists but is not known, or Δ is known and has to hold from some point on.

Our results determine precisely the maximum resiliency possible in cases where communication is partially synchronous, for four interesting fault models. For fail-stop or omission faults, we show that t -resilient consensus is possible iff $N \geq 2t+1$. For Byzantine faults with authentication, we show that t -resilient consensus is possible iff $N \geq 3t+1$. Also, for Byzantine faults without authentication, we show that t -resilient consensus is possible iff $N \geq 3t+1$. (The "only if" direction follows immediately from the result for the completely synchronous case in [LSP].) For the first three types of faults, the time required for all correct processors to reach consensus is (1) a polynomial in N and Δ , for the model in which Δ is unknown, and (2) GST plus a polynomial in N and Δ , for the GST model. On the other hand, our algorithm for the unauthenticated Byzantine case uses an amount of time which is exponential in t . We also have a t -resilient consensus protocol for Byzantine faults without authentication, which uses a polynomial amount of time but which requires $N \geq 4t+1$. We do not know whether it is possible to obtain such a protocol for $3t+1 \leq N \leq 4t$. All of our protocols which reach consensus within time polynomial in parameters such as N and Δ also have the property that the total number of message bits sent (after GST) is also bounded above by a polynomial in the same parameters.

Table 1 shows the maximum resiliency in various cases and compares our results with previous work. The results where communication is partially synchronous and processors are synchronous are shown in column 3 of the table; the results in columns 4 and 5 will be explained shortly. In each case, the table gives N_{\min} , the smallest value of N ($N \geq 2$) for which there is a t -resilient protocol ($t \geq 1$). (Some of the lower bounds on N_{\min} in the "Partially Synchronous Processors" column of the table have slightly stronger constraints on t and N which are given in the formal statements of the theorems.) Results in the synchronous column are due to [LSP, DS, DFFLS] while those in the asynchronous column are due to [FLP, DDS]. A table entry which is a closed interval $[a,b]$ means that $a \leq N_{\min} \leq b$. Except where indicated, by "(exp)", the algorithms require a polynomial amount of time.

It is interesting to note that for fail-stop, omission and Byzantine faults with authentication, the maximum resiliency for partially synchronous communication lies strictly between the maximum resiliency for the synchronous and asynchronous cases. It is also interesting to note that for partially synchronous communication, authentication does not improve resiliency.

Our protocols use variations on a common method: a processor p tries to get other processors to change to some value v which p has found to be "acceptable"; p decides v if it receives suffi-

Failure Mode	Synchronous	Asynchronous	Partially Synchronous Communication	Partially Synchronous Communication & Processors	Partially Synchronous Processors
Fail-Stop	t	∞	$2t+1$	$2t+1$	t
Omission	t	∞	$2t+1$	$2t+1$	$[2t, 2t+1]$
Byzantine Authentication	t	∞	$3t+1$	$3t+1$	$2t+1$
Byzantine	$3t+1$	∞	$3t+1$ (exp) $[3t+1, 4t+1]$	$3t+1$ (exp) $[3t+1, 4t+1]$	$3t+1$ (exp) $[3t+1, 4t+1]$

Table 1. Smallest number of processors N_{\min} for which a t -resilient consensus protocol exists.

ciently many acknowledgements from others that they have changed their value to v , so that a value different from v will never be found acceptable at a later time. Similar methods have already appeared in the literature (see, for example, Skeen [Sk], Bracha and Toueg [BT]). Reischuk [R] and Pinter [P] have also obtained consensus results which treat message and processor faults separately.

1.3. Partially Synchronous Communication and Processors

It is easy to extend the models described in §1.2 to allow processors, as well as communication, to be partially synchronous. That is, Φ (the upper bound on relative processor speed) can exist but be unknown, or Φ can be known but actually hold only from some time GST onward. We obtain results which completely characterize the resiliency in cases where both communication and processors are partially synchronous, for all four of the classes of faults. In such cases, we assume that communication and processors possess the same type of partial synchrony, that is, either both Φ and Δ are unknown, or both hold from some time GST on.

Surprisingly, the bounds we obtain are exactly the same as for the case where communication alone is partially synchronous; see column 4 of Table 1. (The only difference is that in this case the polynomial bounds on time depend on N , Δ and Φ .) In the earlier case, the fact that Φ was equal to 1 implied that each processor could maintain a local time that was guaranteed to be perfectly synchronized with the local times of other processors. In this case, no such notion of time is available. We give two new protocols allowing processors to simulate *distributed clocks*. (These are fault-tolerant variations on the clock used by Lamport in [L2].) One uses $2t+1$ processors and tolerates t fail-stop, omission, or authenticated Byzantine faults, while the other uses $3t+1$ processors and tolerates t unauthenticated Byzantine faults. When the appropriate clock is combined with each of our protocols for the case where only communication is partially synchronous, the result is a new protocol for the case where both communication and processors are partially synchronous.

1.4. Partially Synchronous Processors

In analogy to our treatment of partial communication synchrony, it is easy to define models where processors are partially synchronous and communication is synchronous (Δ exists and is known *a priori*). Column 5 of Table 1 summarizes our results for this case. Once again, time is polynomial (this time in N , Δ and Φ) unless otherwise indicated. The basic strategy used in constructing the protocols for this case also involves combining a consensus protocol which assumes processor synchrony with a distributed clock protocol. For fail-stop faults and Byzantine faults with authentication, either the distributed clock or the consensus protocol can tolerate more failures than the corresponding clock or consensus protocol used for the case where both communication and processors are partially synchronous, so we obtain better resiliencies.

Technical remarks:

1. Our protocols assume that an atomic step of a processor is to either receive messages or send a message to a single processor, but not both; there is no atomic receive/send operation nor an atomic broadcast operation. We adopt this rather weak definition of a processor's atomic step in this paper because it is realistic in practice and seems consistent with assumptions made in much of the previous work on distributed agreement. However, our lower bound arguments are still valid if a processor can receive messages and broadcast a message to all processors in a single atomic step.
2. The strong unanimity condition requires that if all initial values are the same, say v , then v must be the common decision. Weak unanimity requires this condition to hold only if no processor is faulty. Unless noted otherwise, our consensus protocols achieve strong unanimity, and our lower bounds hold even for weak unanimity. However, in the case of Byzantine faults with authentication and partially synchronous processors, the upper bound $2t+1$ in column 5 of Table 1 holds for strong unanimity only if the initial values are signed by a distinguished "sender". This assumption is also used in the algorithm of [DS] for the completely synchronous case. (For weak unanimity, the upper bound $2t+1$ in column 5 holds even without signed

initial values.) We discuss this further in Section 6 which is the first place where the issue of whether the initial values are signed has any effect on our results.

3. Our consensus protocols are designed for an arbitrary value domain V , whereas our lower bounds hold even for the case $|V| = 2$.

The remainder of the paper is organized as follows. Section 2 contains definitions. Section 3 contains our basic protocols, presented in a basic round model which has more power than the models in which we are really interested. Section 4 contains our results for the model in which processors are synchronous and communication is partially synchronous. In particular, the protocols of Section 3 are adapted to this model. The distributed clocks are defined in Section 5, where we also discuss how to combine our results of Section 3 with the clocks to produce protocols for the model in which both processors and communication are partially synchronous. Section 6 contains our results for the case where processors are partially synchronous and communication is synchronous.

2. DEFINITIONS

2.1. Model of Computation

Our formal model of computation is based on the models of [FLP, DDS]. Here we review the basic features of the model informally. The communication system is modeled as a collection of N sets of messages, called *buffers*, one for each processor. The buffer of p_i represents messages which have been sent to p_i but not yet received. Each processor follows a deterministic protocol involving the receipt and sending of messages. Each processor p_i can perform one of the following instructions in each *step* of its protocol:

Send(m, p_j): places message m in p_j 's buffer;

Receive(p_i): removes some (possibly empty) set S of messages from p_i 's buffer and delivers them to p_i .

In the Send(m, p_j) instruction, p_j can be any processor, i.e., the communication network is completely connected. A processor's protocol is specified by a state transition diagram; the number of states can be infinite. The instruction to be executed next depends on the current state, and the execution causes a state transition. For a Send instruction, the next state depends only on the current state, while for a Receive instruction, the next state depends also on the set S of delivered messages. The initial state of a processor p_i is determined by its initial value v_i in V . At some point in its computation, a processor can irreversibly *decide* on a value in V .

For subsequent definitions, it is useful to imagine that there is a *real-time clock* outside the system that measures time in discrete integer-numbered steps. At each tick of real time, some

processors each take one step of their protocols. A *run* of the system is described by specifying the initial states for all processors, and by specifying, for each real-time step:

- (1) which processors take steps,
- (2) the instruction which each processor executes, and
- (3) for each Receive instruction, the set of messages delivered.

Runs can be finite or infinite. Given an infinite run R , the message m is *lost* in run R if m is sent by some $\text{Send}(m, p_j)$, p_j executes infinitely many Receive instructions in R , and m is never delivered by any $\text{Receive}(p_j)$.

2.2. Failures

A processor *executes correctly* if it always performs instructions of its protocol (transition diagram) correctly. A processor is *correct* if it executes correctly and takes infinitely many steps in any infinite run. We consider four types of increasingly destructive faulty behavior of processor p_i .

Fail-stop: The processor p_i executes correctly but can stop at any time. Once stopped it cannot restart.

Omission: Faulty processor p_i follows its protocol correctly, but $\text{Send}(m, p_j)$ might not place m in p_j 's buffer and $\text{Receive}(p_i)$ might cause only a subset of the delivered messages to actually be received by p_i . In other words, an omission fault on reception occurs when some set S of messages is delivered to p_i and all messages in S are removed from p_i 's buffer, but p_i follows a state transition as though some (possibly empty) subset S' of S were delivered.

Authenticated Byzantine: Arbitrary behavior. However, messages can be signed with the name of the sending processor in such a way that this signature cannot be forged by any other processor.

Byzantine: Arbitrary behavior, and no mechanism for signatures. However, we assume that the receiver of a message knows the identity of the sender.

2.3. Partial Synchrony

Let I be an interval of real time and let R be a run. We say that the communication bound Δ *holds in I for run R* provided that if message m is placed in p_j 's buffer by some $\text{Send}(m, p_j)$ at a time s_1 in I , and if p_j executes a $\text{Receive}(p_j)$ at a time s_2 in I with $s_2 \geq s_1 + \Delta$, then m must be delivered to p_j at time s_2 or earlier. This says intuitively that Δ is an upper bound on message transmission time in the interval I . The processor bound Φ *holds in I for R* provided that in any contiguous subinterval of I containing Φ real time steps, every correct processor must take at least one step. This implies that no correct processor can run more than Φ times slower than another in the interval I .

The following conditions, which define varying degrees of communication synchrony, place constraints on the kinds of runs that are allowed. In these definitions, Δ denotes some particular positive integer.

- (1) Δ is known: the communication bound Δ holds in $[1, \infty)$ for every run R ;
delta is known: Δ is known, for some fixed Δ .

This is the usual definition of *synchronous communication*.

- (2) *delta is unknown*: for every run R there is a Δ which holds in $[1, \infty)$.

- (3) Δ holds eventually: for every run R , there is a time T such that Δ holds in $[T, \infty)$.
delta holds eventually: Δ holds eventually, for some fixed Δ .

Such a time T is called *GST*, the *Global Stabilization Time*.

If either (2) or (3) holds, we say that *communication is partially synchronous*.

It is helpful to view each situation as a game between a protocol designer and an adversary. If *delta* is known, the adversary names an integer Δ and the protocol designer must supply a consensus protocol which is correct if Δ always holds. If *delta* is unknown, the protocol designer supplies the consensus protocol first, then the adversary names a Δ , and the protocol must be correct if that Δ always holds. If *delta* holds eventually, the adversary picks Δ , the designer supplies a consensus protocol, and the adversary picks a time T when Δ must start holding.

By replacing Δ by Φ and replacing "delta" by "phi" above, (1) defines *synchronous processors*, and (2) and (3) define two types of *partially synchronous processors*.

2.4. Correctness of a Consensus Protocol

Given assumptions A about processor and communication synchrony, given a fault mode F , and given a number N of processors and an integer t with $0 \leq t \leq N$, correctness of a *t-resilient consensus protocol* is defined as follows.

For any set C containing at least $N-t$ processors and any run R satisfying A and in which the processors in C are correct and the behavior of the processors not in C is allowed by the fault mode F , the protocol achieves:

Consistency: No two different processors in C decide differently,

Termination: If R is infinite then every processor in C makes a decision,

Unanimity: There are two types:

Strong Unanimity: if all initial values are v and if any processor in C decides, then it decides v .

Weak Unanimity: if all initial values are v , if C contains all processors and if any processor decides, then it decides v .

3. A BASIC ROUND MODEL

In this section we define the basic round model and present preliminary versions of our algorithms in this model. In the following sections, we show how each of our models can simulate the basic model.

3.1. Definition of the Model

In the basic round model, processing is divided into synchronous *rounds* of message exchange. Each round consists of a *Send subround*, a *Receive subround* and a *computation subround*. In a Send subround, each processor sends messages to any subset of the processors. In a Receive subround, some subset of the messages sent to the processor during the corresponding Send subround are delivered. In a computation subround, each processor executes a state transition based on the set of messages just received. Not all messages that are sent need arrive, some can be lost; however, we assume that there is some round GST, such that all messages sent from correct processors to correct processors at round GST or afterwards are delivered during the round at which they were sent. Although all processors have a common numbering for the rounds, they do not know when round GST occurs. The various kinds of faults are defined for the basic model as for the earlier models.

3.2. Protocols in the Basic Round Model

In the remainder of this section, we show how the consensus problem can be solved for the basic model, for each of the fault types.

To argue that our protocols achieve strong unanimity, we use the notion of a *proper* value defined as follows: if all processors start with the same value v , then v is the only proper value; if there are at least two different initial values, then all values in V are proper. In all protocols, each processor will maintain a local variable PROPER, which contains a set of values which the processor knows to be proper. Processors will always piggyback their current PROPER sets on all messages. The way of updating the PROPER sets will vary from algorithm to algorithm. If only weak unanimity is desired, the PROPER sets are not needed and the protocols can be simplified somewhat; we leave these simplifications to the interested reader.

3.2.1. Fail-Stop and Omission Faults.

The first algorithm is used for either fail-stop or omission faults. It achieves strong unanimity for an arbitrary value domain V .

Algorithm 1: $N \geq 2t + 1$

Initially, each processor's set *PROPER* contains just its own initial value. Each processor attaches its current value of *PROPER* to every message that it sends. Whenever a processor p receives a *PROPER* set from another processor that contains a particular value, v , then p puts v into its own *PROPER* set. It is easy to check that each *PROPER* set always contains only proper values.

The rounds are organized into alternating *trying* and *lock release* phases, where each trying phase consists of 3 rounds and each lock release phase consists of 1 round. Each pair of corresponding phases is assigned an integer, starting with 1. We say that phase h belongs to processor $h \bmod N$.

At various times during the algorithm, a processor may *lock* a value v . A *phase number* is associated with every lock. If p locks v with associated phase number $k = i \bmod N$, it means that p thinks that processor p_i might decide v at phase k . Processor p only releases a lock if it learns that its supposition was false. A value v is *acceptable* to p if p does not have a lock on any value other than v . Initially, no value is locked.

We now describe the processing during a particular trying phase k . Let $s = 4k - 3$ be the number of the first round in phase k , and assume $k = i \bmod N$. At round s , each processor (including p_i) sends a list of all its acceptable values which are also in its proper set to processor p_i (in the form of a (list, k) message). (If V is very large or infinite, it is more efficient to send a list of proper values and a list of unacceptable values. Given these lists, the proper acceptable values are easily deduced.) Just after round s , i.e., during the computation subround between rounds s and $s+1$, processor p_i attempts to choose a value to *propose*. In order for processor p_i to propose v , it must have heard that at least $N-t$ processors (possibly including itself) find value v acceptable and proper at the beginning of phase k . There might be more than one possible value which processor p_i might propose; in this case, processor p_i will choose one arbitrarily. Processor p_i then broadcasts a message $(\text{lock } v, k)$ at round $s+1$.

If any processor receives a $(\text{lock } v, k)$ message at round $s+1$, it locks v , associating the phase number k with the lock, and sends an acknowledgement to processor p_i (in the form of an (ack, k) message), at round $s+2$. In this case, any earlier lock on v is released. (Any locks on other values are not released at this time.)

If processor p_i receives acknowledgements from at least $t+1$ processors at round $s+2$, then processor p_i decides v . After deciding v , processor p_i continues to participate in the algorithm.

Lock-release phase k occurs at round $s+3 = 4k$. At round $s+3$, each processor p broadcasts the message (v, h) for all v and h such that p has a lock on v with associated phase h . If any processor has a lock on some value v with associated phase h , and receives a message (w, h') with $w \neq v$ and $h' \geq h$, then the processor releases its lock on v .

Lemma 3.1.1. It is impossible for two distinct values to acquire locks with the same associated phase.

Proof. In order for two values v and w to acquire a lock at trying phase k , the processor to which phase k belongs must send conflicting (lock v , k) and (lock w , k) messages, which it will never do in this fault model. \square

Lemma 3.1.2. Suppose that some processor decides v at phase k , and k is the smallest numbered phase at which a decision is made. Then at least $t+1$ processors lock v at phase k . Moreover, each of the processors that locks v at phase k will, from that time onward, always have a lock on v with associated phase number at least k .

Proof. It is clear that at least $t+1$ processors lock v at phase k . Assume that the second conclusion is false. Then let ℓ be the first phase at which one of the locks on v set at phase k is released without immediately being replaced by another, higher-numbered lock on v . In this case the lock is released during lock-release phase ℓ , when it is learned that some processor has a lock on some $w \neq v$ with associated phase h , where $k \leq h \leq \ell$. Lemma 3.1.1 implies that no processor has a lock on any $w \neq v$ with associated phase k . Therefore, some processor has a lock on w with associated phase h , where $k < h \leq \ell$. Thus, it must be that w is found acceptable to at least $N-t$ processors at the first round of some phase numbered h , $k < h \leq \ell$, which means that at least $N-t$ processors do not have v locked at the beginning of that phase. Since $t+1$ processors have v locked at least through the first round of ℓ , this is impossible. \square

Lemma 3.1.3. Immediately after any lock-release phase which occurs at or after GST, the set of values locked by correct processors contains at most one value.

Proof. Straightforward from the lock-release rule. \square

Theorem 3.1. Assume the basic model with fail-stop or omission faults. Assume $N \geq 2t+1$. Then Algorithm 1 achieves consistency, termination and strong unanimity for an arbitrary value domain.

Proof. First, we show consistency. Suppose that some correct processor p_i decides v at phase k , and this is the smallest numbered phase at which a decision is made. Then Lemma 3.1.2 implies that at all times after phase k , at least $t+1$ processors have v locked. In consequence, at no later phase can any value other than v ever be acceptable to $N-t$ processors, so no processor will ever decide any value other than v .

Next, we argue strong unanimity. If all the initial values are v , then v is the only value which is ever in the PROPER set of any processor. Thus, v is the only possible decision value.

Finally, we argue termination. Consider any trying phase k belonging to a correct processor p_i which is executed after a lock-release phase, both occurring at or after round GST. We claim that processor p_i will reach a decision at trying phase k (if it has not done so already). By Lemma

3.1.3, there is at most one value locked by correct processors at the start of trying phase k . If there is such a locked value, v , then sufficient communication has occurred by the beginning of trying phase k so that v is in the PROPER set of each correct processor. Moreover, any initial value of a correct processor is in the PROPER set of each correct processor at the beginning of trying phase k . It follows that a proper, acceptable value will be found for processor p_i to propose, and that the proposed value will be decided upon by processor p_i at trying phase k . \square

It is easy to see that all correct processors make decisions by round $GST + 4(N+1)$.

3.2.2. Byzantine Faults with Authentication

The second algorithm achieves strong unanimity for an arbitrary value set V , in the case of Byzantine faults with authentication.

Algorithm 2: $N \geq 3t+1$

Initially, each processor's PROPER set contains just its own initial value. Each processor attaches its PROPER set and its initial value to every message it sends. If a processor p ever receives $2t+1$ initial values from different processors, among which there are not $t+1$ with the same value, then p puts all of V (the total value domain) into its PROPER set. (Of course, p would actually just set a bit indicating that PROPER contains all of V .) When a processor p receives claims from at least $t+1$ other processors that a particular value v is in their PROPER sets, then p puts v into its own PROPER set. It is not difficult to check that each PROPER set for a correct processor always contains only proper values.

Processing is again divided into alternating trying and lock release phases, with phases numbered as before and of the same length as before. As before, at various times during the algorithm, processors may lock values. In Algorithm 2, not only is a phase number associated with every lock, but also a *proof of acceptability* of the locked value, in the form of a set of signed messages, sent by $N-t$ processors, saying that the locked value is acceptable and in their PROPER sets at the beginning of the given phase. A value v is *acceptable* to p if p does not have a lock on any value other than v .

We now describe the processing during a particular trying phase k . Let $s = 4k-3$ be the first round of phase k , and assume $k = i \bmod N$. At round s , each processor p_j (including p_i) sends a list of all its acceptable values which are also in its PROPER set to processor p_i , in the form $E_j(\text{list}, k)$, where E_j is an authentication function. Just after round s , processor p_i attempts to choose a value to propose. In order for processor p_i to propose v , it must have heard that at least $N-t$ processors find value v acceptable and proper at phase k . Again, if there is more than one possible value which processor p_i might propose, then it will choose one arbitrarily. Processor p_i then broadcasts a message $E_i(\text{lock } v, k, \text{proof})$, where the proof consists of the set of signed messages $E_j(\text{list}, k)$ received from the $N-t$ processors which found v acceptable and proper.

If any processor receives an $E_i(\text{lock } v, k, \text{proof})$ message at round $s+1$, it decodes the proof to check that $N-t$ processors find v acceptable and proper at phase k . If the proof is valid, it locks v , associating the phase number k and the message $E_i(\text{lock } v, k, \text{proof})$ with the lock, and sends an acknowledgement to processor p_i . In this case, any earlier lock on v is released. (Any locks on other values are not released at this time.) If the processor should receive such messages for more than one value v , it handles each one similarly. The entire message $E_i(\text{lock } v, k, \text{proof})$ is said to be a *valid lock* on v at phase k .

If processor p_i receives acknowledgements from at least $2t+1$ processors, then processor p_i decides v . After deciding v , processor p_i continues to participate in the algorithm.

Lock-release phase k occurs at round $s+3 = 4k$. Processors broadcast messages of the form $E_i(\text{lock } v, h, \text{proof})$, indicating that the sender has a lock on v with associated phase h and the given associated proof, and processor p_i sent the message at phase h which caused the lock to be placed. If any processor has a lock on some value v with associated phase h , and receives a properly signed message $E_j(\text{lock } w, h', \text{proof}')$ with $w \neq v$ and $h' \geq h$, then the processor releases its lock on v .

Lemma 3.2.1. It is impossible for two distinct values to acquire valid locks at the same trying phase, if that phase belongs to a correct processor.

Proof. In order for different values v and w to acquire valid locks at trying phase k , the processor p_i to which phase k belongs must send conflicting $E_i(\text{lock } v, k, \text{proof})$ and $E_i(\text{lock } w, k, \text{proof}')$ messages, which correct processors can never do. \square

Lemma 3.2.2. Suppose that some correct processor decides v at phase k , and k is the smallest numbered phase at which a decision is made by a correct processor. Then at least $t+1$ correct processors lock v at phase k . Moreover, each of the correct processors that locks v at phase k will, from that time onward, always have a lock on v with associated phase number at least k .

Proof. Since at least $2t+1$ processors send an acknowledgement that they locked v at phase k , it is clear that at least $t+1$ correct processors lock v at phase k . Assume that the second conclusion is false. Then let ℓ be the first phase at which one of the locks on v set at phase k is released without immediately being replaced by another, higher-numbered lock on v . Then the lock is released during lock-release phase ℓ , when it is learned that some processor has a valid lock on some $w \neq v$ with associated phase h , where $k \leq h \leq \ell$. Lemma 3.2.1 implies that no processor has a valid lock on any $w \neq v$ with associated phase k . Therefore, some processor has a valid lock on w with associated phase h , where $k < h \leq \ell$. Thus, it must be that w is found acceptable to all but at most t of the correct processors at the first round of some phase numbered h , $k < h \leq \ell$, which means that at most t correct processors have v locked at the beginning of that phase. Since $t+1$ correct processors have v locked at least through the first round of ℓ , this is impossible. \square

Lemma 3.2.3. Immediately after any lock-release phase which occurs at or after GST, the set of values locked by correct processors contains at most one value.

Proof. Straightforward from the lock-release rule. \square

Theorem 3.2. Assume the basic model with Byzantine faults and authentication. Assume $N \geq 3t+1$. Then Algorithm 2 achieves consistency, termination and strong unanimity for an arbitrary value domain.

Proof. First, we show consistency. Suppose that some correct processor p_i decides v at phase k , and this is the smallest numbered phase at which a decision is made by a correct processor. Then Lemma 3.2.2 implies that at all times after phase k , at least $t+1$ correct processors have v locked. In consequence, at no later phase can any value other than v ever be acceptable to all but t of the correct processors, so no correct processor will ever decide any value other than v .

Next, we argue strong unanimity. If v is the initial value of all the processors, then v is the only value which is ever put into any correct processor's PROPER set, and thus is the only possible decision value for a correct processor.

Finally, we argue termination. Consider any trying phase k belonging to a correct processor p_i which is executed after a lock-release phase, both occurring at or after GST. We claim that processor p_i will reach a decision at trying phase k (if it has not done so already). By Lemma 3.2.3, there is at most one value locked by correct processors at the start of trying phase k . If there is such a locked value v , then v was found to be proper to at least $N-t$ processors, of which $N-2t \geq t+1$ must be correct. Therefore, by the beginning of trying phase k , these $t+1$ correct processors have communicated to all correct processors that v is proper, so by the way the set PROPER is augmented every correct processor will have v in its PROPER set by the beginning of trying phase k . Next consider the case that no value is locked at the beginning of trying phase k (so all values are acceptable). If there are at least $t+1$ correct processors with the same initial value v , then v is in the PROPER set of each correct processor at the beginning of trying phase k . On the other hand, if this is not the case, then all values in the value set are in the PROPER set of all correct processors at the beginning of trying phase k . It follows that a proper, acceptable value will be found for processor p_i to propose, and that the proposed value will be decided upon by processor p_i at trying phase k . \square

As in the previous case, $GST + 4(N+1)$ is an upper bound on the number of rounds required for all the correct processors to reach decisions.

3.2.3. Byzantine Faults without Authentication

Here, we will describe two protocols. The first, simpler, protocol is t -resilient and uses $4t+1$ processors. It uses a number of rounds which is given by GST plus a constant times N . The second protocol needs only $3t+1$ processors, but it does not use a polynomial number of rounds. Both algorithms achieve strong unanimity for arbitrary value domains. In both algorithms, the processors' PROPER sets are handled exactly as in Algorithm 2.

Algorithm 3: $N \geq 4t+1$

Processing is again divided into alternating trying and lock release phases, with phases numbered as before. Now, however, the trying phases contain 4 rounds rather than 3. As before, at various times during the algorithm, processors may lock values. In Algorithm 3, only a phase number is associated with every lock. As before, a value v is *acceptable* to p if p does not have a lock on any value other than v .

We now describe the processing during a particular trying phase k . Let $s = 5k-4$ be the first round in phase k , and assume $k = i \bmod N$. At round s , each processor broadcasts a list of all its acceptable values which are also in its PROPER set, in the form (list, k) . At round $s+1$, each processor p broadcasts a vector which says, for each processor q , which values p received from q at the preceding round. Just after round $s+1$, processor p_i attempts to choose a value to propose. In order for processor p_i to propose v , it must have heard that each of at least $N-2t$ processors claims that at least $N-2t$ processors find value v acceptable and proper at phase k . As before, ambiguities are resolved arbitrarily. Processor p_i then broadcasts a message $(\text{lock } v, k)$.

If any processor receives a $(\text{lock } v, k)$ message at round $s+2$, and also has heard that each of at least $N-3t$ processors claims that at least $N-2t$ processors find value v acceptable and proper at phase k , it locks v , associating the phase number k with the lock, and sends an acknowledgement to processor p_i at round $s+3$. Release of earlier locks on v are handled as before. If processor p_i receives acknowledgements from at least $3t+1$ processors, then processor p_i decides v . After deciding v , processor p_i continues to participate in the algorithm.

The lock-release phase k occurs at round $s+4 = 5k$. Processors broadcast messages of the form (v, h) , indicating that the sender has a lock on v with associated phase h . If any processor has a lock on some value v with associated phase h , and receives $t+1$ messages indicating that $t+1$ distinct processors all have locks of the form (w, h') with $w \neq v$ and $h' \geq h$, then the processor releases its lock on v . (The values of w and h' need not be the same in all of these locks.)

Lemma 3.3.1. It is impossible for two distinct values to acquire locks by correct processors at the same trying phase, if that phase belongs to a correct processor.

Proof. Similar to Lemma 3.2.1. \square

Lemma 3.3.2. Suppose that some correct processor decides v at phase k , and k is the smallest numbered phase at which a decision is made by a correct processor. Then at least $2t+1$ correct processors lock v at phase k . Moreover, each of the correct processors that locks v at phase k will, from that time onward, always have a lock on v with associated phase number at least k .

Proof. It is clear that at least $2t+1$ correct processors lock v at phase k . Assume that the second conclusion is false. Then let ℓ be the first phase at which one of the locks on v set at phase k is released without immediately being replaced by another, higher-numbered lock on v . Then the lock is released during lock-release phase ℓ , when it is learned that at least $t+1$ processors have locks

on values $w \neq v$ with associated phases h , where $k \leq h \leq \ell$. Therefore, at least one correct processor, say p_j , has such a lock. Lemma 3.3.1 implies that no correct processor has a lock on any $w \neq v$ with associated phase k . Therefore, the correct processor p_j has a lock on $w \neq v$ with associated phase h , where $k < h \leq \ell$. In order for p_j to place this lock on w , at least $N-3t$ processors each claim that at least $N-2t$ processors find w acceptable at the first round of phase h . Since $N-3t \geq t+1$, at least one correct processor makes this claim, so at least $N-2t$ processors actually find w acceptable. Since $2t+1$ correct processors have v locked at least through the first round of ℓ , this is impossible. \square

Lemma 3.3.3. Immediately after any lock-release phase which occurs at or after GST, either no value is locked by a correct processor or there exists some value v locked by a correct processor such that at most t correct processors hold locks on values other than v .

Proof. Straightforward from the lock-release rule. (Consider some v whose lock is from the earliest phase from which any lock persists.) \square

Theorem 3.3. Assume the basic model with Byzantine faults without authentication. Assume $N \geq 4t+1$. Then Algorithm 3 achieves consistency, termination and strong unanimity for an arbitrary value domain.

Proof. The proof of consistency follows easily from Lemma 3.3.2 as in the previous two subsections. The proof of strong unanimity also follows as before.

Next, we argue termination. Consider any trying phase k belonging to a correct processor p_i which is executed after a lock-release phase, both occurring at or after GST. We claim that processor p_i will reach a decision at trying phase k (if it has not done so already). There are two cases. If some value v is locked by a correct processor at the beginning of trying phase k , then by Lemma 3.3.3, there is some locked value v such that at most t correct processors have values other than v locked at the start of trying phase k . Therefore, v is acceptable to at least $N-2t$ correct processors. Since v is locked by a correct processor, then as in the proof of Theorem 3.2, v is also in the PROPER set of all correct processors. In the second case, no value is locked by a correct processor, so all values are acceptable, and as in the proof of Theorem 3.2 some value is in the PROPER set of all correct processors. It follows in either case that a proper, acceptable value will be found for processor p_i to propose.

Moreover, any value v which is proposed by processor p_i must have had $N-2t$ processors tell p_i that $N-2t$ processors found v to be acceptable and proper. Then at least $N-3t$ processors must tell all other processors that $N-2t$ processors found v to be acceptable and proper, so that all the correct processors will acknowledge the proposal. Thus, the proposed value will be decided upon by processor p_i at trying phase k . \square

For this algorithm, $GST + 5(N+1)$ is an upper bound on the number of the round by which all correct processors must reach decisions.

The second protocol of this subsection uses only $N \geq 3t+1$ processors, but the number of rounds after GST to reach a decision grows roughly as N^4 in the worst case.

Algorithm 4: $N \geq 3t+1$

Instead of rotating processors in successive phases, we cycle repeatedly through all pairs (S,i) , where S is a size $N-t$ subset of the set of processors and p_i is a distinguished processor in that set. Each phase k is *owned* by the corresponding S , and the distinguished processor p_i plays the role of the coordinator. Processing is again divided into alternating trying and lock release phases, where each trying phase consists of four rounds and each lock release phase consists of three rounds.

We first describe the processing during a particular trying phase k . Assume that phase k is owned by the set S of $N-t$ processors and that p_i is the distinguished processor. During the first round each processor in S broadcasts a list of all its acceptable values which are also in its PROPER set, in the form $(list,k)$. Based on this information, processor p_i attempts to choose a value to propose. In order for processor p_i to propose v , it must have heard that all processors in S find v to be acceptable and proper. As before, ambiguities are resolved arbitrarily. During the second round, processor p_i broadcasts a message $(propose\ v,k)$. If a processor p_j in S receives a message $(propose\ v,k)$ from p_i and if p_j heard from all processors in S during the first round that v is acceptable and proper, then p_j broadcasts $(lock\ v,k)$ during the third round. If a processor in S receives $(lock\ v,k)$ messages from all in S , then it locks v and sends an acknowledgement to processor p_i . If processor p_i receives acknowledgements from all in S , then p_i decides v . After deciding, processor p_i continues to participate in the algorithm.

Now we describe processing during the lock release phase. During the first round of the lock release phase, all processors broadcast messages of the form (v,h) indicating that the sender has a lock on v at associated phase h . If a processor receives a message (v,h) , then during the next two rounds it checks if (v,h) is *valid* by determining the set S of processors that owns phase h , and asking each processor in S whether it sent a message $(lock\ v,h)$ at phase h . If at least $N-2t$ processors in S respond affirmatively by the end of the third round then (v,h) is valid; otherwise it is not valid. If a processor has a lock on v with associated phase h and it receives a valid message (w,h') with $w \neq v$ and $h' \geq h$, then it releases the lock on v .

Lemma 3.4.1. Suppose that some correct processor decides v at phase k , and k is the smallest numbered phase at which a decision is made by a correct processor. Then at least $t+1$ correct processors lock v at phase k . Moreover, each of the correct processors that locks v at phase k will, from that time onward, always have a lock on v with associated phase number at least k .

Proof. It is clear that at least $t+1$ correct processors lock v at phase k . Assume that the second conclusion is false. As before, let ℓ be the first phase at which one of the locks on v set at phase k is released without immediately being replaced by another, higher-numbered lock on v . Therefore, some correct processor received a valid message (w,h) during lock-release phase ℓ , where $w \neq v$ and $k \leq h \leq \ell$. Since (w,h) is valid, at least $N-2t \geq t+1$ processors said that they sent a message $(lock\ w,h)$ at phase h . Therefore, at least one correct processor p_j actually sent

(lock w, h). If $h = k$, then p_j would have sent both (lock w, k) and (lock v, k), which is impossible. Therefore, $k < h \leq \ell$. Since p_j sent (lock w, h), p_j heard during phase h that $N-t$ processors (namely, the set that owns phase h) found w to be acceptable at phase h . But since at least $t+1$ correct processors have v locked at least through the first round of trying phase ℓ , this is impossible. \square

Lemma 3.4.2. Immediately after any lock release phase which occurs at or after GST, the set of values locked by correct processors contains at most one value.

Proof. Say that processor p_i has a lock on v with associated phase h and processor p_j has a lock on w with associated phase h' where $v \neq w$. Say that $h' \geq h$. During the lock release phase, p_i will receive the message (w, h') from p_j . Since p_j received the message (lock w, h') from at least $N-t$ processors during trying phase h' and since at least $N-2t$ of these are correct, p_i will determine that (w, h') is valid. Therefore p_i will release the lock on v . \square

Theorem 3.4. Assume the basic model with Byzantine faults without authentication. Assume $N \geq 3t+1$. Then Algorithm 4 achieves consistency, termination and strong unanimity for an arbitrary value domain.

Proof. The arguments for consistency and strong unanimity are similar to before. We argue termination. Consider any trying phase k belonging to a set S consisting entirely of correct processors, such that this trying phase and the preceding lock release phase occur at or after GST. Assume p_i is the distinguished processor at phase k . We claim that processor p_i will reach a decision at trying phase k (if it has not done so already). By Lemma 3.4.2, it follows as in previous proofs that a proper, acceptable value will be found for processor p_i to propose. Moreover, since all processors in S are correct, it is obvious that the entire trying phase k will complete successfully, and processor p_i will make a decision at the end. \square

An upper bound on the number of rounds required is $\text{GST} + 7(P+1)$, where P is the number of subsets S , i.e., $P = \binom{N}{t}$, or approximately N^t .

Remark 1. Algorithms 1, 2 and 3 have the property that all correct processors make a decision within $O(N)$ rounds after GST. The time to reach agreement after GST can be improved to $O(t)$ rounds by some simple modifications. The bound $O(t)$ is optimal to within a constant factor since [FLa] shows that $t+1$ rounds are necessary even in case communication and processors are both synchronous and failures are fail-stop. A modification to all the algorithms is to have a processor repeatedly broadcast the message "Decide v " after it decides v . For Algorithm 1 (fail-stop and omission faults) a processor can decide v when it receives any "Decide v " message. For Algorithms 2 and 3 (Byzantine faults), a processor can decide v when it receives $t+1$ "Decide v " messages from different sources. Easy arguments show that the modified algorithms are still correct and that all correct processors make a decision within $O(t)$ rounds after GST, and these arguments are left to the reader.

4. PARTIALLY SYNCHRONOUS COMMUNICATION AND SYNCHRONOUS PROCESSORS

In this section we assume that processors are completely synchronous ($\Phi = 1$) and communication is partially synchronous. We show how to use these models to simulate the basic model of Section 3.1, and thus to solve the same consensus problem.

Since processors operate in lock-step synchrony, it is useful to imagine that each (correct) processor has a clock which is perfectly synchronized with the clocks of other correct processors. Initially, the clock is 0, and a processor increments its clock by 1 every time it takes a step. The assumption $\Phi = 1$ implies that the clocks of all correct processors are exactly the same at any real time step.

As presented in the definitions section, there are two different definitions of partially synchronous communication: (a) delta is unknown, and (b) delta holds eventually. We consider these two cases separately. Section 4.1 describes the upper bound results for the model in which delta holds eventually. Section 4.2 describes the upper bound results for delta unknown. Finally, Section 4.3 contains the lower bound results.

4.1. Upper Bounds When Delta Holds Eventually

We first consider the model in which delta holds eventually. Fix any of the four possible fault models. We show that if there is a t -resilient consensus protocol in the basic model, then there is one in the model where delta holds eventually. To see the implication, fix Δ , and assume algorithm A works for the basic model. From A we will define an algorithm A' for the model in which Δ holds eventually.

Let $R = N + \Delta$. Each processor divides its steps into groups of R each, and uses each group to simulate its own actions in a single round of algorithm A . More specifically, the processor uses the first N steps of group r to send its round r messages to the N processors, sending to one processor at a time, and uses the last Δ steps to perform Receive operations. The state transition for round r is simulated at the last step of group r . (The number R is large enough to allow all processors to exchange messages within a single group of steps, once GST has been reached.) Each processor always attaches a round identifier (number) to messages, and any message sent during a round r which arrives late during some round $r' > r$ is ignored. Thus, communication during each round is independent of communication during any other round.

For any run e' of A' , it is easy to show that there exists a corresponding run e of A with the following properties:

- (1) all processors which are correct in e' are also correct in e ;
- (2) the types of faults exhibited by the faulty processors are the same in e' as in e ;
- (3) every state transition of a correct processor in e is simulated by the corresponding correct processor in e' .

Since A is assumed to be a t -resilient consensus protocol for the basic round model, consensus is eventually reached in e , and so in e' , as needed.

By applying the transformation just described to Algorithms 1-4, we obtain algorithms 1^1 - 4^1 respectively. We immediately obtain the following result.

Theorem 4.1. Assume that processors are completely synchronous ($\Phi = 1$) and communication is partially synchronous (Δ holds eventually).

- (a) For the fail-stop or omission fault model, if $N \geq 2t+1$, then Algorithm 1^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (b) For the authenticated Byzantine fault model, if $N \geq 3t+1$, then Algorithm 2^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (c) For the unauthenticated Byzantine fault model, if $N \geq 4t+1$, then Algorithm 3^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (d) For the unauthenticated Byzantine fault model, if $N \geq 3t+1$, then Algorithm 4^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.

It is easy to see that Algorithms 1^1 and 2^1 guarantee that decisions are reached by all correct processors within time $4(N+1)(N+\Delta)$ after GST. The corresponding bound for Algorithm 3^1 is $5(N+1)(N+\Delta)$. For Algorithm 4^1 , an upper bound of $GST+7(P+1)(N+\Delta)$ holds, where P is approximately N^t . Thus, the time for Algorithms 1^1 - 3^1 is bounded above by GST plus a polynomial in N and Δ , while the time for Algorithm 4^1 is bounded above by GST plus a quantity which is exponential in t . Remark 1 at the end of Section 3 shows how these time bounds can be improved.

4.2. Upper Bounds for Delta Unknown

Now we consider the model in which delta is unknown. Fix any of the four possible fault models. We show that if there is a t -resilient protocol in the basic model, then there is one in the model where delta is unknown. Let algorithm A work for the basic model. As before, we define A' from A so that every execution of A' is a simulation of an execution of A .

Let $R_r = N+r$. Each processor in A' divides its steps into groups so that its r^{th} group contains exactly R_r steps. As before, the processor uses each group to simulate its own actions in a single round of algorithm A . Thus, the processor uses the first N steps of group r to send its round r messages to the N processors, one processor at a time, and uses the last r steps to perform Receive operations. The round r state transition is simulated at the last step of group r . Again, each processor always attaches a round identifier (number) to messages, and any message sent during a round r which arrives late during some round $r' > r$ is ignored.

Now consider any run e' of A' , and assume that the communication bound Δ holds in e' . As before, it is easy to define a corresponding run e of A . All processors which are correct in e' are also correct in e , and the types of faults exhibited by the faulty processors are the same in both

cases. Moreover, the number of steps in e' which are allotted for the simulation of any round $r \geq \Delta$ is sufficient to allow all messages which are sent during round r to get received. Thus, e is an allowable run of A (with Δ as its GST round). Since A is assumed to be a t -resilient consensus protocol for the basic model, consensus is eventually reached in e , and so in e' , as needed.

By applying this transformation to Algorithms 1-4, we obtain algorithms 1²-4² respectively, and immediately obtain the following result.

Theorem 4.2. In the model in which processors are completely synchronous ($\Phi = 1$) and communication is partially synchrononous (δ is unknown) claims (a)-(d) of Theorem 4.1 hold for algorithms 1²-4², respectively.

We now bound the time required by Algorithms 1²-4². Consider Algorithm 1², for example, and fix any execution e with corresponding message bound Δ . Then round Δ is the GST for the execution of Algorithm 1 simulated by e . It requires at most time $\Delta(N+\Delta)$ for processors to complete their simulation of the first Δ rounds of Algorithm 1 (Δ rounds, with $N+\Delta$ as the maximum time to simulate a single round). Then an additional $4(N+1)$ rounds, at most, must be simulated. These additional rounds require at most time $4(N+1)(N+\Delta+4(N+1))$, where the term $(N+\Delta+4(N+1))$ represents the maximum time to simulate one of these rounds (the last and largest one). Thus, the total time is bounded by $\Delta(N+\Delta) + 4(N+1)(N+\Delta+4(N+1))$, or $O(N^2 + \Delta^2)$. The same bound holds for Algorithm 2². The corresponding bound for Algorithm 3² is $\Delta(N+\Delta) + 5(N+1)(N+\Delta+5(N+1))$, and for Algorithm 4², an upper bound of $\Delta(N+\Delta) + 7(N+1)(N+\Delta+7(N+1))$ holds, where P is approximately N^1 . Thus, the time for Algorithms 1²-3² is bounded above by a polynomial in N and Δ , while the time for Algorithm 4² is bounded above by a quantity which is exponential in t . Again, these bounds can be improved using the ideas in Remark 1 at the end of Section 3.

Remark 2. If we strengthen the model where δ holds eventually to require that no messages are ever lost, but that messages sent before GST can arrive late, then we can modify Algorithms 1¹-4¹ to allow processors to terminate. Specifically, we use the ideas described in Remark 1 at the end of Section 3. However, in the present case, each processor need only broadcast a single "Decide v " message, at the time when it decides v . This message is not tagged with a round number, and other processors should accept a "Decide v " message at any time. For fail-stop or omission faults, a processor can stop participating in the algorithm immediately after it broadcasts its "Decide v " message. Further, it can decide v immediately after receiving a "Decide v " message. For Byzantine faults, a processor can decide v after receiving $t+1$ "Decide v " messages, but it cannot stop participating in the algorithm until after it has broadcast its "Decide v " message and has received "Decide v " messages from a total of $2t+1$ processors.

If messages can be lost before GST, it is not hard to argue that in any consensus protocol resilient to one fail-stop fault, at least one correct processor must continue sending messages forever. The argument is similar to those for Theorems 4.3 and 4.4 in the next subsection.

Remark 3. All the results of this section have assumed $\Phi = 1$. If processors are synchronous with $\Phi > 1$, and communication is partially synchronous, we would hope to obtain the same results. We show that this extension holds by proving a more general set of results: in Section 5 we show that the resiliency achieved by the protocols of this section can also be achieved if both processors and communication are partially synchronous. These stronger results imply that the same resiliency is achievable if communication is partially synchronous and processors are synchronous with $\Phi > 1$.

4.3. Lower Bounds

In this section, we give our lower bound results for partially synchronous communication and completely synchronous processors. The first lower bound shows that the resiliency of Theorems 4.1 and 4.2, part (a), cannot be improved, even for weak unanimity and a binary value domain.

Theorem 4.3. Assume the model with fail-stop or omission faults, where the processors are synchronous and communication is partially synchronous (either δ holds eventually or δ is unknown). Assume $t \geq 1$ and $2 \leq N \leq 2t$. Then there is no t -resilient consensus protocol which achieves weak unanimity for binary values.

Proof. The proof is the same for both definitions of partially synchronous communication. Assume the contrary, that there is an algorithm immune to fail-stop faults satisfying the required properties. We will derive a contradiction.

Divide the processors into two groups, P and Q, each with at least 1 and at most t processors. First consider the following Scenario A: all initial values are 0, the processors in Q are initially dead and all messages sent from processors in P to processors in P are delivered in exactly time 1. By t -resiliency, the processors in P must reach a decision; say that this occurs within time T_A . The decision must be 0. For if it were 1, we could modify the Scenario to one where the processors in Q are alive, but all messages sent from Q to P take more than time T_A to be delivered. In the modified Scenario, the processors in P still decide 1, contradicting weak unanimity.

Consider Scenario B: all initial values are 1, the processors in P are initially dead, and messages sent from Q to Q are delivered in exactly time 1. By a similar argument, the processors in Q decide 1 within T_B steps for some finite T_B .

Consider Scenario C (for Contradiction): processors in P have initial values 0, processors in Q have initial values 1, all processors are alive, messages sent from P to P or from Q to Q are delivered in exactly time 1, and messages sent from P to Q or from Q to P take more than $\max(T_A, T_B)$ steps to be delivered. The processors in group P (resp., group Q) act exactly as they do in Scenario A (resp., Scenario B). This yields a contradiction. \square

The following lower bound result again applies in the case of weak unanimity and a binary value domain. It shows that the resiliency of Theorems 4.1 and 4.2, part (b), cannot be improved, even for the case of weak unanimity and a binary value domain.

Theorem 4.4. Assume the model with Byzantine faults and authentication, where the processors are synchronous and communication is partially synchronous (either δ holds eventually or δ is unknown). Assume $t \geq 1$ and $2 \leq N \leq 3t$. Then there is no t -resilient consensus protocol which achieves weak unanimity for binary values.

Proof. Again, the proof is the same for both definitions of partially synchronous communication. Assume the contrary. We will derive a contradiction.

If $N = 2$, then the theorem follows from the previous lower bound, Theorem 4.3. Assume then that $N \geq 3$. Divide the processors into three groups, P , Q , and R , each with at least one and at most t processors. First consider the following Scenario A: all initial values are 0, the processors in R are initially dead and all messages sent from processors in $P \cup Q$ to processors in $P \cup Q$ are delivered in exactly time 1. By t -resiliency, the processors in $P \cup Q$ must reach a decision; say that this occurs within time T_A . As in the previous lower bound proof, the decision must be 0.

Consider Scenario B: all initial values are 1, the processors in P are initially dead, and messages sent from $Q \cup R$ to $Q \cup R$ are delivered in exactly time 1. By a similar argument, the processors in $Q \cup R$ decide 1 within T_B steps for some finite T_B .

Consider Scenario C: processors in P have initial values 0, processors in R have initial values 1, and processors in Q are faulty. The processors in Q behave with respect to those in P exactly as they do in Scenario A, and with respect to those in R exactly as they do in Scenario B. The messages sent from P to $P \cup Q$ and from R to $R \cup Q$ are delivered in exactly time 1, but all messages from P to R or from R to P take more than $\max(T_A, T_B)$ steps to be delivered. The processors in group P (resp., group R) act exactly as they do in Scenario A (resp., Scenario B). This yields a contradiction. \square

The preceding lower bound is tight, for the case of unauthenticated Byzantine faults with no further restrictions (Theorems 4.1 and 4.2, part (d)). If we consider the problem with the requirement that time be bounded by a polynomial, then we do not know how to close the gap between Theorem 4.4 and Theorems 4.1 and 4.2, part (c).

5. PARTIALLY SYNCHRONOUS COMMUNICATION AND PROCESSORS

In this section, we consider the case in which both communication and processors are partially synchronous. We show the existence of protocols with the same resiliencies as in the previous section, where only communication was partially synchronous. Moreover, the algorithms for corresponding cases still require similar amounts of time (polynomial or exponential) to the earlier

case. Again, we proceed by showing how to use the models of this section to simulate the basic model of Section 3.

In the previous section, the processors had a common notion of time which allowed time to be divided into rounds. In this case, where ϕ does not always hold, or is unknown, no such common notion of time is available. Therefore, our first task is to describe protocols which give the processors some approximately common notion of time. We call such protocols *distributed clocks*.

Our distributed clocks do not use explicit knowledge of Δ or Φ . They are designed to be used in either kind of partially synchronous model, δ and ϕ holding eventually or δ and ϕ unknown. However, the properties that the clocks exhibit do depend on the particular bounds Δ and Φ which hold (eventually) during the particular run.

Each processor maintains a private (software) clock. The private clocks grow at a rate which is within some constant factor of real time, and the private clocks remain within a constant of each other. For the model with δ and ϕ unknown, these conditions hold at all times. For the GST model, however, these conditions are only guaranteed to hold after some constant amount of time after GST. The three "constants" here depend polynomially on N , Φ and Δ . We have made no effort to optimize these constants, as this would obfuscate an already difficult and technical argument. In addition, the number of message bits sent by correct processors is polynomially bounded in N , Δ , Φ , and GST.

Once we have defined the distributed clocks, the protocols of Section 3 are simulated by letting each processor use its private clock to determine which round it is in. Several "ticks" of each private clock are used for the simulation of each round in the basic model. In order to use a distributed clock in such simulations, we need to interleave the steps of the distributed clock algorithm with steps belonging to the underlying algorithm being simulated. Moreover, the distributed clock algorithm itself is conveniently described as interleaving Receive steps, which increase the recipient's knowledge of other processors' local clocks, with Send steps, which allow the sender to inform others about its local clock. To be specific, we assume that processors alternately execute a Receive operation for the clock, a Send operation for the clock, and a step of the algorithm being simulated.

In this section, we describe what happens during the clock maintenance steps for two different distributed clocks. The first, presented in Section 5.1, handles Byzantine faults without authentication and requires $N \geq 3t+1$. The second, presented in Section 5.2, handles Byzantine faults with authentication and requires $N \geq 2t+1$. This clock obviously handles fail-stop and omission faults as well. In Section 5.3, the upper bounds for the model in which δ and ϕ hold eventually are given. In Section 5.4, we present the upper bound results for the model in which δ and ϕ are unknown. We do not prove lower bounds in this section, since the lower bounds obtained in Section 4 apply to the current models.

5.1. A Distributed Clock for Byzantine Faults without Authentication

Throughout this section we assume that $N \geq 3t+1$. We again assume that real times are numbered $0, 1, 2 \dots$. Processors participate in our distributed clock protocols by sending ticks to one another. As an expositional convenience, we define a master clock whose value at any time s depends on the past global behavior of the system and is a function of the ticks that have been sent before s . Even approximating the value of the master clock requires global information about what ticks have been sent to which processors. We therefore introduce a second type of message, called a *claim*, in which processors make assertions about the ticks they have sent.

An *i-tick* is the message "i". An i^+ -tick is a j -tick for any $j \geq i$. We say p has *broadcast an i-tick* if it has sent an i^+ -tick to all N processors.

An *i-claim* is the message "I have broadcast an i-tick". An i^+ -claim is a j -claim for any $j \geq i$. We say p has *broadcast an i-claim* if it has sent an i^+ -claim to all N processors.

We adopt the convention that all processors have exchanged ticks and claims of size 0 before time 0. These messages are not actually sent, but they are considered to have been sent and received. When we say that a certain event, such as the receipt of a certain message, has occurred "by time s " we mean that the event has occurred at some real time step $\leq s$.

The master clock, $C: N \rightarrow N$, is defined at any real time s by

$$C(s) = \text{maximum } j \text{ such that } t+1 \text{ correct processors have broadcast a } j\text{-tick by time } s.$$

Since all processors are assumed by convention to have broadcast a 0-tick before time 0, $C(0) = 0$. Note that $C(s)$ is a nondecreasing function of s .

For each processor p_i the private clock, $c_i: N \rightarrow N$, is defined by

$$c_i(s) = \text{maximum } j \text{ such that by time } s, p_i \text{ has received either}$$

- (1) messages from $2t+1$ processors where each message is a j^+ -claim, or
- (2) messages from $t+1$ processors, where each message is either a $(j+1)^+$ -tick or a $(j+1)^+$ -claim.

Since p_i is assumed to have received 0-claims from all N processors before time 0, $c_i(0) = 0$ for all correct p_i . Note that $c_i(s)$ is nondecreasing for all correct p_i .

Let p_i be a correct processor. In sending ticks p_i 's goal is to increment the master clock, so ideally we would like p_i to send a $(C(s)+1)$ -tick at time s . However, knowing $C(s)$ requires global information. Instead, p_i uses c_i , its view of C , to compute its next tick, sending a $(c_i(s)+1)$ -tick at time s . We will show in Lemma 5.1 that $c_i(s) \leq C(s)$, so p_i will never force the master clock to skip a value. We will also show that "soon" after GST for the GST model, the value of the master clock exceeds those of the private clocks by only a constant amount, so that p_i will not be pushing the master clock far ahead of the private clocks of the other processors.

Each processor p_i repeatedly cycles through all N processors, broadcasting, in different cycles, ticks and claims. The private clock of p_i is stored in a local variable c_i . Processor p_i updates its private clock every time it executes a Receive operation in the clock protocol by considering all the ticks and claims it has received and updating its private clock according to the definition of the private clock given above (thus, the private clock is updated every second clock step, i.e. every third step, that p_i takes). The following two programs describe the tick and claim broadcasting procedures. A processor begins the distributed clock protocol by setting c_i to 0 and calling $TICK(0)$, where $TICK(b)$ is the protocol shown in Figure 1. Note that the value of c_i may change during an execution of $TICK(b)$, but only a $(b+1)$ -claim (rather than a (c_i+1) -claim) is sent during execution of $CLAIM(b)$. This is consistent with our definition of what it means to have broadcast a $(b+1)$ -tick.

```

TICK(b):
  for j = 1,...,N do
    send  $(c_i+1)$ -tick to  $p_j$ ;
  CLAIM(b).

CLAIM(b):
  for j = 1,...,N do
    send  $(b+1)$ -claim to  $p_j$ ;
  if  $c_i > b$  then  $TICK(c_i)$  else  $CLAIM(b)$ .

```

Figure 1.

The following lemmas describe limitations on the rates of the master clock and the local clocks. The first three lemmas do not involve Δ and Φ , and so apply to either partially synchronous model (delta and phi holding eventually or delta and phi unknown).

Lemma 5.1. For all $s \geq 0$ and for all i such that p_i is correct, $c_i(s) \leq C(s)$.

Proof. Let $j = c_i(s)$. By the definition of the private clock, there are two possibilities.

(1) p_i has received j^+ -claims from $2t+1$ different processors. Since at least $t+1$ of these j^+ -claims are from correct processors, $C(s) \geq j$ by definition of the master clock.

(2) p_i has received messages from $t+1$ different processors, each of which is either a $(j+1)^+$ -tick or a $(j+1)^+$ -claim. Consider the earliest real time, s_0 , when some correct processor, say p_k , sends a $(j+1)^+$ -tick. By definition of the protocol, $c_k(s_0) \geq j$, so case (1) applies to p_k , and therefore

$$C(s) \geq C(s_0) \geq c_k(s_0) \geq j. \quad \square$$

Lemma 5.2. For all $s \geq 0$ the largest tick sent by a correct processor at real time s has size at most $C(s)+1$.

Proof. Immediate from the protocol and Lemma 5.1. \square

Lemma 5.3. For all $s, x \geq 0$, $C(s+x) \leq C(s) + x$.

Proof. The proof is by induction on x . For the basis, let $x = 1$. By Lemma 5.2 the largest tick sent by a correct processor by time s has size at most $C(s)+1$, so the maximum tick that can be broadcast by $t+1$ processors by time $s+1$ is a $(C(s)+1)$ -tick. Thus, $C(s+1) \leq C(s)+1$. Assume the lemma holds for some x . Then

$$\begin{aligned} C(s + (x + 1)) &= C((s + 1) + x) \leq C(s + 1) + x \quad \text{by the induction hypothesis} \\ &\leq C(s) + (x + 1) \quad \text{by the basis. } \square \end{aligned}$$

The preceding lemmas are independent of both communication and processor synchrony. Now we give several lemmas that assume such synchrony. We would like to state the lemmas in a way which applies to both kinds of partially synchronous models (delta and phi holding eventually and delta and phi unknown). So fix Δ and Φ (for either case). Also fix GST for the model in which Δ and Φ hold eventually. For the model in which delta. and phi are unknown, define GST = 0, for uniformity.

The next few lemmas discuss the behavior of the clocks after GST. Lemma 5.4 says that the private clocks increase at most a constant factor more slowly than real time. Lemmas 5.5 and 5.6 are technical lemmas used to prove the lemmas following them. Lemma 5.7 has two parts. The first says that, at any particular real time, the master clock exceeds the value of the private clocks by at most an additive constant. The second part of Lemma 5.7 says that, at least after GST, the master clock runs at a rate at most a constant factor slower than real time.

Let $D = \Delta + 3\Phi$. Note that if a message is sent to a correct processor p at time $s \geq \text{GST}$, then p will receive the message by time $s+D$: the message will be delivered by time $s+\Delta$, and within an additional time 3Φ , p will execute a Receive operation in the clock protocol.

Lemma 5.4. Assume $s \geq \text{GST}$ and let $s' = s + 12N\Phi + D$. Let j be such that $c_i(s) \geq j$ for all correct p_i . Then $c_i(s') \geq j+1$ for all correct p_i .

Proof. At time s , p_i could be executing TICK(b) for some $b < j$. However, within time $6N\Phi$ after s , p_i will call TICK(b') or CLAIM(b') for some $b' \geq j$, and within an additional $6N\Phi$ steps, p_i will broadcast a $(j+1)$ -claim. Therefore, every correct processor will broadcast a $(j+1)$ -claim by time $s + 12N\Phi$. By time s' , each correct p_i will receive at least $2t+1$ $(j+1)^+$ -claims, so $c_i(s') \geq j+1$. \square

Lemma 5.5. Assume $s \geq \text{GST}$, and let $s' = s + 39N\Phi + 4D$. Then $C(s') \geq C(s)+2$.

Proof. Let $j = C(s)$. By definition of the master clock, $t+1$ correct processors have broadcast a j -tick by time s . These $t+1$ processors send a tick or claim of size at least j to every processor within the first $3N\Phi$ steps after time s . Since these messages are sent after GST, they are received

within D steps, so $c_i(s+3N\Phi+D) \geq j-1$ for all correct p_i . By 3 applications of Lemma 5.4, $c_i(s') \geq j+2$. So $C(s') \geq j+2$ by Lemma 5.1. \square

Lemma 5.6. Let s_0 be the minimum time such that $C(s_0) \geq C(\text{GST})+2$. (Time s_0 exists by Lemma 5.5.) Let $s \geq s_0+D$. Then $c_i(s) \geq C(s-D) - 1$ for all correct p_i .

Proof. Let $j = C(s-D)$. Then $t+1$ correct processors broadcast a j -tick by $s-D$. By Lemma 5.2, the largest tick sent by a correct processor by GST is a $(C(\text{GST})+1)$ -tick. Since $j \geq C(\text{GST})+2$, the j -ticks from correct processors are broadcast entirely after GST, so they are received by time s . Thus, for all correct p_i , $c_i(s) \geq j-1$. \square

Lemma 5.7. Let s_0 be the minimum time such that $C(s_0) \geq C(\text{GST})+2$.

(1) For all $s \geq s_0+D$ and for all correct processors p_i , $c_i(s) \geq C(s) - D - 1$.

(2) For all $s \geq s_0$ and for $s' = s + 24N\Phi + 3D$, $C(s') \geq C(s)+1$.

Proof. (1) Lemma 5.6 implies that $c_i(s) \geq C(s-D) - 1$. By Lemma 5.3, $C(s) \leq C(s-D) + D \leq c_i(s) + 1 + D$. Thus, $c_i(s) \geq C(s) - D - 1$.

(2) Let $x = s+D$. Lemma 5.6 implies that $c_i(x) \geq C(s) - 1$ for all correct p_i . By two applications of Lemma 5.4, $c_i(s') \geq C(s)+1$. So $C(s') \geq C(s)+1$ by Lemma 5.1. \square

5.2. A Distributed Clock for Byzantine Faults with Authentication

The new clock is very similar to the one just described. We only explain the differences. Here we assume $N \geq 2t+1$.

An i -claim is a signed message "I have broadcast an i -tick". An i^+ -claim is a j -claim for any $j \geq i$. For $i \geq 1$, an i -tick is the message " $\langle i, i\text{-proof} \rangle$ " where a 1 -proof is the empty string and where an i -proof ($i > 1$) is a list of $t+1$ $(i-1)^+$ -claims each signed by a different processor. An i^+ -tick is a j -tick for any $j \geq i$. The definitions of *broadcast an i -tick* and *broadcast an i -claim* are the same as before.

The master clock $C: \mathbb{N} \rightarrow \mathbb{N}$ is defined by

$C(s) =$ maximum j such that some correct processor has broadcast a j -tick by time s .

The private clock $c_i: \mathbb{N} \rightarrow \mathbb{N}$ is defined by

$c_i(s) =$ maximum j such that p_i has received $t+1$ j^+ -claims (from different sources), either directly, or indirectly as part of a tick, by time s .

The definition of the clock protocol is the same as before with the addition that whenever a processor sends a $(b+1)$ -claim in the procedure CLAIM(b), it attaches the largest size tick which it can construct (this will always be a $(b+1)^+$ -tick). A correct processor will ignore any received

j -claim if it does not come with an attached j^+ -tick. The reason for this modification is so that correct processors will not accept claims that are much too large from faulty processors and incorporate these large claims into proofs.

Lemma 5.8. Lemmas 5.1-5.7 hold for the authenticated Byzantine clock.

Proof. The proofs are virtually identical to the proofs for the unauthenticated Byzantine clock, and most details are omitted. The major differences are the following.

The proof of Lemma 5.1 is easier since there is only one case. Letting $j = c_i(s)$, processor p_i has received $t+1$ j^+ -claims from different processors, at least one of which must be correct. Since a correct processor sends a j^+ -claim only after it has broadcast a j -tick, we have $C(s) \geq j$ by definition of the master clock.

The proofs of Lemmas 5.2 and 5.3 are unchanged.

In the proof of Lemma 5.4, change " $2t+1$ " to " $t+1$ ".

In the proof of Lemma 5.5, letting $j = C(s)$, we can only say that at least one correct processor has broadcast a j -tick by time s . However, this j -tick contains a j -proof consisting of $t+1$ $(j-1)^+$ -claims, so we can conclude that $c_i(s+3N\Phi+D) \geq j-1$ for all correct p_i as before. The proof of Lemma 5.6 is changed similarly.

The proof of Lemma 5.7 follows from previous lemmas by calculations, and this proof is unchanged. \square

We need one more lemma to support our claim that the number of message bits sent by correct processors is bounded above by a polynomial in GST , N , Δ and Φ .

Lemma 5.9. For all $s \geq 0$, the largest tick sent by any processor (correct or faulty) at real time s has size at most $C(s)+2$.

Proof. A j -tick sent at time s contains $t+1$ $(j-1)^+$ -claims, at least one of which was sent by a correct processor. The conclusion now follows from Lemma 5.2. \square

From this lemma and the definition of the protocol it follows easily that any tick or claim sent by a correct processor at time s can be encoded in $O(t \log C(s))$ bits.

Remark 4. The clocks of sections 5.1 and 5.2 are similar to the one discovered independently by Attiya, Dolev and Gil [ADG].

5.3. Upper Bounds When Delta and Phi Hold Eventually

We now present our upper bound results for partially synchronous communication and processors, for the model where delta and phi hold eventually. Fix any of the four possible fault models. We show that if there is a t -resilient protocol in the basic model, then there is one in the model where delta and phi hold eventually. To see the implication, fix Δ and Φ , and assume algorithm A works for the basic model. We define A' from A as follows, so that A' works for the model where Δ and Φ hold after GST.

As described above, two out of every three steps of each processor are used to maintain a distributed clock, and the other step is used to simulate Algorithm A . For fail-stop or omission faults, we use the authenticated Byzantine clock, simplified appropriately because the signatures are not needed and because we cannot assume the authentication capability. Note that the consensus protocol and the distributed clock protocol have the same constraint on the number of processors, $N \geq 2t+1$. For unauthenticated Byzantine faults, we use the unauthenticated Byzantine clock. For authenticated Byzantine faults, either clock could be used.

The Receive steps of Algorithm A' are designated as belonging to either the clock simulation or the algorithm simulation. However, each time a Receive step of A' occurs, it is possible that messages for either or both simulations will be received. We assume that each processor maintains a pair of message buffers, one for each of the two simulations it is carrying out. When the processor does a Receive step that belongs to the clock simulation, it saves any messages for the algorithm simulation in the algorithm message buffer, and vice versa. Also, each time the processor does a Receive step that belongs to the clock simulation, it collects not only the new incoming messages, but all those in the clock message buffer, to use in its clock simulation step; analogous assumptions are made for the algorithm simulation.

Fix $R = 3N\Phi + 2D + 2$, where as before, $D = \Delta + 3\Phi$. Each processor uses its private clock to determine the round of Algorithm A currently being simulated. Namely, if $(r-1)R \leq c_i(s) < rR$, then processor p_i determines at real time s that the current round is r . Processors label messages with round numbers. As long as a processor determines that the current round is r , it uses its protocol simulation steps to simulate steps of round r in the basic model. The first N protocol simulation steps are used for sending the round r messages to all the processors, and the remaining steps are spent executing Receive operations. Unlike in the simulations in Section 4, it is possible that there will be insufficient time for a processor to actually send all of its round r messages.

Processor p_i simulates its state transition for round r at its first algorithm simulation step at which it decides the current round is strictly greater than r . More specifically, assume that processor p_i has reached an algorithm simulation step, s , at which the current round is k , and assume that the round at processor p_i 's last algorithm simulation step was $h < k$. Then processor p_i simulates its state transitions for rounds $h, h+1, \dots, k-1$, all at the beginning of step s . In simulating these state transitions, processor p_i simulates all of its sending steps for these rounds, that is, it

makes the appropriate state transitions but does not actually send any messages, and it simulates the receipt of all the messages which are in the algorithm message buffer.

For any run e' of A' , it is easy to define a corresponding run e of A . We see that all processors which are correct in e' are also correct in e , and that the types of faults exhibited by the faulty processors are the same in both cases. We must argue that, within a short time after GST, the number of ticks in e' which are allotted for the simulation of any round r is sufficient to allow all round r messages to be sent and received. More precisely, the "short time after GST" is chosen so that parts (1) and (2) of Lemma 5.7 hold.

We must first show that there is sufficient time for each correct processor p_i to send all its round r messages and then to do at least one Receive operation. Assume that s is the first real time at which processor p_i 's private clock reaches or exceeds $(r-1)R$. Then processor p_i would finish sending all its round r messages and doing one Receive operation by real time $s + 3(N+1)\Phi$. We must show that processor p_i 's clock up to real time $s + 3(N+1)\Phi$ remains less than rR , i.e. that

$$(5.3.1) \quad c_i(s + 3(N+1)\Phi) < rR.$$

We must also show that there is sufficient time for all the round r messages sent by processor p_i to be received. Fix a correct processor p_j . We show that processor p_j has sufficient time to receive a round r message from processor p_i , before going on to simulate round $r+1$. Again letting s be the first real time for which $c_i(s) \geq (r-1)R$, p_i will send the message to p_j by real time $s + 3N\Phi$, and p_j will receive the message by real time $s + 3N\Phi + D$. Therefore, we must show that

$$(5.3.2) \quad c_j(s + 3N\Phi + D) < rR.$$

Since $D \geq 3\Phi$ and since clocks are nondecreasing, we can prove both (5.3.1) and (5.3.2) by showing that, for any correct processor p_k ,

$$c_k(s + 3N\Phi + D) < rR.$$

This follows because

$$\begin{aligned} c_k(s + 3N\Phi + D) &\leq C(s + 3N\Phi + D) \quad \text{by Lemma 5.1} \\ &\leq C(s - 1) + 3N\Phi + D + 1 \quad \text{by Lemma 5.3} \\ &\leq c_i(s - 1) + 3N\Phi + 2D + 2 \quad \text{by Lemma 5.7(1)} \\ &< (r - 1)R + 3N\Phi + 2D + 2 \quad \text{by assumption} \\ &= rR. \end{aligned}$$

Since A is assumed to be a t -resilient consensus protocol for the basic model, consensus is eventually reached in e , and so in e' , as needed.

By applying the transformation just described to Algorithms 1-4, we obtain algorithms 1³-4³ respectively. We immediately obtain the following result.

Theorem 5.1. Assume that communication and processors are partially synchronous (Δ and Φ hold eventually).

- (a) For the fail-stop or omission fault model, if $N \geq 2t+1$, then Algorithm 1³ achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (b) For the authenticated Byzantine fault model, if $N \geq 3t+1$, then Algorithm 2³ achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (c) For the unauthenticated Byzantine fault model, if $N \geq 4t+1$, then Algorithm 3³ achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (d) For the unauthenticated Byzantine fault model, if $N \geq 3t+1$, then Algorithm 4³ achieves consistency, termination and strong unanimity for an arbitrary value domain.

As before, we claim that Algorithms 1³-3³ reach agreement within a polynomial (in N , Δ and Φ) amount of time after GST, while the time for Algorithm 4³ is bounded above by GST plus a quantity which is exponential in t . Our claims of polynomial time performance follow from the fact that the master clock, a short time after GST, runs at a rate no slower than $1/(24N\Phi + 3(\Delta+3\Phi))$ times real time (see Lemma 5.7(2)). Finally, the total number of message bits sent by correct processors is polynomially bounded in N , Δ , Φ , and GST, since the number of bits in each message sent by a correct processor is polynomially bounded in these quantities.

5.4. Upper Bounds for Delta and Phi Unknown

Next, we present our upper bound results for partially synchronous communication and processors, for the model where Δ and Φ are unknown. The ideas are a simple combination of ideas from Sections 4.2 and 5.3. Again, fix any of the fault models. We show that if there is a t -resilient protocol A in the basic model, then there is a t -resilient protocol A' in the model where Δ and Φ are unknown. In A' , processor steps are interleaved exactly as in the previous construction, and the same clocks are used for the same fault models as before. This time, define bound $R_r = 3Nr + 8r + 2$ to describe the number of ticks to be used for the simulation of round r . (This bound is obtained from the previous bound by replacing both Δ and Φ by r .) The simulation of the message exchange is carried out just as before. For any run e' of A' , we again define a corresponding run e of A . We must argue that, for sufficiently large rounds r , the number of ticks in e' which are allotted for the simulation of round r is sufficient to allow all round r messages to be sent and received. The argument is as in Section 5.3, once r has exceeded Δ and Φ . Thus, e is an allowable run of A , and since consensus is eventually reached in e , it is also reached in e' .

By applying this transformation to Algorithms 1-4, we obtain algorithms 1⁴-4⁴ respectively.

Theorem 5.2. Assume that communication and processors are partially synchronous (Δ and Φ are unknown). Then claims (a)-(d) of Theorem 5.1 hold for Algorithms 1⁴-4⁴, respectively.

As before, we claim that Algorithms 1⁴-3⁴ reach agreement within a polynomial (in N , Δ and Φ) amount of time, while the time for Algorithm 4⁴ is bounded above by an exponential in t .

Remark 5. In the simulation of the basic model described in Sections 5.3 and 5.4, we have said that if the round number of processor p_i 's last algorithm simulation step was h and processor p_i updates its clock and finds that it is now simulating some round $k > h$, then all the state transitions in rounds h through $k-1$ are simulated (except that no messages are sent). For a general simulation of the basic model, these transitions must all be simulated since they may involve state transitions which processor p_i must make in order that the simulation of the algorithm in the basic model be correct. However, it is not hard to see that for the particular Algorithms 1-4 designed for the basic model in Section 3.2, processor p_i can just simulate the state transition for round h and continue the simulation at round k , without simulating the "missed" transitions in rounds $h+1$ through $k-1$. This can be done since the state information in Algorithms 1-4 (not including the current round number) consists of the PROPER sets, which values are locked, and other information associated with each lock. Changes in this state information are caused only by the receipt of certain messages. Since we have shown consistency for Algorithms 1-4 even if messages are lost before GST, it follows that the algorithms remain consistent if processors, including correct ones, skip state transitions before GST.

6. PARTIALLY SYNCHRONOUS PROCESSORS AND SYNCHRONOUS COMMUNICATION

In this section we consider models where processors are partially synchronous and communication is synchronous, that is, there is a fixed upper bound Δ on message transmission time which always holds (in particular, no messages are lost). Of course, the protocols of the previous section with their associated resiliencies work for such models, but by using the fact that communication is now synchronous, we can achieve higher resiliencies in some cases.

It is convenient to base our consensus algorithms on another basic model, which we call the *basic model with signals*. In section 6.1, we define this new basic model and give consensus algorithms which are designed to work in the basic model with signals. We then show how to use the eventual ϕ and unknown ϕ models to simulate the basic model with signals. As in Section 5, we use distributed clocks to give the processors some approximately common notion of time. The clocks are discussed in Section 6.2. Section 6.3 contains algorithms for the case where ϕ holds eventually, and Section 6.4 contains algorithms for ϕ unknown. Section 6.5 contains lower bounds.

6.1. A Basic Model With Signals

The basic model with signals is just like the basic model, except that the Receive subround also includes the possible receipt of a *signal* by each processor. In any round r , the receipt of a signal by processor p_i implies that all processors receive the round r messages sent to them by processor p_i . The non-receipt of a signal does not imply anything. At round GST and afterwards, we as-

sume that all correct processors receive signals at each round. The next two subsections, 6.1.1 and 6.1.2, give consensus protocols for the basic model with signals which are resilient to two types of faults.

6.1.1. Fail-Stop Faults

The next algorithm achieves strong unanimity for an arbitrary value domain V .

Algorithm 5: $N \geq t$

Each processor has a local variable $VALUE$, initialized at its initial value. We say that each round $k = i \bmod N$ belongs to processor p_i . Processing in an arbitrary round k is as follows:

Processing for p_i , where round k belongs to p_i :

Broadcast $VALUE$;

If a signal is received, then decide on $VALUE$.

Processing for p_j , where round k does not belong to p_j :

If a message is received with contents v , then set $VALUE := v$.

Lemma 6.1. Assume that processor p_i decides v at round k , and this is the smallest numbered round at which a decision is made. Then no message containing value $w \neq v$ is ever sent at any round $\geq k$.

Proof. Assume for the sake of contradiction that the lemma is false, and let h be the smallest numbered round $\geq k$ when a message containing value $w \neq v$ is sent. It is clear that $h \neq k$ since faults are fail-stop. Let p_j be the processor that owns round h .

Since processor p_i receives a signal at round k , it must be the case that processor p_j receives value v from processor p_i at round k , and therefore sets its $VALUE$ to v . By assumption, no message with value different from v is sent at rounds after k and before h . Therefore, processor p_j 's $VALUE$ remains equal to v until the beginning of round h . This contradicts the assumption that processor p_j sends w at round h . \square

Theorem 6.1. Assume the basic model with signals, with fail-stop faults. Assume $N \geq t$. Then Algorithm 5 achieves consistency, termination and strong unanimity for an arbitrary value domain.

Proof. First, we show consistency. Suppose that some correct processor p_i decides v at round k , and this is the smallest numbered round at which a decision is made. Then Lemma 6.1 implies that no message containing value $w \neq v$ is ever sent at any round $\geq k$. But a processor can decide on a value w only if it first sends out messages containing w . Therefore, no processor ever decides on a value $w \neq v$.

Strong unanimity is obvious, since a message with contents v is only sent if v was the initial value of some processor.

Since a signal is received by each correct processor at every round on or after GST, by definition of the basic model with signals, it is obvious that each round on or after GST results in a decision for its owner if that owner has not already decided. \square

6.1.2. Authenticated Byzantine Faults

The next algorithm, Algorithm 6, achieves weak unanimity for an arbitrary value domain.

Algorithm 6: $N \geq 2t+1$

The protocol is similar to Algorithm 2 of Section 3.2.2, with a few changes as indicated below. Because we are only dealing with weak unanimity, the PROPER sets are not used. This time, the rounds are divided into trying phases of 2 rounds each and lock release phases of one round each. A trying phase of Algorithm 6 is the same as the first two rounds of the corresponding trying phase of Algorithm 2 except that if a processor, during one of its trying phases, is choosing a value to propose and if several values are acceptable, the processor chooses its own initial value if that value is acceptable, or chooses arbitrarily otherwise. The third round is omitted; processor p_i does not wait for messages from others claiming that they have responded to a message $E_i(\text{lock } v, k, \text{proof})$ by locking v . Instead, it checks that a signal has been received at the second round of the trying phase. If a signal is received, then processor p_i decides v .

In Algorithm 2, processor p_i needed at least $2t+1$ acknowledgement messages to conclude that at least $t+1$ correct processors actually locked v at phase k . Now we can argue that if a signal is received, then all correct processors will have actually locked v at phase k , and since $N \geq 2t+1$ there are at least $t+1$ correct processors.

The proof of the following theorem is very similar to that of Theorem 3.2 (the result about Algorithm 2), and details are left to the reader.

Theorem 6.2. Assume the basic model with signals, with authenticated Byzantine faults. Assume $N \geq 2t+1$. Then Algorithm 6 achieves consistency, termination and weak unanimity for an arbitrary value domain.

One version of the consensus problem studied in the literature supposes that a distinguished processor, called the "general", gives the initial values v_i to all the processors. In the case of Byzantine faults with authentication, it is usually assumed that the general signs these initial values with its own unforgeable signature. Thus, if the general is correct, there is a single value v such that the general gives a signed v to every processor; in this case, strong unanimity requires that v is the value decided by all correct processors. If the general is faulty, the general can give out different values, and can even give two different values, both signed, to the same processor; in this

case, strong unanimity does not require any particular value to be the decision value. This issue was not raised earlier because it is irrelevant to the results of Sections 3, 4 and 5; that is, our protocols for the authenticated Byzantine case are designed to work even if the general does not sign the initial values, and our lower bound Theorem 4.4 is still valid if the general does sign the initial values. (If the general does sign the initial values, updating of the PROPER sets in Algorithm 2 can be considerably simplified.) This distinction is important in the completely synchronous case: N -resilient strong unanimity is possible in the authenticated Byzantine case (column 1, row 3 of Table 1) only if the general signs the initial values.

This distinction also matters in this section of the paper. Namely, consider the basic model with signals, with authenticated Byzantine faults, where the general signs the initial values, and where $N \geq 2t+1$. Then a slight variant of Algorithm 6 achieves consistency, termination and *strong* unanimity for an arbitrary value domain.

Algorithm 7: $N \geq 2t+1$

The algorithm is identical to Algorithm 6 except that PROPER sets are used. Initially, the PROPER set of processor p_i contains its initial value v_i which is signed by the general. Each processor piggybacks its initial value, signed by the general, on all messages. If p_i ever receives a value different from v_i which is also signed by the general, then p_i puts all of V in its PROPER set. It is clear that a correct processor's PROPER set always contains proper values.

6.2. Distributed Clocks

Recall that in this section there is some known communication bound Δ which always holds. Because the previous clocks have limited resiliency, we first describe a distributed clock which is resilient to any number of fail-stop faults. The general form of the clock is similar to the clocks of Sections 5.1 and 5.2.

As in Section 5.1, an *i-tick* is the message "i" and an *i-claim* is the message "I have broadcast an i-tick". The definitions of *i⁺-tick*, *i⁺-claim*, *broadcast an i-tick* and *broadcast an i-claim* are also the same as in Section 5.1. The clock protocol is given by TICK(b) and CLAIM(b) as in Figure 1.

The master clock is:

$$C(s) = \text{maximum } j \text{ such that some processor has broadcast a } j\text{-tick by time } s.$$

The private clock c_i is:

$$c_i(s) = \text{maximum } j \text{ such that } p_i \text{ has received either a } j^+\text{-claim or a } (j+1)^+\text{-tick by time } s.$$

We claim that the new fail-stop clock and the authenticated Byzantine clock of Section 5.2 when used in the model of this section have the following properties.

- (A1) For all s and all correct p_i , $c_i(s) \leq C(s)$.
- (A2) For all $s, x \geq 0$, $C(s+x) \leq C(s) + x$.
- (A3) Consider a run in which the processor bound Φ holds after time GST (in the unknown ϕ model, GST = 0 for uniformity as explained before), and let $D = 3\Phi + \Delta$. There are constants a_1 and a_2 depending polynomially on Δ and Φ such that:
- (A3.1) For all correct p_i and all $s \geq \text{GST} + a_1$, $c_i(s) \geq C(s) - D - 1$.
- (A3.2) For all $s \geq \text{GST} + a_1$, $C(s+a_2) \geq C(s) + 1$.
- (A4) For all s at which the correct processor p_i executes a Receive operation in the clock protocol,

$$C(s) - (\Delta + 1) \leq c_i(s).$$

In the fail-stop case, we consider a processor to be "correct" up until the time that it fails. This should cause no confusion since in the fail-stop model, a faulty processor acts as a correct processor up until the time when it fails.

For the authenticated Byzantine clock, we have already proved (A1), (A2) and (A3) in Lemmas 5.1, 5.3 and 5.7, respectively, with modifications as described in the proof of Lemma 5.8. To prove that these properties hold for the fail-stop clock, we first note that Lemmas 5.1-5.4 hold for the fail-stop clock; the proofs are very similar to the proofs given in Section 5.1 and are left to the reader. Lemma 5.5 is not needed. Since Δ always holds, we can prove a stronger version of Lemma 5.6 for the fail-stop clock.

Lemma 5.6'. For all $s \geq \text{GST} + D$ and all correct p_i , $c_i(s) \geq C(s - D) - 1$.

Proof. Let $j = C(s - D)$. By definition of the master clock, some processor has broadcast a j -tick by time $s - D$, so every correct processor will receive a j^+ -tick by time s . Therefore, $c_i(s) \geq j - 1$, by definition of the private clock. \square

Now Lemma 5.7 follows from Lemma 5.6' and previous Lemmas as before. (However, we only need $s \geq \text{GST} + D$ for part (1) and $s \geq \text{GST}$ for part (2)).

The proof of (A4) is similar for both clocks. Let $j = C(s - \Delta)$. A j -tick has been broadcast by time $s - \Delta$, so processor p_i , by time s , will receive a j^+ -tick. For the authenticated Byzantine clock, this j^+ -tick contains $t + 1$ $(j - 1)^+$ -claims. For either clock, by definition of the private clock and by property (A2),

$$c_i(s) \geq j - 1 = C(s - \Delta) - 1 \geq C(s) - \Delta - 1.$$

6.3. Upper Bounds When Phi Holds Eventually

The only improvements over the case in which both phi and delta hold eventually are for fail-stop faults and authenticated Byzantine faults (the latter either for weak unanimity or for strong unanimity with a general signing the initial values). Fix one of these fault models. We show that if there is a t -resilient protocol in the basic model with signals, then there is one in the model where phi holds eventually. Fix Δ and Φ and assume Algorithm A works for the basic model with signals. Define A' as follows.

Two out of every three steps of each processor are used to maintain a distributed clock, and the other step is used to simulate Algorithm A. For fail-stop faults we use the new fail-stop clock of Section 6.2, while for authenticated Byzantine faults we use the authenticated Byzantine clock. Message buffers are maintained as before.

Fix $R = 3N\Phi + (2D+2) + (\Delta+1)$ where, as before, $D = \Delta + 3\Phi$. Each processor determines the current round being simulated and conducts the rest of the simulation exactly as in Section 5.3. We must describe how signals are simulated. If a processor p_i has sent all of its messages for a particular round r , performed a Receive operation in the clock protocol and updated its private clock, and if the clock then satisfies

$$c_i < rR - (2\Delta+1),$$

then p_i assumes that it has received a signal for round r .

For any run e' of A' , we define a corresponding run e of A. Again, faults are preserved. Since the R in this section is larger than the R used in Section 5.3, it follows as in Section 5.3 that, within a short time after GST, the number of ticks in e' which are allotted for the simulation of any round r is sufficient to allow all round r messages to be sent and received. It remains to show that signals behave correctly:

- (a) whenever a correct processor p_i receives a signal at any round r , it means that all of the messages sent by processor p_i at round r to correct processors actually get received, and
- (b) within a short time after GST, all correct processors receive signals at all rounds.

We first show (a). Assume that correct processor p_i receives a signal at round r , that p_i sends a message to correct processor p_j at round r , and that s is the real time when the message is sent. Then the message arrives at processor p_j by real time $s+\Delta$. Processor p_j might not actually receive the message at this time since it is not executing a Receive operation at this time. However, the key fact for the simulation is that the message will be received the next time that p_j executes a Receive operation, and that when this Receive occurs, p_j has not yet started any round greater than r . That is, we must show that

$$c_j(s+\Delta) < rR.$$

To show this, first note that since processor p_i receives a signal for round r there must be a real time s' with $s' > s$ such that p_i executes a Receive operation in the clock protocol at time s' and

$$c_i(s') < rR - (2\Delta+1).$$

Now

$$\begin{aligned}
c_j(s + \Delta) &\leq C(s + \Delta) \text{ by (A1)} \\
&\leq C(s' + \Delta) \text{ since } s' > s \\
&\leq C(s') + \Delta \text{ by (A2)} \\
&\leq c_i(s') + 2\Delta + 1 \text{ by (A4)} \\
&< rR \text{ by the condition defining signalling.}
\end{aligned}$$

Next, we show (b). Fix some round number r after GST, and let s be the earliest time at which p_i 's private clock reaches or exceeds $(r-1)R$. Processor p_i can broadcast a message to all processors and execute a Receive operation in the clock protocol within $3(N+1)\Phi$ steps after s . Therefore, we must show that

$$c_i(s+3(N+1)\Phi) < rR - (2\Delta+1).$$

This is true because

$$\begin{aligned}
c_i(s + 3(N + 1)\Phi) &\leq C(s + 3(N + 1)\Phi) \text{ by (A1)} \\
&\leq C(s - 1) + 3(N + 1)\Phi + 1 \text{ by (A2)} \\
&\leq c_i(s - 1) + 3(N + 1)\Phi + D + 2 \text{ by (A3.1)} \\
&< (r - 1)R + 3(N + 1)\Phi + D + 2 \text{ by assumption.} \\
&\leq rR - (2\Delta + 1) \text{ by calculation.}
\end{aligned}$$

Thus, e is an allowable run of A , so consensus is reached in e , and so in e' , as needed. By applying this transformation to Algorithms 5, 6 and 7, we obtain Algorithms 5^1 , 6^1 and 7^1 , respectively.

Theorem 6.3. Assume that communication is synchronous, and processors are partially synchronous (ϕ holds eventually).

- (a) For the fail-stop model, if $N \geq t$, then Algorithm 5^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.
- (b) For Byzantine faults with authentication, if $N \geq 2t+1$, then Algorithm 6^1 achieves consistency, termination and weak unanimity for an arbitrary value domain.
- (c) For Byzantine faults with authentication, if $N \geq 2t+1$ and if the general signs the initial values, then Algorithm 7^1 achieves consistency, termination and strong unanimity for an arbitrary value domain.

6.4. Upper Bounds for Phi Unknown

The strategy is the same as in Sections 4.2 and 5.4. Namely, we use the algorithm of Section 6.3 where $R_r = 3Nr + 6r + 3\Delta + 3$ steps are allowed for the simulation of round r , where R_r is obtained from the R of Section 6.3 by replacing Φ by r . It is important to note that the verification of (a) in Section 6.3 (namely that if a signal is received by p_i at round r then all messages sent by p_i during round r to correct processors arrive before the other processor starts any round greater

than r) did not depend in any way on Φ . Therefore, (a) holds even for rounds r where r is smaller than the actual (unknown) Φ which holds in the run. Applying this transformation to Algorithms 5, 6 and 7, we obtain Algorithms 5^2 , 6^2 and 7^2 , respectively.

Theorem 6.4. Assume that communication is synchronous and processors are partially synchronous (ϕ is unknown). Then claims (a), (b) and (c) of Theorem 6.3 hold for Algorithms 5^2 , 6^2 and 7^2 , respectively.

Our claim of a polynomial time bound (after GST) for the algorithms of Sections 6.3 and 6.4 follows from clock property (A3.2) which states that the master clock runs fast enough after GST.

We should also mention that Remark 5 at the end of Section 5 does not apply to the simulations of Sections 6.3 and 6.4. Here, if a processor's clock makes a big jump so that rounds are missed, all steps of the consensus protocol during the missed round(s) must be simulated. If the correct p_i sends a message to a correct p_j and receives a signal during round r , then p_j must receive the message and make the appropriate state transition caused by this reception, even if p_j 's clock makes a large jump which causes it to miss round r .

6.5. Lower Bounds

The following lower bound shows that the resiliency of Theorems 6.3 and 6.4, parts (b) and (c), cannot be improved. The method used to prove this lower bound was suggested by Dolev [D].

Theorem 6.5. Assume the model with Byzantine faults with authentication, synchronous communication, and partially synchronous processors. Assume $t \geq 2$ and $4 \leq N \leq 2t$. Then there is no t -resilient consensus protocol which achieves weak unanimity for binary values, even if the general signs the initial values.

Proof. Assume to the contrary that a consensus algorithm exists. The proof is identical for both variations of partially synchronous processors. In the following, we assume without loss of generality that all messages are delivered in one real time step. Divide the processors into four groups P , Q , $\{b\}$ and $\{r\}$, where groups P and Q each contain at least 1 and at most $t-1$ processors and where b and r are single processors. We say that a processor *wakes up* at real time s if it takes the first step of its protocol at real time s . We say that a processor *runs fast* in the real time interval $[s_1, s_2]$ if it takes a step of its protocol at each real time step in the interval.

Consider Scenario CP where the processors in $P \cup \{b\}$ have initial values 0, wake up at time 1 and run fast in the interval $[1, \infty)$, and the other processors are initially dead. By t -resiliency, the processors in P make some decision within some finite time T_p . We claim the decision must be 0. For if it were 1, we could modify the scenario to one in which all initial values are 0, and the processors in $Q \cup \{r\}$ are correct but do not wake up until after time T_p . The processors in P still decide 1 in the modified scenario, which contradicts weak unanimity.

Consider the analogous Scenario CQ where the processors in $P \cup \{r\}$ are initially dead, and the processors in $Q \cup \{b\}$ wake up at time 1 with initial values 1 and run fast in the interval $[1, \infty)$. Therefore, the processors in Q decide 1 after some finite time T_Q .

Consider the following Scenario BP. Processors in $P \cup \{b\}$ are Byzantine. The processors in P have value 0, and b has both 0 and 1 (so the general is Byzantine). They wake up at time 1, with b acting as if its value is 0, and they run fast in the interval $[1, T_p]$. They send the same messages to r as are sent in Scenario CP, but no messages are sent to Q. After time T_p , the processors in P die. The processors in Q are correct. They wake up at time T_p+1 and run fast thereafter. Starting at time T_p+1 , the Byzantine processor b starts behaving towards Q and r exactly as it does in Scenario CQ, as if its value is 1, except that a message sent at real time s in Scenario CQ is sent at time T_p+s in Scenario BP. Since Q has received no messages from P, the processors in Q decide 1 at time T_p+T_Q , and they all behave exactly as in Scenario CQ except that everything happens T_p real time steps later. At time T_p+T_Q+1 , the correct processor r wakes up and runs fast thereafter. The initial value of r is irrelevant. Note that at most t processors are faulty in this run. In the model where phi is unknown, the processor bound $\Phi = T_p+T_Q+1$ holds in this run; in the model where phi holds eventually, the processor bound $\Phi=1$ holds after $GST = T_p+T_Q+1$. Since the correct processors in Q have already decided 1 before r wakes up, r must decide 1 at some real time T_r .

Consider now Scenario BQ. The processors in P are correct and begin with value 0. They run fast in the interval $[1, T_p]$, but take no more steps until after time T_r . In the time interval $[1, T_p]$, the Byzantine processor b behaves towards P and r exactly as it does in Scenario CP, acting as if it has initial value 0. Therefore, at time T_p the processors in P decide 0. The processors in Q are Byzantine. They wake up at time T_p+1 with value 1 and behave with respect to r exactly as they do in scenario BP, that is, the messages which have been sent from P to Q during the interval $[1, T_p]$ are ignored by Q. At time T_p+1 , b starts acting toward r exactly as it does in Scenario BP, as if it had initial value 1. The correct processor r wakes up at time T_p+T_Q+1 and runs fast thereafter. It is easy to see that the messages received by r between time T_p+T_Q+1 and time T_r are exactly the same in scenario BQ as in scenario BP. Therefore, r decides 1 at time T_r , which is a contradiction because the correct processors in P decided 0. \square

In the preceding proof, note that the processors in P and Q exhibit only omission faults: P fails to send messages to Q in Scenario BP and Q fails to receive messages from P in Scenario BQ. Processor b is the only one which exhibits Byzantine behavior stronger than omission faults. Therefore, it can be checked that the same proof can be carried out for omission faults with three groups of processors, P, Q and $\{r\}$ where P and Q each contain at least 1 and at most $t-1$ processors. This proves the following, which shows that the resiliency of Theorems 5.1 and 5.2, part (a), when applied to the case of omission faults and partially synchronous processors, cannot be improved by more than 1.

Theorem 6.6. Assume the model with omission faults, synchronous communication, and partially synchronous processors. Assume $t \geq 2$ and $3 \leq N \leq 2t - 1$. Then there is no t -resilient consensus algorithm which achieves weak unanimity for binary values.

For the case of strong unanimity and Byzantine faults with authentication, but where the initial values are not signed by a general, Theorems 5.1 and 5.2, part (b), give consensus algorithms if $N \geq 3t+1$. The following shows that this resiliency is the best possible for this case.

Theorem 6.7. Assume the model with Byzantine faults with authentication, synchronous communication, and partially synchronous processors. Assume $t \geq 1$ and $3 \leq N \leq 3t$. If the general does not sign the initial values, there is no t -resilient consensus protocol which achieves strong unanimity for binary values.

Proof. Assume $N \leq 3t$. Divide the processors into three groups, P, Q and R, each containing at least 1 and at most t processors.

Consider the following Scenario A. Processors in P have initial values 0, processors in Q have initial values 1, processors in P \cup Q wake up at time 1 and run fast thereafter, and processors in R are initially dead. Therefore, the processors in P \cup Q must make some decision after some finite time. By symmetry we can assume, without loss of generality, that they decide 1 within time T_A .

Consider Scenario B. All processors have initial value 0, processors in R are correct but do not wake up until after time T_A , and processors in Q are Byzantine and behave with respect to P exactly as they do in Scenario A. The processors in group P act exactly as they do in Scenario A, so they decide 1. This contradicts strong unanimity. \square

7. OPEN QUESTIONS

(1) For the case of partially synchronous models and Byzantine faults without authentication, we have shown that any t -resilient consensus protocol which reaches consensus within polynomial time (after GST) needs at most $4t+1$ processors, and it is known that $3t+1$ processors are required (since [LSP] proves this even in a completely synchronous system without the polynomial time constraint). It would be interesting to close this gap.

(2) We have noted in Remark 1 at the end of Section 3 that the basic consensus Algorithms 1, 2, and 3, with minor modifications, have the property that the number of rounds required to reach agreement after round GST is optimal to within constant factors (at most 10). We have not tried to reduce these constants. Some reduction is probably possible, say by overlapping trying phases with lock release phases, although it would be surprising if the number of rounds could be made to match the known lower bound of $t+1$ rounds. On the other hand, partial synchrony might provide a model for which the lower bound $t+1$ could be strengthened to something larger.

(3) A general direction for future research is to study other distributed computing problems in partially synchronous models.

Acknowledgment. Joe Halpern asked whether the impossibility results of [FLP, DDS] would continue to hold in case the parameters Φ or Δ exist but are not known *a priori*, and this led to the formulation of the version of partial synchrony where Φ or Δ are unknown. We are grateful to Jennifer Lundelius who read a draft of this paper and provided many helpful comments.

REFERENCES

- [ADG] Attiya, A., Dolev, D., and Gil, J., "Asynchronous Byzantine consensus," *Proc. 3rd ACM Symp. on Principles of Distributed Computing*, 1984, pp. 119-133.
- [BT] Bracha, G., and Toueg, S., "Resilient consensus protocols," *Proc. 2nd ACM Symp. on Principles of Distributed Computing*, 1983, pp. 12-26.
- [D] Dolev, D., personal communication.
- [DDS] Dolev, D., Dwork, C., and Stockmeyer, L., "On the minimal synchronism needed for distributed consensus," *Proc. 24th Symp. on Foundations of Computer Science*, 1983, pp. 393-402.
- [DFFLS] Dolev, D., Fischer, M.J., Fowler, R., Lynch, N.A., and Strong, H.R., "Efficient Byzantine agreement without authentication," *Information and Control* 52 (1982), pp. 257-274.
- [DLPSW] Dolev, D., Lynch, N., Pinter, S., Stark, E. and Weihl, W., "Reaching approximate agreement in the presence of faults," *Proc. 3rd Annual IEEE Symposium on Reliability in Distributed Software and Database Systems*, 1983; revised version to appear in *J.ACM*.
- [DLS] Dwork, C., Lynch, N., Stockmeyer, L., "Consensus in the Presence of Partial Synchrony," *Proc. 3rd ACM Symp. on Principles of Distributed Computing*, 1984, pp. 103-118.
- [DS] Dolev, D. and Strong, H. R., "Authenticated algorithms for Byzantine agreement," *SIAM J. Computing* 12 (1983), pp. 656-666.
- [F] Fischer, M.J., "The consensus problem in unreliable distributed systems (a brief survey)," Report YALEU/DCS/RR-273, Yale University, June, 1983.
- [FLa] Fischer, M.J. and Lamport, L., "Byzantine generals and transaction commit protocols," manuscript.
- [FL] Fischer, M.J. and Lynch, N., "A lower bound for the time to assure interactive consistency," *Info. Proc. Lett.* 14, 4 (1982), pp. 183-186.
- [FLP] Fischer, M.J., Lynch, N. A. and Paterson, M.S., "Impossibility of distributed consensus with one faulty process," *J.ACM* 32 (1985), pp. 374-382.
- [G] Garcia-Molina, H., Pittelli, F., and Davidson, S., "Is Byzantine agreement useful in a distributed database?," *Proc. 3rd SIGACT-SIGMOD Symp. on Principles of Database Systems*, 1984, pp. 61-69.
- [GLPT] Gray, J.N., Lorie, R.A., Putzulo, G.R., and Traiger, I.L., "Notes on Database Operating Systems," *Operating Systems: An Advanced Course, Lecture Notes in Computer Science*, Vol 60, Springer-Verlag, New York, 1978, pp. 393-481.

- [L1] Lamport, L., "The weak Byzantine generals problem," *J.ACM* 30 (1983), pp. 668-676.
- [L2] Lamport, L., "Time, clocks, and the ordering of events in a distributed system," *CACM* 21, No. 7, (1978), pp. 558-564.
- [LSP] Lamport, L., Shostak, R., and Pease, M., "The Byzantine generals problem," *ACM Trans. on Programming Languages and Systems* 4 (1982), pp. 382-401.
- [P] Pinter, S., "Distributed computation systems: modelling, verification and algorithms," PhD Thesis, Boston University, 1984.
- [R] Reischuk, R., "A new solution for the Byzantine generals problem," IBM Report RJ3673, November, 1982.
- [Sc] Schneider, F.B., "Byzantine generals in action: implementing fail-stop processors," Manuscript, August, 1983.
- [Sk] Skeen, D., "A quorum based commit protocol," Report TR 82-483, Dept. of Computer Science, Cornell Univ., Feb., 1982.

