

# UC San Diego

## Technical Reports

### Title

Synchronous Consensus for Dependent Process Failures

### Permalink

<https://escholarship.org/uc/item/19h274ng>

### Authors

Junqueira, Flavio  
Marzullo, Keith

### Publication Date

2002-10-03

Peer reviewed

# Synchronous Consensus for Dependent Process Failures <sup>\*</sup>

Flavio P. Junqueira  
flavio@cs.ucsd.edu

Keith Marzullo  
marzullo@cs.ucsd.edu

University of California, San Diego  
Department of Computer Science and Engineering  
9500 Gilman Drive  
La Jolla, CA

## Abstract

We present a new abstraction to replace the  $t$  of  $n$  assumption used in designing fault-tolerant algorithms. This abstraction models dependent process failures yet it is as simple to use as the  $t$  of  $n$  assumption. To illustrate this abstraction, we consider Consensus for synchronous systems with both crash and arbitrary process failures. By considering failure correlations, we are able to reduce latency and enable the solution of Consensus for system configurations in which it is not possible when forced to use protocols designed under the  $t$  of  $n$  assumption. We give lower bounds for the number of rounds and replication requirements that are sufficient to solve Consensus. We show that, in general, the lower bound for number of rounds in the worst case assuming crash failures is different from the lower bound assuming arbitrary failures given the same system configuration. This is in contrast with the traditional result under the  $t$  of  $n$  assumption.

**Keywords:** Distributed Systems, Fault Tolerance, Correlated Failures, Consensus, Time Complexity

## 1 Introduction

Most fault-tolerant protocols are designed assuming that out of  $n$  components, no more than  $t$  can be faulty. For example, solutions to the Consensus problem are usually developed assuming no more than  $t$  of the  $n$  processes are faulty where “being faulty” is specialized by a failure model. We call this the  $t$  of  $n$  *assumption*. It is a convenient assumption to make. For example, bounds are easily expressed as a function of  $t$ : if processes can fail only by crashing, then the Consensus problem is

---

<sup>\*</sup>This work was developed in the context of RAMP, which is DARPA project number N66001-01-1-8933.

solvable when  $t < n$  if the system is synchronous and when  $t < 2n$  if the system is asynchronous extended with a failure detector of the class  $\diamond W$ . [1, 2]

The use of the  $t$  of  $n$  assumption dates back to the earliest work on fault-tolerant computing. [3] It was first applied to distributed coordination protocols in the SIFT project [4] which designed a fly-by-wire system. The reliability of systems like this is a vital concern, and using the  $t$  of  $n$  assumption allows one to represent the probabilities of failure in a simple manner. For example, if each process has a probability  $p$  of being faulty, and processes fail independently, then the probability  $P(t)$  of no more than  $t$  out of  $n$  processes being faulty is:

$$P(t) = \sum_{i=0}^t \binom{n}{i} p^i (1-p)^{n-i}$$

If one has a target reliability  $R$  then one can choose the smallest value of  $t$  that satisfies  $P(t) \geq R$ .

The  $t$  of  $n$  assumption is best suited for components that have identical probabilities of failure and that fail independently. For embedded systems built using rigorous software development this is often a reasonable assumption, but for most modern distributed systems it is not. Process failures can be correlated because, for example, the same buggy software was used. [5] Computers in the same room are subject to correlated crash failures in the case of a power outage.

That failures can have different probabilities and can be dependent is not a novel observation. The continued popularity of the  $t$  of  $n$  assumption, however, implies that it is an observation that is being overlooked by protocol designers. If one wishes to apply, for example, a Consensus protocol in some real distributed system, one can use one of two approaches:

1. Use some off-line analysis technique, such as fault tree analysis [6] to identify how processes fail in a correlated manner. For those that do not fail independently or fail with different probabilities, re-engineer the system so that failures are independent and identically distributed (IID).
2. Use the same off-line analysis technique to compute what the maximum number of faulty processes can be, given a target reliability. Use this value for  $t$  and compute the value of  $n$  that, under the  $t$  of  $n$  assumption, is required to implement Consensus. Replicate to that degree.

Both of these approaches are used in practice. [6] This paper advocates a third approach:

3. Use the same off-line analysis to identify how processes fail in a correlated manner. Represent this using our abstraction for dependent failures, and replicate in a way that satisfies our

replication requirement and that minimizes the number of replicas. Instantiate the appropriate dependent failure protocol.

We believe that our approach and protocols are amenable to on-line adaptive replication techniques as well.

In this paper we propose an abstraction that exposes dependent failure information for one to take advantage of in the design of a protocol. Like the  $t$  of  $n$  assumption, it is expressed in a way that hides its underlying probabilistic nature in order to make it more generally applicable.

We then apply this abstraction to the Consensus problem under the synchronous system assumption and for both crash and arbitrary failures. We derive a new lower bound for the number of rounds to solve Consensus and show that it takes, in general, fewer rounds for crash failures in the worst case than for arbitrary failures. This is in contrast to the  $t$  of  $n$  assumption, where the number of rounds required in the worst case is the same for both failure models. [1, 7] We show that these bounds are tight by giving protocols that meet them. These protocols are fairly simple generalizations of well-known protocols, which is convenient. Proving them correct gave us new insight into the structure of the original protocols. We also show that expressing process failure correlations with our model enables the solution of Consensus in some systems in which it is impossible when making the  $t$  of  $n$  assumption.

There has been some work in providing abstractions more expressive than the  $t$  of  $n$  assumption. The hybrid failure model (for example, [8]) generalizes the  $t$  of  $n$  assumption by providing a separate  $t$  for different classes of failures. Using a hybrid failure model allows one to design more efficient protocols by having sufficient replication for masking each failure class. It is still based on failures in each class being independent and identically distributed. In this paper, however, we do not consider hybrid failure models.

Byzantine Quorum systems have been designed around the abstraction of a *Fail-prone System* [9]. This abstraction allows one to define quorums that take correlated failures into account. This abstraction has been used to express a sufficiency condition for replication. Our work can be seen as generalizing this work, which applies only to Quorum Systems.

The remainder of this paper is divided as follows. Section 2 presents our assumptions for the system model and introduces our abstraction that models dependent process failures. Section 3 defines the distributed Consensus problem. In Section 4, we state a theorem that generalizes the lower bound on the number of rounds in our model. Sections 5 and 6 describe tight weakest replication requirements

and algorithms for Consensus on the crash and arbitrary failure models, respectively. A discussion on the implementation of our abstraction in real systems is provided in Section 7. Finally, we draw conclusions and discuss future work in Section 8.

Due to lack of space, we give proof sketches for some lemmas and theorems and omit them entirely for others. Detailed proofs can be found in [10, 11, 12].

## 2 System Model

A system is composed of a set  $\Pi$  of processes, numbered from 1 to  $n = |\Pi|$ . The number assigned to a process is its process *id*, and it is known by all the other processes. In the rest of the paper, every time we refer to a process with id  $i$ , we use the notation  $p_i$ . Additionally, we define  $Pid$  as the set of process id's, i.e.,  $Pid = \{i : p_i \in \Pi\}$ . We use this set to define a sequence  $w$  of process id's. Such a sequence  $w$  is an element of  $Pid^*$ .

A process communicates with others by exchanging messages. Messages are transmitted through point-to-point reliable channels, and each process is connected to every other process through one of these channels. Processes, on the other hand, are not assumed to be reliable. We consider both crash and arbitrary process failures. In contrast to most previous works in fault-tolerant distributed systems, process failures are allowed to be correlated.

Each process  $p \in \Pi$  executes a deterministic automaton as part of the distributed computation [2, 13]. A deterministic automaton is composed of a set of states, an initial state, and a transition function. The collection of the automata executed by the processes is defined as a distributed algorithm. An execution of a distributed algorithm proceeds in steps of the processes. In a step, a process may: 1) receive a message; 2) undergo a state transition; 3) send a message to a single process. Steps are atomic, and steps of different processes are allowed to overlap in time. We assume that there is an external device that provides the time a process takes a step. The time a process takes a step can be used in proofs, but processors do not have access to this device. The range of time is the non-negative integers.

As discussed later in this section, we assume that the computation can be split into synchronous rounds. The algorithms we describe here proceed in rounds.

## 2.1 Replacing the $t$ of $n$ Assumption: Cores and Survivor Sets

In our model, process failures are allowed to be correlated, which means that the failure of a process may indicate an increase in the failure probability of another process.

Assuming that failed processes do not recover, to achieve fault-tolerance in a system with a set of processes  $\Pi$ , it is necessary to guarantee in every execution that a non-empty subset of  $\Pi$  survives. A process is said to survive an execution if and only if it is correct throughout that execution. Thus, we would like to distinguish subsets of processes such that the probability of all processes in such a subset failing is negligible. Moreover, we want these subsets to be minimal in that removing any process of such a subset  $c$  makes the probability that the remaining processes in  $c$  fail not negligible. We call these minimal subsets *cores*. Cores can be extracted from the information about process failure correlations. In this paper, however, we assume that the set of cores is provided as part of the system's specification. We present in Section 7 a discussion on the problem of finding cores.

By assumption, at least one process in each core will be correct in an execution. Thus, a subset of processes that has a non-empty intersection with every core contains processes that are correct in some execution. If such a subset is minimal, then it is called a survivor set. Notice that in every run of the system there is at least one survivor set that contains only correct processes. The definition of survivor sets is equivalent to the definition of a *fail-prone system*  $\mathcal{B}$  [9]: the set of all survivor sets is the complement of  $\mathcal{B}$ .

We now define cores and survivor sets more formally. Let  $R$  be a rational number expressing a desired reliability, and  $r(x)$ ,  $x \subseteq \Pi$ , be a function that evaluates to the reliability of the subset  $x$ . We define cores and survivor sets as follows:

**Definition 2.1** Given a set of processes  $\Pi$  and rational target degree of reliability  $R \in [0, 1]$ , the set of processes  $c$  is a *core* of  $\Pi$  if and only if:

1.  $c \subseteq \Pi$ ;
2.  $r(c) \geq R$ ;
3.  $\forall p \in c, r(c - \{p\}) < R$ .

Given a set of processes  $\Pi$  and a set of cores  $C_\Pi$ ,  $s$  is a survivor set if and only if:

1.  $s \subseteq \Pi$ ;
2.  $\forall c \in C_\Pi, s \cap c \neq \emptyset$ ;
3.  $\forall p_i \in s, \exists c \in C_\Pi$  such that  $p_i \in c$  and  $(s - \{p_i\}) \cap c = \emptyset$ .

We define  $C_\Pi$  and  $S_\Pi$  to be the set of cores and survivor sets of  $\Pi$ , respectively.

The function  $r(\cdot)$  and the target degree of reliability  $R$  are used at this point only to formalize the idea of a core. In reality, reliability does not necessarily need to be expressed as a probability. If this information is known by other means, then cores can be directly determined. For example, consider the following six process system:

**Example 2.2 :**

- $\Pi = \{ph_1, ph_2, pl_1, pl_2, pl_3, pl_4\}$
- $C_\Pi = \{\{ph_1, ph_2, pl_1\}, \{ph_1, ph_2, pl_2\}, \{ph_1, ph_2, pl_3\}, \{ph_1, ph_2, pl_4\}\}$
- $S_\Pi = \{\{ph_1\}, \{ph_2\}, \{pl_1, pl_2, pl_3, pl_4\}\}$

In this system,  $ph_1$  and  $ph_2$  are highly reliable and both fail independently of every other  $p \in \Pi$ . On the other hand, processes  $pl_1, pl_2, pl_3, pl_4$  fail dependently among each other. That is, for every pair of processes  $pl_i, pl_j$ ,  $1 \leq i, j \leq 4$  and  $i \neq j$ , we have that if  $pl_i$  is faulty in some execution of the system, then  $pl_j$  is also faulty. Thus, a subset with maximum reliability contains processes  $ph_1, ph_2$ , and at least one process  $pl_i$ . Suppose that the maximum reliability achievable for a subset of processes satisfies the intuitive notion of target degree of reliability for this system. We can therefore infer that for each  $i$ ,  $1 \leq i \leq 4$ ,  $\{ph_1, ph_2, pl_i\}$  is a core. From the set  $C_\Pi$  of cores, it is straightforward to identify the survivor sets in  $S_\Pi$ .

In the following sections, we assume that these subsets are provided as part of the system's representation. A system is henceforth described by a triple  $\langle \Pi, C_\Pi, S_\Pi \rangle$ , where  $\Pi$  is a set of processes,  $C_\Pi$  is a set of cores of  $\Pi$ , and  $S_\Pi$  is a set of survivor sets of  $\Pi$ . From this point on, we call  $\langle \Pi, C_\Pi, S_\Pi \rangle$  a *system representation*.

## 2.2 Failure Models

We assume two different models for process failures: crash model and arbitrary model. In the crash model, processes fail by crashing. A crashed process does not send or receive messages. We say that a process is *alive* at time  $t$  either if it is correct in the run or it has not crashed at any time  $t' \leq t$ . In the arbitrary model, a faulty process can take any action, including not receiving messages, sending messages that are not legal under the protocol specification, and sending correct messages at incorrect times. The arbitrary model is strictly weaker than the crash model.

Independently of the assumption for process failures, channels are assumed to be reliable. A reliable channel is one that satisfies the following properties:

**Validity:** If  $p, q \in \Pi$  are correct processes and  $p$  sends a message  $m$  to  $q$ , then  $q$  eventually delivers  $m$ ;

**Integrity:** A process  $p \in \Pi$  receives a message  $m$  from process  $q \in \Pi$  if and only if process  $q$  sent it to  $p$ .

Our Consensus algorithm for crash failures relies on the Validity property to detect crashed processes. This property enables a solution that requires fewer steps of processes in the case of a small number of failures. For arbitrary process failures, the Integrity Property prevents faulty processes from impersonating other ones. Consequently, a faulty process  $p_i$  cannot send a valid message with the id of another process  $p_j$  to other processes.

### 2.3 Synchronous Model

A synchronous system imposes bounds on message delay, process speed, and clock drift. These bounds, however, are not necessarily based on absolute time. As in the model of Dolev *et al.* [14], steps of an algorithm are used to define these bounds. One can then organize an execution into rounds of message exchange. In each round, a process: 1) sends messages at the beginning of the round; 2) receives messages that other processes send at the beginning of the round; 3) changes its state.

The algorithms for synchronous systems described in Sections 5 and 6 are round-based. This format facilitates understanding, since it abstracts several details of the system model. The algorithms are also not described in an automaton format, since the description would be longer and would not improve clarity. Instead, we use sequential code to present the algorithms. States and transitions, however, are easily observed from the changes of the values stored by the variables of the algorithm.

### 2.4 Executions

We define an *execution*  $\alpha$  of an algorithm  $\mathcal{A}$  with the tuple  $\langle F_\alpha, I_\alpha, E_\alpha, T_\alpha \rangle$ . This definition is based on the one by Chandra and Toueg [2] and Charon-Bost *et al.* [15].  $F_\alpha(t)$  evaluates to the subset of processes that have failed by time  $t$ . A direct implication of this definition is that  $F_\alpha(t) \subseteq F_\alpha(t+1)$ . Because an execution depends on the initial state of the processes, we have that  $I_\alpha$  provides the initial configuration of the system. This initial configuration depends on the problem being solved. The



Consensus problem, for example, requires every process to have an initial proposed value.  $E_\alpha$  is an infinite sequence of steps of the processes in  $\Pi$ . The time  $t$  at which a step  $e \in E_\alpha$  is executed is given by  $T_\alpha(e)$ . For every correct process  $p_i$  in  $\alpha$ , we assume that  $E_\alpha$  contains an infinite number of steps of  $p_i$ .<sup>1</sup>

Although we do not use explicitly this definition of execution throughout the paper, we refer several times to executions of algorithms. Therefore, this definition makes clear to the reader what we mean by an execution.

### 3 Consensus

The Consensus problem in a fault-tolerant message-passing distributed system consists, informally, in reaching agreement among a set of processes upon a value. Each process starts with a proposed value and the goal is to have all non-faulty processes decide on the same value. The set of possible decision values is denominated  $V$  throughout this paper. For many applications, a binary set  $V$  is sufficient, but we assume a set  $V$  of arbitrary size, to keep the definition as general as possible.

In the crash failure model, Consensus is often specified in terms of the following three properties [16]:

**Validity** If some non-faulty process  $p \in \Pi$  decides on value  $v$ , then  $v$  was proposed by some process  $q \in \Pi$ ;

**Agreement** If two non-faulty processes  $p, q \in \Pi$  decide on values  $v_p$  and  $v_q$  respectively, then  $v_p = v_q$ ;

**Termination** Every correct process eventually decides.

The Validity property as specified above assumes that no process will ever try to “cheat” on its proposed value. This is true in the crash failure model, but unrealistic in the arbitrary failure model. Although a faulty process might not be able to prevent agreement by cheating on its proposed value, it may prevent progress of the system as whole. For example, assuming that the only possible decision values are either write or abort, with the above Validity definition, a faulty process may prevent correct processes from writing and consequently making progress. Thus, in a byzantine model, Strong Validity is usually considered instead of Validity [17, 18]. Strong Validity is stated as follows:

---

<sup>1</sup>Distributed Consensus requires that processes eventually decide. Because we are assuming that every correct process takes an infinite number of steps, every correct process executes null steps once it halts.

**Strong Validity** If the proposed value of process  $p$  is  $v$ , for all  $p \in \Pi$ , then the only possible decision value is  $v$ .

Strong Validity only considers the case in which all processes have the same initial value. Intuitively, this is sufficient to prevent a byzantine process from disrupting the normal behavior of a system when all non-faulty processes are enabled to make progress. When the system is facing problems and not all of the processes propose the same value, however, this property allows the decision value to be arbitrary in the set of possible decision values. That is, the decision value  $v \in V$  of correct processes can be either the value proposed by a faulty process or even a value that was not proposed by any process, assuming the set of decision values is not binary.

## 4 Lower Bound on the Number of Rounds

Consider a synchronous system in which the  $t$  of  $n$  assumption holds for process failures. In such a system,  $t$  is the maximum number of process failures among all possible executions and  $f$  is the number of failures of a particular execution. It is well known that for every synchronous Consensus algorithm  $\mathcal{A}$ , there is some execution in which some correct process does not decide earlier than  $f + 1$  rounds, where  $f \leq t \leq n - 2$ . [13, 19, 20] Furthermore, there is some execution in which some correct process does not stop earlier than  $\min(t + 1, f + 2)$  rounds, for  $t \leq n - 2$ . [21]

These lower bounds were originally proved for crash failures, but they have to hold for arbitrary failures as well because the arbitrary model is strictly weaker than the crash model. In our model for dependent failures, however, the lower bound on the number of rounds in general differs between these two models.

Before generalizing the lower bound on the number of rounds for our model of dependent failures, we define the term *subsystem*. Let  $\Lambda$  be some predicate that defines the process replication requirement for a given failure model. For example, assuming  $t$  of  $n$  arbitrary process failures, the replication requirement is  $n > 3t$ . Examples of such predicates in our model for dependent failures are provided in Sections 5 and 6. A subsystem of a system that satisfies  $\Lambda$  is then defined as follows:

**Definition 4.1** Let  $\Lambda$  be a replication requirement and  $sys = \langle \Pi, C_\Pi, S_\Pi \rangle$  be a system representation. A system  $sys'$  represented by  $\langle \Pi', C'_\Pi, S'_\Pi \rangle$  is a subsystem of  $sys$  if and only if  $\Pi' \subseteq \Pi$ ,  $C'_\Pi \subseteq C_\Pi$ , and  $sys'$  satisfies  $\Lambda$ .

A subsystem  $sys'$  represented by  $\langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  is minimal if and only if there is no other subsystem  $sys''$  represented by  $\langle \Pi'', C''_{\Pi}, S''_{\Pi} \rangle$  of  $sys$  such that  $|\Pi''| < |\Pi'|$  or  $|C''_{\Pi}| < |C'_{\Pi}|$ .

The following theorem generalizes the lower bound on the number of rounds:

**Theorem 4.2** *Let  $sys = \langle \Pi, C_{\Pi}, S_{\Pi} \rangle$  be the representation of a synchronous system,  $sys' = \langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  be the representation of a minimal subsystem of  $sys$ ,  $\mathcal{A}$  be a Consensus algorithm, and  $\kappa = |\Pi'| - \min\{|s| : s \in S'_{\Pi}\}$ . There are two cases to be considered:*

- i. If  $|\Pi| - \kappa > 1$ , then there is an execution of  $\mathcal{A}$  in which  $f \leq \kappa$  processes are faulty and some correct process takes at least  $f + 1$  rounds to decide;*
- ii. If  $|\Pi| - \kappa = 1$ , then there is an execution of  $\mathcal{A}$  in which  $f \leq \kappa$  processes are faulty and some correct process takes at least  $\min(\kappa, f + 1)$  rounds to decide.*

To illustrate the utilization of this theorem, consider a system  $sys = \langle \Pi, C_{\Pi}, S_{\Pi} \rangle$  under the  $t$  of  $n$  assumption crash failure model. If we assume that  $|\Pi| = n \geq t + 2$ , then  $|C_{\Pi}| \geq 2$  and every core has size  $t + 1$ . A minimal subsystem represented by  $sys' = \langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  has  $n' = |\Pi'| = t + 1$ ,  $|C'_{\Pi}| = 1$ , and  $|S'_{\Pi}| = t + 1$  (each survivor set  $s \in S'_{\Pi}$  contains a single process). From Theorem 4.2(i), for every Consensus algorithm  $\mathcal{A}$ , there is an execution with  $f \leq \kappa$  failures in which no process decides before round  $f + 1$ . The value of  $\kappa$  is  $|\Pi'| - \min\{|s| : s \in S'_{\Pi}\} = t + 1 - 1 = t$ . This result matches the one given by theorem 3.2 in [20]. If we instead assume that  $|\Pi| = t + 1$ , then  $C_{\Pi}$  contains a single core and  $sys$  is already minimal. By Theorem 4.2(ii), we have that  $\kappa = |\Pi| - 1$ . For some execution  $\alpha$  of  $\mathcal{A}$  with  $f \leq \kappa$  failures, there is some correct process that does not decide earlier than round  $\min(|\Pi| - 1, f + 1)$ .

We use Theorem 4.2 in Sections 5 and 6 to derive lower bounds on the number of rounds for the crash and arbitrary models respectively.

## 5 Synchronous Consensus with Crash Failures

Consensus in a synchronous system with crash process failures is solvable for any number of failures. [20] In the case that all processes may fail in some execution before agreement is reached, though, it is often necessary to recover the latest state prior to total failure for recovery purposes. [22] Since we assume that failed processes do not recover, we don't consider total failure in this work. That is, we assume that the following condition holds for a system representation  $\langle \Pi, C_{\Pi}, S_{\Pi} \rangle$ :

**Property 5.1**  $C_{\Pi} \neq \emptyset$ 

Property 5.1 implies that there is at least one correct process in any execution. A core is hence a minimal subsystem in which Consensus is solved. Consider a synchronous system with crash failures  $sys = \langle \Pi, C_{\Pi}, S_{\Pi} \rangle$ , and a subsystem  $sys' = \langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  of  $sys$  such that  $\Pi' = c_{min}$ ,  $c_{min} \in C_{\Pi}$  and  $(\forall c' \in C_{\Pi}, |c_{min}| \leq |c'|)$ . By definition,  $sys'$  is minimal. From Theorem 4.2(ii), if  $|\Pi| = |\Pi'|$ , then there is some execution with  $f \leq |\Pi| - 1$  process failures such that a correct process does not decide earlier than round  $\min(\kappa, f + 1)$ . On the other hand, if  $sys$  is not minimal, then there is some execution with  $f \leq |\Pi'| - 1$  process failures such that a correct process does not decide earlier than round  $f + 1$  by Theorem 4.2(i).

We now describe a protocol for a synchronous system represented by  $\langle \Pi, C_{\Pi}, S_{\Pi} \rangle$ , assuming that Property 5.1 holds for this system. The protocol is based on the early-deciding protocols discussed by Charron-Bost and Schiper [20] and by Lamport and Fischer [19]. Algorithms that take the actual number of failures into account are important because they reduce the latency on the common case in which just a few process failures occur. An important observation made by Charron-Bost and Schiper [20] is that there is a fundamental difference between early-deciding protocols and early-stopping protocols for Consensus. In a early-deciding protocol, a process may be ready to decide, but may not be ready to halt, whereas an early-stopping protocol is concerned about the round in which a process is ready to halt. One consequence of this difference, which was already noted in Section 4, is that the lower bounds for deciding and for stopping are not the same.

Our algorithm **SyncCrash** differentiates the processes of a chosen core  $c \in C_{\Pi}$  from the rest of the processes in  $\Pi - c$ . In a round, every process in  $c$  broadcast its knowledge of proposed values to all the other processes, while processes in  $\Pi - c$  just listen to these messages. Processes in  $c$  from which a message is received at round  $r$ , but from which no message is received at round  $r + 1$ , are known to have crashed before sending all the messages of round  $r + 1$ . This observation is used to detect a round in which no process crashes. Process  $p_i \in \Pi$  keeps track of the processes in  $c$  that have crashed in a round, and as soon as  $p_i$  detects a round with no crashes,  $p_i$  can decide. An important observation is that when such a round  $r$  with no crashes happens (by assumption it eventually happens), all alive processes are guaranteed to have the same array of proposed values. Once each process  $p_i$  in  $c$  decides, it broadcasts a message announcing its decision value  $v_i$ . All undecided processes receiving this message decide on  $v_i$  as well. Thus, only two types of messages are necessary in the protocol: messages containing the array of proposed values and decision messages. Because processes in  $c$

broadcast at most one message in every round to all the processes in  $|\Pi|$ , the message complexity is  $O(|c| * |\Pi|)$ . This is, in general, better than the protocols in [19, 20], designed with the  $t$  of  $n$  failure assumption, which have message complexity  $O(|\Pi|^2)$ .

If  $c = \Pi$ , then in every execution of **SyncCrash** with  $f$  process crashes, every correct process decides in at most  $\min(|c| - 1, f + 1)$  rounds. Otherwise, every correct process decides in at most  $f + 1$  rounds. Thus, the lower bound on the number of rounds discussed in Section 4 is tight for crash failures.

The idea of using a subset of processes to reach agreement on behalf of the whole set of processes is not new. The Consensus Service proposed by Guerraoui and Schiper utilizes this concept. [23] Their failure model, however, assumes  $t$  of  $n$  failures, and consequently the subset used to reach agreement is not chosen based on information about correlated failures. This is the main point where our work differs.

Before presenting a pseudo-code of the algorithm, we show a table describing the variables used in the protocol. Table 1 describes the variables, and the pseudo-code of **SyncCrash** is presented in Figure 1. A detailed proof of correctness for **SyncCrash** is provided in [12].

$c \in C_\Pi$	Core set chosen as the one responsible for the decision.
$dec_i \in V \cup \{\perp\}$	A process $p_i$ decides once it sets $dec_i$ .
$d \in \{true, false\}$	Boolean variable indicating whether the process decided in the previous round or not.
$v_i[1 \dots  c ], v_i[j] \in V$	Array of proposed values.
$e_i[1 \dots ( c  - 1)], e_i[r] \subset c$	Array of failed processes. $e_i[r]$ stores the subset of processes detected by $p_i$ as crashed at round $r$ .

Table 1: Variables used in the algorithm **SyncCrash**

The set of rounds assigned to processes in  $|\Pi| - c$  is only effective if this subset is not empty, and sending a message to an empty set of processes is a no-op.

By characterizing correlated process failures with cores and survivor sets, we improve performance both in terms of message and time complexity. For example, consider again the six process system described in Example 2.2. By assuming  $t$  of  $n$  failures,  $t$  must be as large as the maximum number of failures possible in any execution, which is five. Thus, it is necessary to have at least five rounds to solve Consensus in the worst case. By executing **SyncCrash** with a minimum-sized core as  $C$ , only three rounds are necessary in the worst case. In addition, no messages are broadcast by the processes in  $\Pi - c$ . This is different from most protocols designed under the  $t$  of  $n$  assumption [19, 20, 21],

**Algorithm *SyncCrash* for process  $p_i$ :**

**Input:** set  $\Pi$  of processes; set  $C_\Pi$  of cores; initial value  $v \in V$

**Initialization:**  $c \in C_\Pi$ ;  $dec_i \leftarrow \perp$ ;  $d \leftarrow false$

$v_i[1 \dots |c|]$ ,  $v_i[k] = \perp$ ,  $\forall k \in [1 \dots |c|]$ ,  $k \neq i$ . If  $p_i \in c$ ,  $v_i[i] \leftarrow v$   
 $e_i[1 \dots (|c| - 1)]$ ,  $e_i[k] = c$ ,  $\forall k \in [1 \dots (|c| - 1)]$

**Round  $1 \leq r < |c|$ ,  $\forall p_i \in c$ :**

```

if ( $d = false$ ) then
  send( $i, v_i$ ) to all process in  $c$ 
  send( $i, v_i$ ) to all process in  $\Pi - c$ 
else
  send(Decide,  $dec_i$ ) to all processes in  $c$ 
  send(Decide,  $dec_i$ ) to all processes in  $\Pi - c$ 
  halt
upon reception of ( $m = (Decide, dec_j)$ ) do
   $dec_i \leftarrow dec_j$ 
   $d \leftarrow true$ 
upon reception of ( $m = (j, v_j)$ ) do
   $e_i[r] \leftarrow e_i[r] - \{j\}$ 
  for  $k = 1$  to  $|\Pi|$  do
    if ( $v_j[k] \neq \perp$ ) then  $v_i[k] \leftarrow v_j[k]$ 
if ( $((e_i[r - 1] = e_i[r]) \wedge (d = false)) \vee (r = |c| - 1)$ ) then
   $dec_i \leftarrow \min(v_i[k])$ 
   $d \leftarrow true$ 

```

**Round  $|c|$ ,  $\forall p_i \in c$ :**

```

send(Decide,  $dec_i$ ) to all processes in  $\Pi - c$ 
halt

```

**Round  $1 \leq r \leq |c|$ ,  $\forall p_i \in \Pi - c$ :**

```

upon reception of ( $m = (Decide, dec_j)$ ) do
   $dec_i \leftarrow dec_j$ 
  halt
upon reception of ( $m = (j, v_j)$ ) do
   $e_i[r] \leftarrow e_i[r] \cup \{j\}$ 
  for  $k = 1$  to  $|\Pi|$  do
    if ( $v_j[k] \neq \perp$ ) then  $v_i[k] \leftarrow v_j[k]$ 
if ( $(e_i[r - 1] = e_i[r])$ ) then
   $dec_i \leftarrow \min(v_i[k])$ 
  halt

```

Figure 1: Synchronous Consensus for Dependent Crash Failures

although the same idea can be applied by having only a specific subset of  $t + 1$  processes broadcasting messages.

## 6 Synchronous Consensus with Arbitrary Failures

Given a system representation  $\langle \Pi, C_\Pi, S_\Pi \rangle$ , consider the following properties:

**Property 6.1 (Byzantine Partition)** *For every partition  $(A, B, C)$  of  $\Pi$ , at least one of  $A$ ,  $B$ , and*

$C$  contain a core.

**Property 6.2 (Byzantine Intersection)**  $\forall s_i, s_j \in S_\Pi, \exists c_k \in C_\Pi, c_k \subseteq (s_i \cap s_j)$ .

The following theorem states that these two properties are equivalent.

**Theorem 6.3** *Byzantine Partition  $\equiv$  Byzantine Intersection.*

**Proof sketch:**

- Byzantine Partition  $\Rightarrow$  Byzantine Intersection.

We prove the contrapositive. Assume that there are two survivor sets  $s_i, s_j \in S_\Pi$  such that  $(s_i \cap s_j)$  does not contain a core. Consider the following partitioning:  $A = \Pi - s_i$ ,  $B = (s_i \cap s_j)$ , and  $C = (s_i - B)$ . Subset  $A$  cannot contain a core because it has no element from  $s_i$ . By assumption,  $B$  does not contain a core. Because  $C$  contains no elements from  $s_j$ , we have that  $C$  does not contain a core. Thus, none of  $A$ ,  $B$ , or  $C$  contain a core.  $\therefore$

To prove the other direction, we make use of two observations. First, if the Byzantine Intersection property holds, then every survivor set  $s$  contains at least one core. Otherwise the intersection between  $s$  and some other survivor set  $s' \in S_\Pi, s \neq s'$ , cannot contain a core. Second, if a subset  $A$  of processes contains at least one element from every survivor set, then  $A$  contains a core: by definition, in every execution there is at least one survivor set that contains only correct processes. If  $A$  contains at least one element from every survivor set, then in every execution there is at least one correct process in  $A$ .

- Byzantine Intersection  $\Rightarrow$  Byzantine Partition.

We prove this relation by contradiction. Assume that Byzantine Intersection holds and there is a partition  $(A, B, C)$  such that none of  $A$ ,  $B$ , and  $C$  contain a core. If none of these subsets contains a core, then none of them contains either a survivor set or one element from each survivor set  $s' \in S_\Pi$ . Thus, there has to be two distinct survivor sets  $s_1$  and  $s_2$  such that there are no elements of  $s_1$  in  $C$  and no elements of  $s_2$  in  $B$ . Suppose the contrary. If there are no such  $s_1$  or  $s_2$ , then one of two possibilities has to take place, both in which at least one subset contains a core: 1)  $s_1 = s_2$ . In this case,  $A$  contains  $s_1$ ; 2)  $(\forall s \in S_\Pi, (s \cap B) \neq \emptyset) \vee (\forall s \in S_\Pi, (s \cap C) \neq \emptyset)$ .

Assuming therefore that there are such survivor sets  $s_1$  and  $s_2$ , we have that  $(s_1 \cap s_2) \subseteq A$ . By assumption,  $A$  does not contain a core, and consequently  $s_1 \cap s_2$  does not contain a core. This contradicts, however, our assumption that Byzantine Intersection holds.  $\therefore$

□

The utility for having two equivalent conditions becomes clear below. We use the Byzantine Partition property to show that this replication requirement is necessary to solve Consensus in a synchronous arbitrary failure system. The Byzantine Intersection property is assumed by our protocol **SyncByz**.

Byzantine Intersection along with the definition of  $S_{\Pi}$  is equivalent to the replication requirement for blocking writes in Byzantine Quorum Systems identified by Martin *et al.* They show that this requirement is sufficient for such a protocol. [24] Both our requirement and the one identified by Martin *et al.* are weaker than the replication requirement for masking quorum systems. [9] A masking quorum system requires that in every execution at least one quorum contains only correct processes (that is, it contains a survivor set). In addition, for every quorum in a masking quorum system and every pair of failure scenarios, there is at least one process that is not faulty in both scenarios. The Byzantine Intersection property, on the other hand, only requires that the intersection of two survivor sets contains at least one process that is correct in the execution.

## 6.1 Requirement on Process Replication

Byzantine Partition is necessary to solve Strong Consensus in a synchronous system with arbitrary process failures. The informal proof we provide here is based upon the one by Lamport for the  $t$  of  $n$  assumption. [7, 25] We show that, for any algorithm  $\mathcal{A}$ , if there is a partition of the processes into three non-empty subsets such that none of them contain a core, then there is at least one run in which agreement is violated. This is illustrated in figure 2, where we assume the converse and consider three executions  $\alpha$ ,  $\beta$ , and  $\gamma$ .

In execution  $\alpha$ , the initial value of every process is the same, say  $v$ . All the processes in subset  $B$  are faulty, and they all lie to the processes in subset  $C$  about their initial values and the values received from processes in  $A$ . By Strong Validity, running algorithm  $\mathcal{A}$  in such an execution results in all the processes in subset  $C$  deciding  $v$ . Execution  $\beta$  is analogous to execution  $\alpha$ , but instead of every process beginning with a initial value  $v$ , they all have initial value  $v' \neq v$ . Again, by Strong Validity, all processes in  $B$  decide  $v'$ . In execution  $\gamma$ , the processes in subset  $C$  have initial value  $v$ , whereas processes in subset  $B$  have initial value  $v'$ . The processes in subset  $A$  are all faulty and behave for processes in  $C$  as they do in execution  $\alpha$ . For processes in  $C$ , however, processes in  $B$  behave as they do in execution  $\beta$ . Because processes in  $C$  cannot distinguish executions  $\alpha$  from  $\gamma$ , processes in  $C$  must decide  $v$ . At the same time, processes in  $B$  cannot distinguish executions  $\beta$  from  $\gamma$ , and



therefore they must decide  $v'$ . Consequently, there are correct processes which decide differently in execution  $\gamma$ , violating the Agreement property of Strong Consensus.

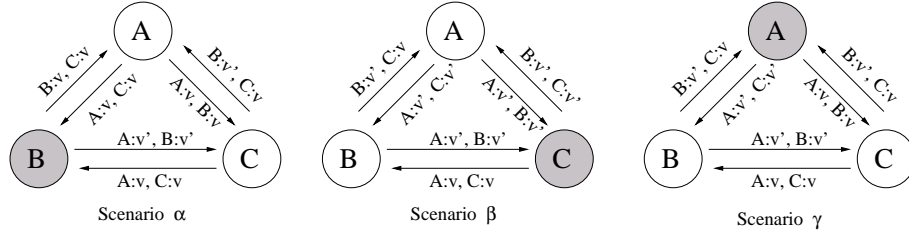


Figure 2: Executions illustrating the violation of Consensus. The processes in shaded subsets are all faulty in the given execution.

## 6.2 Number of Rounds

In every synchronous system with crash failures it suffices to have a single core to solve Consensus. In general, this is not the case for synchronous systems with arbitrary process failures. The only particular case in which Consensus can be solved with a single core is the case that the system has a single reliable process  $p_i$  that does not fail in any execution. For such a system, a minimal subsystem under Byzantine Partition is represented by  $\langle \{p_i\}, \{\{p_i\}\}, \{\{p_i\}\} \rangle$ . In every other case, a system has to contain multiple cores. Although fault-tolerant systems may rely upon a single reliable process, this is a special case.

Assuming a minimal subsystem  $\langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  under Byzantine Partition with multiple cores, every survivor set for such a subsystem contains at least 2 processes. Otherwise, there is a core containing a single process, and it reduces to the particular case described above. By Theorem 4.2(i), the minimum number of rounds required in the worst case is  $\kappa + 1$ , where  $\kappa$  is defined as  $|\Pi'| - \min\{|s| : s \in S'_{\Pi}\}$ . In contrast, all survivor sets of a minimal subsystem have size 1, assuming crash failures.

To illustrate the difference on the total number of rounds in the worst case between the crash and the arbitrary models, consider the following example:

### Example 6.4 :

- $\Pi = \{p_a, p_b, p_c, p_d, p_e\}$
- $C_{\Pi} = \{\{p_a, p_b, p_c\}, \{p_a, p_d\}, \{p_a, p_e\}, \{p_b, p_d\}, \{p_b, p_e\}, \{p_c, p_d\}, \{p_c, p_e\}, \{p_d, p_e\}\}$
- $S_{\Pi} = \{\{p_a, p_b, p_c, p_d\}, \{p_a, p_b, p_c, p_e\}, \{p_a, p_d, p_e\}, \{p_b, p_d, p_e\}, \{p_c, p_d, p_e\}\}$

For the crash model, a minimal subsystem  $\langle \Pi', C'_{\Pi}, S'_{\Pi} \rangle$  is such that  $|\Pi'| = 2$ ,  $|C'_{\Pi}| = 1$ , and a minimum-sized survivor set contains a single process. By Theorem 4.2(i), the lower bound on the number of rounds is 2 in the worst case ( $\kappa = 1$  and  $|\Pi| - \kappa > 1$ ). In the arbitrary model,  $\langle \Pi, C_{\Pi}, S_{\Pi} \rangle$  is already a minimal subsystem: if any process or core is removed, then the remaining subsystem does not satisfy Byzantine Partition. By Theorem 4.2(i), the lower bound on the number of rounds is 3 in the worst case ( $\kappa = 2$  and  $|\Pi| - \kappa > 1$ ). Thus, for the same system configuration, fewer rounds are required assuming crash failures.

### 6.3 An Algorithm to Solve Strong Consensus

We describe an algorithm that solves Strong Consensus in a system  $sys = \langle \Pi, C_{\Pi}, S_{\Pi} \rangle$  that satisfies Byzantine Intersection. This algorithm is based on the one described by Lamport to demonstrate that it is sufficient to have  $3t + 1$  processes ( $t$  is the maximum tolerated number of faulty processes) to have interactive consistency in a setting with arbitrarily faulty processes [7].

In our algorithm, all the processes execute the same sequential code. Every process creates a tree in which each node is labeled with a string  $w$  of distinct process identifiers and in which is stored a value. The value stored in a node labeled  $w$  corresponds to the value forwarded by the sequence of processes named in  $w$ . At round  $r + 1$ , every correct process  $p_j$  sends a message containing the labels and values of the nodes stored at depth  $r$  of the tree to all the other processes. Every correct process  $p_i$  that receives such a message stores the values contained in it in the following manner: if there is a node labeled  $wj$ , with  $w \in \text{Pid}^*$ ,  $|wj| = r + 1$ , then store at this node the value in the message sent by  $p_j$  corresponding to  $w$ .

A simple example will help to clarify the use of the tree. Suppose that a correct process  $p_i$  receives at round three a message from process  $p_j$  that contains the string  $lk$  and the value  $v$  associated to this string. Process  $p_i$  stores the value  $v$  at the node labeled  $lkj$  and forward at round four a message containing the pair  $\langle lkj, v \rangle$  to all the other processes.

The leaves in this tree are survivor sets. More specifically, if we use  $\text{Node}(w)$  to denote the node of the tree labeled with the string  $w$  and  $\text{Processes}(w)$  the set of processes named in  $w$ , then  $\text{Node}(w)$  is a leaf if and only if  $\Pi - \text{Processes}(w)$  does not contain a survivor set. Consequently, if  $\text{Node}(wi)$  is a leaf and we denote with  $\text{Child}(w)$  the set of processes  $\{p_i | \text{Node}(wi) \text{ is a child of } \text{Node}(w)\}$ , then  $\text{Child}(w)$  is a survivor set <sup>2</sup>. For every non-leaf  $\text{Node}(w)$ , we have that  $\Pi - \text{Processes}(w)$  has to contain a survivor

---

<sup>2</sup>Observe that the tree structure is the same for all correct processes, and hence none of  $\text{Processes}(\cdot)$ ,  $\text{Node}(\cdot)$ , or

set. A consequence of this definition is that the depth of the tree is  $|\Pi| - \min\{|s_i| : s_i \in S_\Pi\} + 1$ . Figure 3 gives an example of a tree for the system representation in Example 6.4.

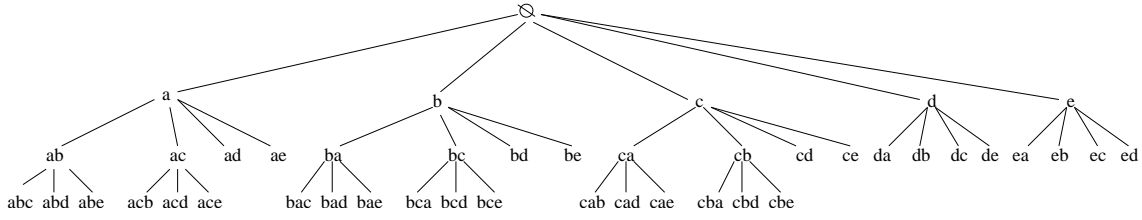


Figure 3: An example of a tree built by each process in the first stage of the algorithm.

The first stage of the algorithm builds and initializes the tree. The second stage runs several rounds of message exchange. In the first round, each process broadcasts its initial value, and in subsequent rounds, each process broadcasts the values it learned in the previous round. As processes receive the messages containing values learned in previous rounds, each node populates the nodes of its tree with these values. Because the depth of the tree is  $(|\Pi| - \min\{|s_i| : s_i \in S_\Pi\} + 1)$ , this is exactly the total number of rounds required for message exchanging. Finally, in the last round, each process traverses the tree visiting the nodes in postorder to decide upon a value. When visiting a leaf, the algorithm does not change the value this node stores. On the other hand, when an internal node of process  $p_i$  with label  $w$  is visited, we use a replacement strategy to determine its value. Suppose there are two survivor sets  $s_1$  and  $s_2$  such that  $(s_1 \cap s_2) \subseteq \text{Child}(w)$  and for every process  $p_j \in (s_1 \cap s_2)$ , we have that  $p_j.\text{Value}(wj) = v$ , for some  $v \in V \cup \{\perp\}$ . In this case, we replace the value of  $\text{Node}(w)$  with  $v$ . Otherwise we replace with the default value ( $\perp$ ). In the original protocol, the replacement strategy is based on the majority. [13]

The pseudo-code of the algorithm is described in Figure 4. In the algorithm, we use  $\text{Value}(w)$  to refer to the value associated to  $w$  both in the tree of a process and in a message some process sends. To differentiate one case from the other, we use a prefix:  $x.\text{Value}(w)$  is the value  $v$  stored at node labeled  $w$  of process  $p_i$  if  $x = p_i$ , whereas it is the value  $v$  in the pair  $\langle w, v \rangle$  in a message  $m$  if  $x = m$ . This is a slight abuse of notation, but it is convenient and the differentiation between the cases will be clear from context.

Instead of providing a formal proof of correctness for **SyncByz**, we illustrate the decision process for the system described in Example 6.4. For a proof of correctness, we point the interested reader  $\text{Child}(\cdot)$  need to be associated with any particular process.

**Algorithm SyncByz for process  $p_i$ :**

**Input:** a set of processes  $\Pi$ ; a set of cores  $C_\Pi$ ; a set of survivor sets  $S_\Pi$ ; an input value  $v_i \in V$

**Variables:**

Let  $s_{min}$  be a smallest survivor set in  $\mathcal{S}$   
 Let  $r$  be the current round number  
 Let  $root$  be a reference to the root of process  $i$ 's tree  
 Let  $M$  be a set of messages  
 Let  $P, P'$  be sets of pairs  $\langle w, v \rangle$ , where  $w \in Pid^*$ , and  $v \in V$

**initialization:**

$root \leftarrow \text{CreateNode}(\emptyset, v_i)$   
 $\text{BuildTree}(root)$   
 $P \leftarrow \{\langle \emptyset, v_i \rangle\}$

**rounds**  $1 \leq r < (|\Pi| - |s_{min}| + 1)$ :

$\text{SendAll}(i, P)$   
**let**  $M$  be the set of messages received by  $p_i$  at round  $r$   
 $P \leftarrow \emptyset$   
**for** every message  $m = (j, P') \in M$  **do**  
   **for** every node at depth  $r$  labeled  $wj$ ,  $w \in Pid^*$ ,  $|w| = r$  **do**  
      $p_i.\text{Value}(wj) \leftarrow m.\text{Value}(w)$   
     **if** node labeled  $wj$  is not a leaf **then**  $P \leftarrow P \cup \{\langle wj, m.\text{Value}(w) \rangle\}$

**round**  $r = (|\Pi| - |s_{min}| + 1)$ :

$\text{SendAll}(i, P)$   
**let**  $M$  be the set of messages received by  $p_i$  at round  $r$   
**for** every message  $m = (j, P') \in M$  **do**  
   **for** every node at level  $r$  labeled  $wj$ ,  $w \in Pid^*$ ,  $|w| = r$ , **do**  
      $p_i.\text{Value}(wj) \leftarrow m.\text{Value}(w)$

Traverse Tree in postorder, executing the following steps when visiting a node labeled  $w$ :

**if**  $\text{Child}(w) \neq \emptyset$   
**then let**  $I \leftarrow \text{Child}(w)$   
   **if**  $(\exists s_1, s_2 \in S \text{ such that } ((s_1 \cap s_2) \subseteq I) \wedge (\forall p_j \in (s_1 \cap s_2), p_i.\text{Value}(wj) = v, v \in V))$   
   **then**  $p_i.\text{Value}(w) \leftarrow v$   
   **else**  $p_i.\text{Value}(w) \leftarrow \perp$

**Auxiliary function**

**Function**  $\text{BuildTree}(w)$   
**let**  $\Gamma \leftarrow \text{Processes}(w)$   
 $\forall p_j \in \Pi$  such that  $p_j \notin \Gamma$   
   **if**  $(\exists s_1 \in S \text{ such that } s_1 \subseteq (\Pi - \Gamma))$   
   **then**  $node \leftarrow \text{CreateNode}(wj, \perp)$   
      $\text{Child}(w) \leftarrow \text{Child}(w) \cup \{node\}$   
      $\text{BuildTree}(wj)$

Figure 4: Synchronous Consensus for Dependent Arbitrary Failures

to [12].

After  $|\Pi| - |s_{min}| + 1 = 5 - 3 + 1 = 3$  rounds of message exchange, every correct process has populated its tree with values received from other processes. The values stored at non-leaf nodes are

not important, because they are replaced according to the strategy defined for the algorithm during the traversal of the tree. We illustrate this procedure for the subtrees rooted at both  $Node(a)$  and  $Node(b)$ . This is shown in Figures 5 and 6. White nodes are the ones that have the same value across all the correct processes, whereas shaded nodes are the ones that have possibly different values across correct processes. A node is shaded if the last node in the string that labels the node is a faulty process. Note that if two nodes  $w$  and  $w'$  are white, it does not mean that they contain necessarily the same value. It only means that every correct process has value  $v$  at node  $w$  and  $v'$  at node  $w'$ .

Consider the particular scenario in which processes  $p_a$  and  $p_c$  are faulty and  $p_b$ ,  $p_d$ , and  $p_e$  are all correct. First, we discuss the subtree rooted at  $Node(a)$ . At Time 1, only the nodes at the last level have been visited. From the algorithm, when a leaf is visited, its value does not change. Thus, the state of the tree at Time 1 is the same state as right before starting the traversal of the tree. Time 2 corresponds to the state of the tree exactly after all the nodes at Level 2 are visited. Because processes  $p_b$ ,  $p_d$ , and  $p_e$  are correct,  $Node(abd)$  and  $Node(abe)$  contain the same value across all correct processes. By the replacement strategy of the algorithm, the new value of node  $ab$  is the value of nodes  $abd$  and  $abe$ , because  $\{p_d, p_e\} \subseteq \{p_a, p_d, p_e\} \cap \{p_b, p_d, p_e\}$  and  $Node(abd)$  contains the same value as node  $abe$ . Similarly, the new value of node  $ac$  is the one of  $acb$ ,  $acd$ , and  $ace$ . The values of  $Node(ad)$  and  $Node(ae)$  across correct processes have to be the same, because  $p_d$  and  $p_e$  are correct. At Time 3, the value of  $Node(a)$  become the same for all correct processes. Since the value of  $Node(ab)$ ,  $Node(ac)$ ,  $Node(ad)$ , and  $Node(ae)$  are the same across all correct processes, the new value of node  $a$  has to be the same.

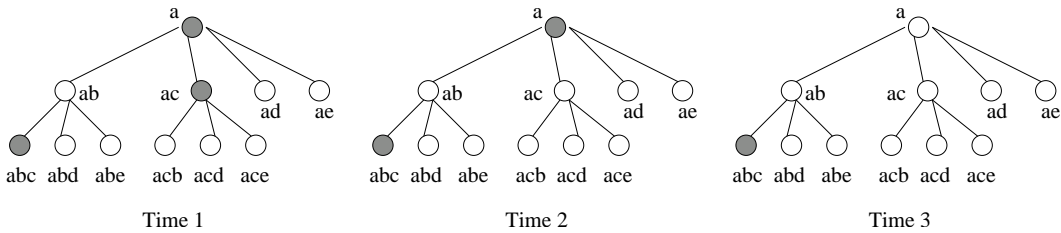


Figure 5: An example of traversing the subtree rooted at  $Node(a)$ . Time  $i$  corresponds to the state of the tree exactly after all the nodes of Level  $4 - i$  are visited.

For the subtree rooted at  $Node(b)$ , the value of  $Node(ba)$  may still not be the same across all correct processes at Time 2. Both  $\{p_d, p_e\}$  and  $\{p_c, p_d\}$  are cores and are subsets of processes contained in some intersection of two survivor sets. Thus, if the value of  $Node(bac)$  is the same of  $Node(bad)$  in a

correct process  $p_i$ , but different in another correct process  $p_j$ , then  $p_i$  and  $p_j$  may replace the value of  $Node(ba)$  with different values, depending on  $Node(bae)$ . Note that one value must be the default  $\perp$  and the other some  $v \in V$ . Similarly for  $Node(bc)$  at Time 2. The values of  $Nodes(bd)$  and  $Nodes(be)$ , however, have to be the same across all correct processes. In Time 3,  $\{p_d, p_e\}$  is the only core in  $Child(b)$  to contain the same value in their respective nodes at Level 1, unless  $ba$  and  $bc$  have the same value as  $Node(bd)$  and  $Node(be)$ . Furthermore, this core is in the intersection of  $\{p_b, p_d, p_e\}$  and  $\{p_a, p_d, p_e\}$ . Consequently, the new value of  $Node(b)$  has to be the same for every correct process at Time 3, by the value replacement strategy.

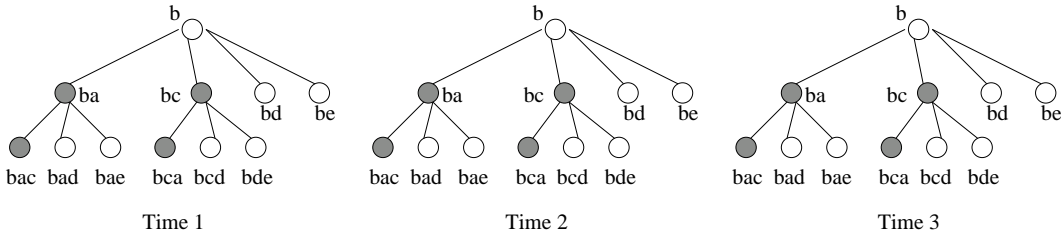


Figure 6: An example of traversing the subtree rooted at  $Node(b)$ . Time  $i$  corresponds to the state of the tree exactly after all the nodes of Level  $4 - i$  are visited.

By doing the same analysis for the subtrees rooted at nodes  $c$ ,  $d$ , and  $e$ , we observe that every node at Level 1 of the tree rooted at  $\emptyset$  has the same value across all correct processes. Therefore, the decision value, which is the value at node  $\emptyset$  after visiting it, has to be the same for every correct process. One important observation is that the value at  $Node(i)$  across all correct processes is the initial value of process  $p_i$ , if  $p_i$  is correct. In the case that every process has the same initial value  $v$ , then the decision value has to be  $v$ .

To illustrate the benefits of using our abstractions, consider once more the five process system of Example 6.4. If the  $t$  of  $n$  failure assumption is used, then Strong Consensus is not solvable: the smallest survivor set contains three processes and so the maximum number of failures in any execution is two. With the  $t$  of  $n$  assumption, the replication requirement is  $|\Pi| \geq 3t + 1$ , and for  $t = 2$ , it is necessary to have at least seven processes. With our model, however, algorithm **SyncByz** solves Strong Consensus.

## 7 Practical Considerations

Two important issues concerning the use of cores and survivor sets are (1) how to extract information about these subsets and (2) how to represent them.

To extract core information (such as finding a smallest core) using failure probabilities is an NP-hard problem in the general case. [10] This result need not be discouraging, however. First, this is a problem that is already addressed for many safety critical systems, and techniques have been developed to extract such information [6]. Furthermore, for many real systems, there are simple heuristics for finding cores that depend on how failure correlations are identified. Suppose a system in which processes are tagged with colors. In this model, all processes have identical probabilities of failing, but those with the same color have highly correlated probabilities of failing while processes with different colors fail independently. A core in such a system is composed of processes with different colors, and the size of a core depends on the probability of having colors failing. To find cores in such a model, one has to interactively add processes with different colors to a subset and verify whether this subset is a core. The verification procedure consists in multiplying the probability of failure for every color that has a representative in the subset. This clearly can be accomplished in polynomial time. For real systems, a color would represent some intrinsic characteristic. For example, all components in a certain part of an airplane are damaged if there is a structural damage on that particular part. Computers in the same room are subject to correlated crash failures in the case of a power outage.

One can go further in extracting cores based on characteristics of the system and propose the utilization of several attributes, instead of one as in the color model. It turns out that in the general case, this problem is also NP-hard. Some simplifying assumptions such as finding orthogonal cores (cores in which processes do not share attributes) make the problem tractable. Finally, fault tree analysis is an option in the design of reliable systems.

Representing cores or survivor sets is relevant for arbitrary failures. As discussed in Section 5, to solve Consensus assuming crash failures, it suffices a single core. Space complexity is hence  $O(1)$  in this case, and consequently is not a problem. For arbitrary failures, however, multiple survivor sets are usually necessary. An important observation is that the number of processes in fault-tolerant systems is usually not large. Thus, for a small number of processes, space complexity is still  $O(1)$ , even if there is an exponential number of survivor sets. In particular cases, it is possible to determine algorithmically whether a subset of processes is a survivor set. Considering the color model once more, a subset  $s$  of processes is a survivor set if  $\Pi - s$  does not contain a core. Whether  $\Pi - s$  contains

a core or not is verifiable in polynomial time, as discussed previously.

## 8 Conclusions and Future Directions

Cores and survivor sets are abstractions capable of using dependent process failure information in the design of fault-tolerant algorithms in a simple manner. We showed this by describing two Consensus algorithms. The main structures of these algorithms were proposed in the literature assuming  $t$  process failures out of  $n$  process. With simple modifications, we obtained new algorithms that in several cases perform better than the original ones. An important observation is that the algorithms we presented improve performance only if there is failure correlation. If all process fail independently, then the protocols behave as the original ones for the  $t$  of  $n$  assumption. In either case, they never have worse performance.

The trade-off, however, is in finding and representing cores or survivor sets. In the general case, finding and representing them require exponential time and space. As we discussed in Section 7, however, this need not hinder the use of cores and survivor sets. Their equivalent are already used in the analysis of safety critical systems. Moreover, there are heuristics that make these problems tractable when the number of cores is not sufficiently small. We believe that many real systems that can benefit from our model either have a small number of cores or are amenable to the application of simplifying heuristics.

So far, we have identified a few real scenarios that would benefit from the application of our model. We believe, however, that our techniques are widely applicable. Aside from identifying other real applications, we are interested in investigating the utilization of cores and survivor sets for other problems of interest in fault-tolerant computing. We already have corresponding results for Consensus in asynchronous systems.

## References

- [1] I. Keidar and S. Rajsbaum, “On the Cost of Fault-Tolerant Consensus When There Are No Faults - A Tutorial,” Tech. Rep. MIT-LCS-TR-821, MIT, May 2001.
- [2] T. Chandra and S. Toueg, “Unreliable Failure Detectors for Reliable Distributed Systems,” *Journal of the ACM*, vol. 43, pp. 225–267, March 1996.



- [3] J. von Neumann, “Probabilistic Logics and the Synthesis of Reliable Organisms from Unreliable Components,” in *Automata Studies*, pp. 43–98, Princeton University Press, 1956.
- [4] J. Wensley, “SIFT: Design and Analysis of a Fault-Tolerant Computer for Aircraft Control,” in *Proceedings of the IEEE*, vol. 66, pp. 1240–1255, October 1978.
- [5] R. Rodrigues, B. Liskov, and M. Castro, “BASE: Using Abstraction to Improve Fault Tolerance,” in *18th ACM Symposium on Operating Systems Principles (SOSP’01)*, vol. 35, (Chateau Lake Louise, Banff, Alberta, Canada), pp. 15–28, October 2001.
- [6] Y. Ren and J. B. Dugan, “Optimal Design of Reliable Systems Using Static and Dynamic Fault Trees,” *IEEE Transactions on Reliability*, vol. 47, pp. 234–244, December 1998.
- [7] L. Lamport, R. Shostak, and M. Pease, “The Byzantine Generals Problem,” *ACM Transactions on Programming Languages and Systems*, vol. 4, pp. 382–401, July 1982.
- [8] P. Thambidurai and Y.-K. Park, “Interactive Consistency with Multiple Failure Modes,” in *IEEE 7th Symposium on Reliable Distributed Systems*, (Columbus, Ohio), pp. 93–100, October 1988.
- [9] D. Malkhi and M. Reiter, “Byzantine Quorum Systems,” in *29th ACM Symposium on Theory of Computing*, pp. 569–578, may 1997.
- [10] F. Junqueira, K. Marzullo, and G. Voelker, “Coping with Dependent Process Failures,” tech. rep., UCSD, La Jolla, CA, December 2001. <http://www.cs.ucsd.edu/users/flavio/Docs/JuMaVo2001.ps>.
- [11] F. Junqueira and K. Marzullo, “Lower Bound on the Number of Rounds for Synchronous Consensus with Dependent Process Failures,” tech. rep., UCSD, La Jolla, CA, September 2002. <http://www.cs.ucsd.edu/users/flavio/Docs/lb.ps>.
- [12] F. Junqueira and K. Marzullo, “Consensus for Dependent Process Failures,” tech. rep., UCSD, La Jolla, CA, September 2002. <http://www.cs.ucsd.edu/users/flavio/Docs/Gen.ps>.
- [13] H. Attiya and J. Welch, *Distributed Computing: Fundamentals, Simulations, and Advanced Topics*. McGraw-Hill, 1998.
- [14] D. Dolev, C. Dwork, and L. Stockmeyer, “On the Minimal Synchronism Needed for Distributed Consensus,” *Journal of the ACM*, vol. 1, pp. 77–97, January 1987.

- [15] B. Charron-Bost, R. Guerraoui, and A. Schiper, “Synchronous System and Perfect Failure Detector: solvability and efficiency issues,” in *IEEE International Conference on Dependable Systems and Networks (DSN’00)*, (New York, NY), pp. 523–532, June 2000.
- [16] A. Doudou and A. Schiper, “Muteness Detectors for Consensus with Byzantine Processes,” in *Proceedings of the 17th ACM Symposium on Principle of Distributed Computing*, (Puerto Vallarta, Mexico), p. 315, July 1998. (Brief Announcement).
- [17] K. Kihlstrom, L. Moser, and P. M. Melliar-Smith, “Solving Consensus in a Byzantine Environment using an Unreliable Failure Detector,” in *Proceedings of the International Conference on Principles of Distributed Systems (OPODIS’97)*, (Chantilly, France), pp. 61–76, December 1997.
- [18] D. Malkhi and M. Reiter, “Unreliable Intrusion Detection in Distributed Computations,” in *Proceedings of the 10th Computer Security Foundations Workshop (CSFW97)*, (Rockport, MA), pp. 116–124, June 1997.
- [19] L. Lamport and M. Fischer, “Byzantine Generals and Transaction Commit Protocols,” tech. rep., SRI International, April 1982.
- [20] B. Charron-Bost and A. Schiper, “Uniform Consensus is Harder Than Consensus,” tech. rep., École Polytechnique Fédérale de Lausanne, Switzerland, May 2000.
- [21] D. Dolev, R. Reischuk, and H. R. Strong, “Early Stopping in Byzantine Agreement,” *Journal of the ACM*, vol. 37, pp. 720–741, October 1990.
- [22] D. Skeen, “Determining the Last Process to Fail,” *ACM Transactions on Computer Systems*, vol. 3, pp. 15–30, February 1985.
- [23] R. Guerraoui and A. Schiper, “Consensus Service: A Modular Approach for Building Fault-tolerant Agreement Protocols in Distributed Systems,” in *26th International Symposium on Fault-Tolerant Computing (FTCS-26)*, (Sendai, Japan), pp. 168–177, June 1996.
- [24] L. Alvisi, P. Martin, and M. Dahlin, “Minimal Byzantine Storage,” in *16th International Symposium on Distributed Computing (DISC 2002)*, (Toulouse, France), October 2002.
- [25] M. Pease, R. Shostak, , and L. Lamport, “Reaching Agreement in the Presence of Faults,” *Journal of the ACM*, vol. 27, pp. pp. 228–234, April 1980.