

The Weakest Failure Detector for Solving Consensus*

Tushar Deepak Chandra[†]

Vassos Hadzilacos[‡]

Sam Toueg[§]

June 1996

Abstract

We determine what information about failures is necessary and sufficient to solve Consensus in asynchronous distributed systems subject to crash failures. In [CT91], it is shown that $\diamond\mathcal{W}$, a failure detector that provides surprisingly little information about which processes have crashed, is sufficient to solve Consensus in asynchronous systems with a majority of correct processes. In this paper, we prove that to solve Consensus, any failure detector has to provide at least as much information as $\diamond\mathcal{W}$. Thus, $\diamond\mathcal{W}$ is indeed the weakest failure detector for solving Consensus in asynchronous systems with a majority of correct processes.

1 Introduction

1.1 Background

The asynchronous model of distributed computing has been extensively studied. Informally, an *asynchronous distributed system* is one in which message transmission times and relative processor speeds are both unbounded. Thus an algorithm designed for an asynchronous system does not rely on such bounds for its correctness. In practice, asynchrony is introduced by unpredictable loads on the system.

Although the asynchronous model of computation is attractive for the reasons outlined above, it is well-known that many fundamental problems of fault-tolerant distributed computing that are solvable in synchronous systems, are unsolvable in asynchronous systems. In particular, it is well-known that *Consensus*, and several forms

*Research supported by NSF grants CCR-8901780 and CCR-9102231, DARPA/NASA Ames grant NAG-2-593, grants from the IBM Endicott Programming Laboratory and Siemens Corp, and a grant from the Natural Sciences and Engineering Research Council of Canada. A preliminary version of this paper appeared in *Proceedings of the Eleventh ACM Symposium on Principles of Distributed Computing*, pages 147–158. ACM press, August 1992.

[†]H2-L10, IBM T. J. Watson Research Center, 30 Saw Mill Road, Hawthorne, NY 10532, USA.

[‡]Computer Systems Research Institute, University of Toronto, 6 King's College Road, Toronto, Ontario M5S 1A1.

[§]Computer Science Department, Upson Hall, Cornell University, Ithaca, NY 14853.

of reliable broadcast, including *Atomic Broadcast*, cannot be solved deterministically in an asynchronous system that is subject to even a single *crash* failure [FLP85, DDS87]. Essentially, these impossibility results stem from the inherent difficulty of determining whether a process has actually crashed or is only “very slow”.

To circumvent these impossibility results, previous research focused on the use of randomization techniques [CD89], the definition of some weaker problems and their solutions [DLP⁺86, ABD⁺87, BW87, BMZ88], or the study of several models of *partial synchrony* [DDS87, DLS88]. Nevertheless, the impossibility of deterministic solutions to many agreement problems (such as Consensus and Atomic Broadcast) remains a major obstacle to the use of the asynchronous model of computation for fault-tolerant distributed computing.

An alternative approach to circumvent such impossibility results is to augment the asynchronous model of computation with a *failure detector*. Informally, a failure detector is a distributed oracle that gives (possibly incorrect) hints about which processes may have crashed so far: Each process has access to a local *failure detector module* that monitors other processes in the system, and maintains a list of those that it currently suspects to have crashed. Each process periodically consults its failure detector module, and uses the list of suspects returned in solving Consensus.

A failure detector module can make mistakes by erroneously adding processes to its list of suspects: i.e., it can suspect that a process p has crashed even though p is still running. If it later believes that suspecting p was a mistake, it can remove p from its list. Thus, each module may repeatedly add and remove processes from its list of suspects. Furthermore, at any given time the failure detector modules at two different processes may have different lists of suspects.

It is important to note that the mistakes made by a failure detector should not prevent any correct process from behaving according to specification. For example, consider an algorithm that uses a failure detector to solve Atomic Broadcast in an asynchronous system. Suppose all the failure detector modules wrongly (and permanently) suspect that a correct process p has crashed. The Atomic Broadcast algorithm must still ensure that p delivers the same set of messages, in the same order, as all the other correct processes. Furthermore, if p broadcasts a message m , all correct processes must deliver m .¹

In [CT91], it is shown that a surprisingly weak failure detector is sufficient to solve Consensus and Atomic Broadcast in asynchronous systems with a majority of correct processes. This failure detector, called the *eventually weak failure detector* and denoted \mathcal{W} here, satisfies only the following two properties:²

1. There is a time after which every process that crashes is always suspected by some correct process.
2. There is a time after which some correct process is never suspected by any correct process.

¹A different approach was taken in [RB91]: a correct process that is wrongly suspected to have crashed, voluntarily leaves the system. It may later rejoin the system by assuming a new identity.

²In [CT91], this is denoted $\diamond\mathcal{W}$.

Note that, at any given time t , processes cannot use \mathcal{W} to determine the identity of a correct process. Furthermore, they cannot determine whether there is a correct process that will not be suspected after time t .

The failure detector \mathcal{W} can make an *infinite* number of mistakes. In fact, it can forever add and then remove some *correct* processes from the lists of suspects (this reflects the inherent difficulty of determining whether a process is just slow or has crashed). Moreover, some correct processes may be erroneously suspected to have crashed by all the other processes throughout the entire execution.

The two properties of \mathcal{W} state that eventually something must hold forever; this may appear too strong a requirement to implement in practice. However, when solving a problem that “terminates”, such as Consensus, it is not really required that the properties hold *forever*, but merely that they hold for a *sufficiently long time*, i.e., long enough for the algorithm that uses the failure detector to achieve its goal. For instance, in practice the algorithm of [CT91] that solves Consensus using \mathcal{W} only needs the two properties of \mathcal{W} to hold for a relatively short period of time.³ However, in an asynchronous system it is not possible to quantify “sufficiently long”, since even a single process step or a single message transmission is allowed to take an arbitrarily long amount of time. Thus it is convenient to state the properties of \mathcal{W} in the stronger form given above.

1.2 The problem

The failure detection properties of \mathcal{W} are *sufficient* to solve Consensus in asynchronous systems. But are they *necessary*? For example, consider failure detector \mathcal{A} that satisfies Property 1 of \mathcal{W} and the following weakening of Property 2:

There is a time after which some correct process is never suspected by at least 99% of the correct processes.

\mathcal{A} is clearly weaker than \mathcal{W} . Is it possible to solve Consensus using \mathcal{A} ? Indeed what is the *weakest* failure detector *sufficient* to solve Consensus in asynchronous systems? In trying to answer this fundamental question we run into a problem. Consider failure detector \mathcal{B} that satisfies the following two properties:

1. There is a time after which every process that crashes is always suspected by *all* correct processes.
2. There is a time after which some correct process is never suspected by a majority of the processes.

It seems that \mathcal{B} and \mathcal{W} are *incomparable*: \mathcal{B} 's first property is stronger than \mathcal{W} 's, and \mathcal{B} 's second property is weaker than \mathcal{W} 's. Is it possible to solve Consensus in an asynchronous system using \mathcal{B} ? The answer turns out to be “yes” (provided that this asynchronous system has a majority of correct processes, as \mathcal{W} also requires). Since \mathcal{W} and \mathcal{B} appear

³In that algorithm processes are cyclically elected as “coordinators”. Consensus is achieved as soon as a correct coordinator is reached, and no process suspects it to have crashed while this coordinator is trying to enforce consensus.

to be incomparable, one may be tempted to conclude that \mathcal{W} cannot be the “weakest” failure detector with which Consensus is solvable. Even worse, it raises the possibility that no such “weakest” failure detector exists.

However, a closer examination reveals that \mathcal{B} and \mathcal{W} are indeed comparable in a natural way: There is a distributed algorithm $T_{\mathcal{B} \rightarrow \mathcal{W}}$ that can transform \mathcal{B} into a failure detector with the Properties 1 and 2 of \mathcal{W} . $T_{\mathcal{B} \rightarrow \mathcal{W}}$ works for any asynchronous system that has a majority of correct processes. We say that \mathcal{W} is *reducible to* \mathcal{B} in such a system. Since $T_{\mathcal{B} \rightarrow \mathcal{W}}$ is able to transform \mathcal{B} into \mathcal{W} in an asynchronous system, \mathcal{B} must provide at least as much information about process failures as \mathcal{W} does. Intuitively, \mathcal{B} is at least as strong as \mathcal{W} .

1.3 The result

In [CT91], it is shown that \mathcal{W} is sufficient to solve Consensus in asynchronous systems if and only if $n > 2f$ (where n is the total number of processes, and f is the maximum number of processes that may crash). In this paper, we prove that \mathcal{W} is reducible to *any* failure detector \mathcal{D} that can be used to solve Consensus (this result holds for any asynchronous system). We show this reduction by giving a distributed algorithm $T_{\mathcal{D} \rightarrow \mathcal{W}}$ that transforms any such \mathcal{D} into \mathcal{W} . Therefore, \mathcal{W} is indeed the weakest failure detector that can be used to solve Consensus in asynchronous systems with $n > 2f$. Furthermore, if $n \leq 2f$, any failure detector that can be used to solve Consensus must be strictly stronger than \mathcal{W} .

The task of transforming any given failure detector \mathcal{D} (that can be used to solve Consensus) into \mathcal{W} runs into a serious technical difficulty for the following reasons:

- To strengthen our result, we do not restrict the output of \mathcal{D} to lists of suspects. Instead, this output can be *any value* that encodes some information about failures. For example, a failure detector \mathcal{D} should be allowed to output any boolean formula, such as “(not p) and (q or r)” (i.e., p is up and either q or r has crashed)—or any *encoding* of such a formula. Indeed, the output of \mathcal{D} could be an arbitrarily complex (and unknown) encoding of failure information. Our transformation from \mathcal{D} into \mathcal{W} must be able to decode this information.
- Even if the failure information provided by \mathcal{D} is not encoded, it is not clear how to extract from it the failure detection properties of \mathcal{W} . Consequently, if \mathcal{D} is given in isolation, the task of transforming it into \mathcal{W} may not be possible.

Fortunately, since \mathcal{D} can be used to solve Consensus, there is a corresponding algorithm, $Consensus_{\mathcal{D}}$, that is somehow able to “decode” the information about failures provided by \mathcal{D} , and knows how to use it to solve Consensus. Our reduction algorithm, $T_{\mathcal{D} \rightarrow \mathcal{W}}$ uses $Consensus_{\mathcal{D}}$ to extract this information from \mathcal{D} and obtain \mathcal{W} .

2 The model

We consider *asynchronous* distributed systems in which there is no bound on message delay, clock drift, or the time necessary to execute a step. Our model of asynchronous

computation with failure detection is patterned after the one in [FLP85]. The system consists of a set of n processes, $\Pi = \{p_1, p_2, \dots, p_n\}$. Every pair of processes is connected by a reliable communication channel.

To simplify the presentation of our model, we assume the existence of a discrete global clock. This is merely a fictional device: the processes do not have access to it. We take the range \mathcal{T} of the clock's ticks to be the set of natural numbers.

2.1 Failures and failure patterns

Processes can fail by *crashing*, i.e., by prematurely halting. A *failure pattern* F is a function from \mathcal{T} to 2^Π , where $F(t)$ denotes the set of processes that have crashed through time t . Once a process crashes, it does not “recover”, i.e., $\forall t : F(t) \subseteq F(t + 1)$. We define $crashed(F) = \bigcup_{t \in \mathcal{T}} F(t)$ and $correct(F) = \Pi - crashed(F)$. If $p \in crashed(F)$ we say p *crashes in* F and if $p \in correct(F)$ we say p *is correct in* F . We consider only failure patterns F such that at least one process is correct, i.e., $correct(F) \neq \emptyset$.

2.2 Failure detectors

Informally, a failure detector provides (possibly incorrect) information about the failure pattern that occurs in an execution. Associated with each failure detector is a (possibly infinite) range \mathcal{R} of values output by that failure detector. A *failure detector history* H with range \mathcal{R} is a function from $\Pi \times \mathcal{T}$ to \mathcal{R} . $H(p, t)$ is the value of the failure detector module of process p at time t . A *failure detector* \mathcal{D} is a function that maps each failure pattern F to a *set* of failure detector histories with range $\mathcal{R}_{\mathcal{D}}$ (where $\mathcal{R}_{\mathcal{D}}$ denotes the range of failure detector outputs of \mathcal{D}). $\mathcal{D}(F)$ denotes the set of possible failure detector histories permitted by \mathcal{D} for the failure pattern F .

For example, consider the failure detector \mathcal{W} mentioned in the introduction. Each failure detector module of \mathcal{W} outputs a *set of processes* that are suspected to have crashed: in this case $\mathcal{R}_{\mathcal{W}} = 2^\Pi$. For each failure pattern F , $\mathcal{W}(F)$ is the set of all failure detector histories $H_{\mathcal{W}}$ with range $\mathcal{R}_{\mathcal{W}}$ that satisfy the following properties:

1. There is a time after which every process that crashes in F is always suspected by some process that is correct in F :

$$\exists t \in \mathcal{T}, \forall p \in crashed(F), \exists q \in correct(F), \forall t' \geq t : p \in H_{\mathcal{W}}(q, t')$$

2. There is a time after which some process that is correct in F is never suspected by any process that is correct in F :

$$\exists t \in \mathcal{T}, \exists p \in correct(F), \forall q \in correct(F), \forall t' \geq t : p \notin H_{\mathcal{W}}(q, t')$$

Note that we *specify* a failure detector \mathcal{D} as a function of the failure pattern F of an execution. However, this does *not* preclude an *implementation* of \mathcal{D} from using other aspects of the execution such as when messages are received. Thus, executions with the

same failure pattern F may still have different failure detector histories. It is for this reason that we allow $\mathcal{D}(F)$ to be a *set* of failure detector histories from which the actual failure detector history for a particular execution is selected non-deterministically.

2.3 Algorithms

We model the asynchronous communication channels as a *message buffer* which contains messages of the form $(p, data, q)$ indicating that process p has sent $data$ addressed to process q and q has not yet received that message. An *algorithm* A is a collection of n (possibly infinite state) deterministic automata, one for each of the processes. $A(p)$ denotes the automaton running on process p . Computation proceeds in *steps of the given algorithm* A . In each step of A , process p performs atomically the following three phases:⁴

Receive phase: p receives a single message of the form $(q, data, p)$ from the message buffer, or a “null” message, denoted λ , meaning that no message is received by p during this step.

Failure detector query phase: p queries and receives a value from its failure detector module. We say that p *sees a value* d when the value returned by p ’s failure detector module is d .

Send phase: p changes its state and sends a message to all the processes according to the automaton $A(p)$, based on its state at the beginning of the step, the message received in the receive phase, and the value that p sees in the failure detector query phase.

The message actually received by the process p in the receive phase is chosen *non-deterministically* from amongst the messages in the message buffer addressed to p , and the null message λ . The null message may be received *even if* there are messages in the message buffer that are addressed to p : the fact that m is in the message buffer merely indicates that m was sent to p . Since ours will be a model of asynchronous systems, where messages may experience arbitrary (but finite) delays, the amount of time m may remain in the message buffer before it is received is unbounded. Indeed, our model will allow a message sent later than another to be received earlier than the other. Though message delays are arbitrary, we also want them to be finite. We model this by introducing a liveness assumption: every message sent will eventually be received, provided its recipient makes “sufficiently many” attempts to receive messages. All this will be made more precise later.

We also remark that the non-determinism arising from the choice of the message to be received reflects the asynchrony of the message buffer — it is not due to non-deterministic choices made by the process. The automaton $A(p)$ is deterministic in the sense that the message that p sends in a step and p ’s new state are uniquely determined from the

⁴Our result also applies to models where steps have finer granularity (see Section 7.1).

present state of p , the message p received during the step and the failure detector value seen by p during the step.

To keep things simple we assume that a process p sends a message m to q at most once. This allows us to speak of the contents of the message buffer as a set, rather than a multiset. We can easily enforce this by adding a counter to each message sent by p to q — so this assumption does not damage generality.

2.4 Configurations, runs and environments

A *configuration* is a pair (s, M) , where s is a function mapping each process p to its local state, and M is a set of triples of the form $(q, data, p)$ representing the messages presently in the message buffer. An *initial configuration of an algorithm A* is a configuration (s, M) , where $s(p)$ is an initial state of $A(p)$ and $M = \emptyset$. A *step* of a given algorithm A transforms one configuration to another. A step of A is uniquely determined by the identity of the process p that takes the step, the message m received by p during that step, and the failure detector value d seen by p during the step. Thus, we identify a step of A with a tuple (p, m, d, A) . If the message received in that step is the null message, then $m = \lambda$, otherwise m is of the form $(-, -, p)$.⁵ We say that a step $e = (p, m, d, A)$ is *applicable to a configuration $C = (s, M)$* if and only if $m \in M \cup \{\lambda\}$. We write $e(C)$ to denote the unique configuration that results when e is applied to C .

A *schedule S of algorithm A* is a finite or infinite sequence of steps of A . S_{\perp} denotes the empty schedule. We say that a schedule S of an algorithm A is *applicable to a configuration C* if and only if (a) $S = S_{\perp}$, or (b) $S[1]$ is applicable to C , $S[2]$ is applicable to $S[1](C)$, etc (we denote by $v[i]$ the i th element of a sequence v). If S is a finite schedule applicable to C , $S(C)$ denotes the unique configuration that results from applying S to C . Note $S_{\perp}(C) = C$ for all configurations C . Let C be any configuration and S be any schedule applicable to C . We say that C' is a *configuration of the pair (S, C)* if there is a finite prefix S' of S such that $C' = S'(C)$.

A *partial run of algorithm A using a failure detector \mathcal{D}* is a tuple $R = \langle F, H_{\mathcal{D}}, I, S, T \rangle$ where F is a failure pattern, $H_{\mathcal{D}} \in \mathcal{D}(F)$ is a failure detector history, I is an initial configuration of A , S is a *finite* schedule of A , and T is a *finite* list of increasing time values (indicating when each step in S occurred) such that $|S| = |T|$, S is applicable to I , and for all $i \leq |S|$, if $S[i]$ is of the form (p, m, d, A) then:

- p has not crashed by time $T[i]$, i.e., $p \notin F(T[i])$
- d is the value of the failure detector module of p at time $T[i]$, i.e., $d = H_{\mathcal{D}}(p, T[i])$

Informally, a partial run of A using \mathcal{D} represents a point of an execution of A using \mathcal{D} .

A *run of an algorithm A using a failure detector \mathcal{D}* is a tuple $R = \langle F, H_{\mathcal{D}}, I, S, T \rangle$ where F is a failure pattern, $H_{\mathcal{D}} \in \mathcal{D}(F)$ is a failure detector history, I is an initial configuration of A , S is an *infinite* schedule of A , and T is an *infinite* list of increasing time values indicating when each step in S occurred. In addition to satisfying the above properties of a partial run, a run must also satisfy the following properties:

⁵Throughout this paper, a “—” in a tuple denotes an arbitrary value of the appropriate type.

- Every correct process takes an infinite number of steps in S . Formally:

$$\forall p \in \text{correct}(F), \forall i, \exists j > i : S[j] \text{ is of the form } (p, -, -, A)$$

- Every message sent to a correct process is eventually received. Formally:

$$\forall p \in \text{correct}(F), \forall C = (s, M) \text{ of } (S, I) : m = (q, \text{data}, p) \in M \Rightarrow (\exists i : S[i] \text{ is of the form } (p, m, -, A))$$

In [CT91], it is shown that any algorithm that uses \mathcal{W} to solve Consensus requires $n > 2f$. With other failure detectors the requirements may be different. For example, there is a failure detector that can be used to solve Consensus only if p_1 and p_2 do not both crash. In general whether a given failure detector can be used to solve Consensus depends upon assumptions about the underlying “environment”. Formally, an *environment* \mathcal{E} (of an asynchronous system) is set of possible failure patterns.⁶

3 The Consensus problem

In the Consensus problem, each process p has an initial value, 0 or 1, and must reach an irrevocable decision on one of these values. Thus, the algorithm of process p , $A(p)$, has two distinct *initial states* σ_0^p and σ_1^p signifying that p 's initial value is 0 or 1. $A(p)$ also has two disjoint sets of *decision states* Σ_0^p and Σ_1^p .

We say that algorithm A uses failure detector \mathcal{D} to solve Consensus in environment \mathcal{E} if every run $R = \langle F, H_{\mathcal{D}}, I, S, T \rangle$ of A using \mathcal{D} where $F \in \mathcal{E}$ satisfies:

Termination: Every correct process eventually decides some value. Formally:

$$\forall p \in \text{correct}(F), \exists C = (s, M) \text{ of } (S, I) : s(p) \in \Sigma_0^p \cup \Sigma_1^p$$

Irrevocability: Once a correct process decides a value, it remains decided on that value.

Formally, let $S[1..i]$ be the prefix of S consisting of the first i elements of S :

$$\forall p \in \text{correct}(F), \forall k \in \{0, 1\}, \forall i \leq i' : (S[1..i](I) = (s, M) \wedge S[1..i'](I) = (s', M') \wedge s(p) \in \Sigma_k^p) \Rightarrow s'(p) \in \Sigma_k^p$$

Agreement: No two correct processes decide differently.

Formally:

$$\forall p, p' \in \text{correct}(F), \forall C = (s, M) \text{ of } (S, I), \forall k, k' \in \{0, 1\} : (s(p) \in \Sigma_k^p \wedge s(p') \in \Sigma_{k'}^{p'}) \Rightarrow k = k'$$

⁶In a synchronous system, assumptions about the underlying environment may also include other characteristics such as the relative process speeds, the maximum message delay, the degree of clock synchronization, etc. In such a system, a more elaborate definition of an environment would be required.

Validity: If a correct process decides v , then v was proposed by some process.

Formally, let $I = (s_0, M_0)$:

$$\forall p \in \text{correct}(F), \forall k \in \{0, 1\} : (\exists C = (s, M) \text{ of } (S, I) : \\ s(p) \in \Sigma_k^p \Rightarrow (\exists q \in \Pi : s_0(q) = \sigma_k^q))$$

4 Reducibility

We now define what it means for an algorithm $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ to transform a failure detector \mathcal{D} into another failure detector \mathcal{D}' in an environment \mathcal{E} . Algorithm $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ uses \mathcal{D} to maintain a variable $output_p$ at every process p . This variable, reflected in the local state of p , emulates the output of \mathcal{D}' at p . Let O_R be the history of all the $output$ variables in run R , i.e., $O_R(p, t)$ is the value of $output_p$ at time t in run R . Algorithm $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ transforms \mathcal{D} into \mathcal{D}' in \mathcal{E} if and only if for every run $R = \langle F, H_{\mathcal{D}}, I, S, T \rangle$ of $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ using \mathcal{D} , where $F \in \mathcal{E}$, $O_R \in \mathcal{D}'(F)$.

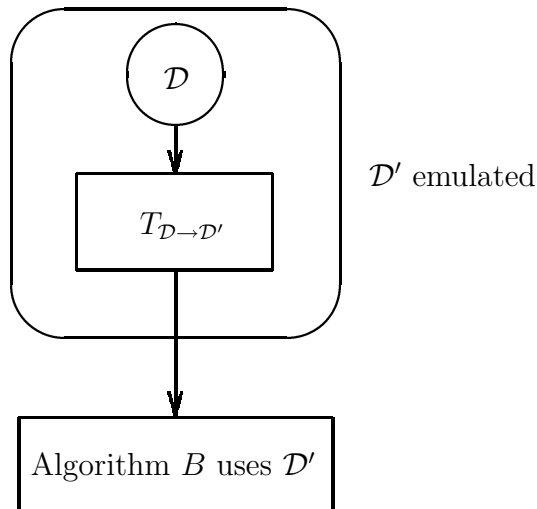


Figure 1: Transforming \mathcal{D} into \mathcal{D}'

Given the reduction algorithm $T_{\mathcal{D} \rightarrow \mathcal{D}'}$, anything that can be done using failure detector \mathcal{D}' in environment \mathcal{E} , can be done using \mathcal{D} instead. To see this, suppose a given algorithm B requires failure detector \mathcal{D}' (when it executes in \mathcal{E}), but only \mathcal{D} is available. We can still execute B as follows. Concurrently with B , processes run $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ to transform \mathcal{D} into \mathcal{D}' . We modify Algorithm B at process p as follows: whenever p is required to query its failure detector module, p reads the current value of $output_p$ (which is concurrently maintained by $T_{\mathcal{D} \rightarrow \mathcal{D}'}$) instead. This is illustrated in Figure 1.

Intuitively, since $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ is able to use \mathcal{D} to emulate \mathcal{D}' , \mathcal{D} provides at least as much information about process failures in \mathcal{E} as \mathcal{D}' does. Thus, if there is an algorithm $T_{\mathcal{D} \rightarrow \mathcal{D}'}$

that transforms \mathcal{D} into \mathcal{D}' in \mathcal{E} , we write $\mathcal{D} \succeq_{\mathcal{E}} \mathcal{D}'$ and say that \mathcal{D}' is *reducible to \mathcal{D} in \mathcal{E}* ; we also say that \mathcal{D}' is *weaker than \mathcal{D} in \mathcal{E}* . Clearly, the reducibility relation $\succeq_{\mathcal{E}}$ is transitive.

Note that, in general, $T_{\mathcal{D} \rightarrow \mathcal{D}'}$ need not emulate *all* the failure detector histories of \mathcal{D}' (in environment \mathcal{E}); what we do require is that all the failure detector histories it emulates be histories of \mathcal{D}' (in that environment).

5 An outline of the result

In [CT91], it is shown that \mathcal{W} can be used to solve Consensus in any environment in which $n > 2f$. We now show that \mathcal{W} is weaker than any failure detector that can be used to solve Consensus. This result holds for any environment \mathcal{E} . Together with [CT91], this implies that \mathcal{W} is indeed the weakest failure detector that can be used to solve Consensus in any environment in which $n > 2f$.

To prove our result, we first define a new failure detector, denoted Ω , that is at least as strong as \mathcal{W} . We then show that any failure detector \mathcal{D} that can be used to solve Consensus is at least as strong as Ω . Thus, \mathcal{D} is at least as strong as \mathcal{W} .

The output of the failure detector module of Ω at a process p is a *single* process, q , that p currently considers to be *correct*; we say that p *trusts* q . In this case, $\mathcal{R}_{\Omega} = \Pi$. For each failure pattern F , $\Omega(F)$ is the set of all failure detector histories H_{Ω} with range \mathcal{R}_{Ω} that satisfy the following property:

- There is a time after which all the correct processes always trust the same correct process:

$$\exists t \in \mathcal{T}, \exists q \in \text{correct}(F), \forall p \in \text{correct}(F), \forall t' \geq t : H_{\Omega}(p, t') = q$$

As with \mathcal{W} , the output of the failure detector module of Ω at a process p may change with time, i.e., p may trust different processes at different times. Furthermore, at any given time t , processes p and q may trust different processes.

Theorem 1: For all environments \mathcal{E} , $\Omega \succeq_{\mathcal{E}} \mathcal{W}$.

PROOF: [Sketch] The reduction algorithm $T_{\Omega \rightarrow \mathcal{W}}$ that transforms Ω into \mathcal{W} is as follows. Each process p periodically sets $\text{output}_p \leftarrow \Pi - \{q\}$, where q is the process that p currently trusts according to Ω . It is easy to see that (in any environment \mathcal{E}) this output satisfies the two properties of \mathcal{W} . \square

Theorem 2: For all environments \mathcal{E} , if a failure detector \mathcal{D} can be used to solve Consensus in \mathcal{E} , then $\mathcal{D} \succeq_{\mathcal{E}} \Omega$.

PROOF: The reduction algorithm $T_{\mathcal{D} \rightarrow \Omega}$ is shown in Section 6. It is the core of our result. \square

Corollary 3: For all environments \mathcal{E} , if a failure detector \mathcal{D} can be used to solve Consensus in \mathcal{E} , then $\mathcal{D} \succeq_{\mathcal{E}} \mathcal{W}$.

PROOF: If \mathcal{D} can be used to solve Consensus in \mathcal{E} , then, by Theorem 2, $\mathcal{D} \succeq_{\mathcal{E}} \Omega$. From Theorem 1, $\Omega \succeq_{\mathcal{E}} \mathcal{W}$. By transitivity, $\mathcal{D} \succeq_{\mathcal{E}} \mathcal{W}$. \square

In [CT91] it is shown that, for all environments \mathcal{E} in which $n > 2f$, \mathcal{W} can be used to solve Consensus. Together with Corollary 3, this shows that:

Theorem 4: For all environments \mathcal{E} in which $n > 2f$, \mathcal{W} is the weakest failure detector that can be used to solve Consensus in \mathcal{E} .

6 The reduction algorithm

6.1 Overview

Let \mathcal{E} be an environment, \mathcal{D} be a failure detector that can be used to solve Consensus in \mathcal{E} , and $\text{Consensus}_{\mathcal{D}}$ be the Consensus algorithm that uses \mathcal{D} . We describe an algorithm $T_{\mathcal{D} \rightarrow \Omega}$ that transforms \mathcal{D} into Ω in \mathcal{E} . Intuitively, this algorithm works as follows.

Fix an arbitrary run of $T_{\mathcal{D} \rightarrow \Omega}$ using \mathcal{D} , with failure pattern $F \in \mathcal{E}$ and failure detector history $H_{\mathcal{D}} \in \mathcal{D}(F)$. Processes periodically query their failure detector \mathcal{D} and exchange information about the values of $H_{\mathcal{D}}$ that they see in this run. Using this information, processes construct a directed acyclic graph (DAG) that represents a “sampling” of failure detector values in $H_{\mathcal{D}}$ and some temporal relationships between the values sampled.

To illustrate this, suppose that process q_1 queries its failure detector \mathcal{D} for the k_1 -th time and sees value d_1 ; q_1 then sends to all processes the message $[q_1, d_1, k_1]$. When a process q_2 receives $[q_1, d_1, k_1]$ it can add vertex $[q_1, d_1, k_1]$ to its (current) version of the DAG: This vertex indicates that q_1 saw d_1 in its k_1 -th failure detector query. When q_2 later queries \mathcal{D} and sees the value d_2 (say this is its k_2 -th query), it adds vertex $[q_2, d_2, k_2]$, and edge $[q_1, d_1, k_1] \rightarrow [q_2, d_2, k_2]$, to its DAG: This edge indicates that q_1 saw d_1 (in its k_1 -th query) *before* q_2 saw d_2 (in its k_2 -th query). By periodically sending its current version of the DAG to all processes, and incorporating all the DAGs that it receives into its own DAG, every correct process can construct ever increasing finite approximations of the same (infinite) limit DAG G .

It turns out that DAG G can be used to simulate runs of $\text{Consensus}_{\mathcal{D}}$ with failure pattern F and failure detector history $H_{\mathcal{D}}$. These are runs that *could have occurred* if processes were running $\text{Consensus}_{\mathcal{D}}$ instead of $T_{\mathcal{D} \rightarrow \Omega}$.

To illustrate this, consider a path of G , say $g = [q_1, d_1, k_1], [q_2, d_2, k_2], [q_3, d_3, k_3], \dots$. We can use this path to simulate schedules of runs of $\text{Consensus}_{\mathcal{D}}$ in which q_1 takes the first step and sees failure detector value d_1 , q_2 takes the second step and sees d_2 , q_3 takes the third step and sees d_3 , etc. We say that such a schedule is “compatible” with path g . Note that there are many schedules of $\text{Consensus}_{\mathcal{D}}$ that are compatible with g : For each step, we have a choice of which message to receive — either one of the messages contained in the simulated buffer (i.e., a message previously sent but not yet received) or the empty message.

Now consider any initial configuration I of $\text{Consensus}_{\mathcal{D}}$. The set of simulated schedules of $\text{Consensus}_{\mathcal{D}}$ that are compatible with some path of G and are also applicable to

I can be organized as a tree: paths in this tree represent simulated runs of $\text{Consensus}_{\mathcal{D}}$ with initial configuration I , and branching occurs at the points where simulated runs diverge. By considering several initial configurations of $\text{Consensus}_{\mathcal{D}}$, we obtain a *forest* of simulated runs of $\text{Consensus}_{\mathcal{D}}$: a tree for each different initial configuration.

Thus, the (infinite) DAG G induces an (infinite) simulation forest Υ of runs of $\text{Consensus}_{\mathcal{D}}$ with failure pattern F and failure detector history $H_{\mathcal{D}}$. Using Υ , we show that it is possible to extract the identity of a process p^* that is correct in F , and we give the extraction algorithm.

The simulation forest Υ , however, is infinite and cannot be computed by any process. Fortunately, the information needed by the extraction algorithm to identify p^* is present in a “crucial” *finite* subgraph of Υ that processes are able to eventually compute. When running $T_{\mathcal{D} \rightarrow \Omega}$, each process p constructs ever increasing finite approximations of the DAG G . Using these approximations, p also constructs ever increasing finite approximations of Υ that eventually include the crucial subgraph needed to extract p^* . At all times, p runs the extraction algorithm on its present finite approximation of Υ to select some process that it considers to be correct: once p ’s approximation of Υ includes the crucial subgraph, the extraction algorithm will select p^* (forever). Thus, there is a time after which all correct processes trust the same correct process, p^* —which is exactly what Ω requires.

Having given an overall account of how the transformation of \mathcal{D} to Ω works, we now provide a roadmap for the rest of this section. We first define the DAGs G that allow us to induce an infinite simulation forest Υ (Section 6.2), and to extract a correct process from Υ (Sections 6.3–6.5). We then show how processes compute ever increasing approximations of such a G and corresponding Υ (Section 6.6.1). Finally, we show that by periodically extracting a process from their current finite approximation of Υ , all correct processes will eventually keep extracting (forever) the same correct process (Section 6.6.2).

We now state some conventions that simplify the discussion that follows. We say that a process is *correct* (*crashes*) if it is correct (crashes) in F . For the rest of this paper, whenever we refer to a run of $\text{Consensus}_{\mathcal{D}}$, we mean a run of $\text{Consensus}_{\mathcal{D}}$ using \mathcal{D} . Furthermore, we only consider schedules of $\text{Consensus}_{\mathcal{D}}$, and so we write (p, m, d) instead of $(p, m, d, \text{Consensus}_{\mathcal{D}})$ to denote a step.

6.2 A DAG and a forest

Let \mathcal{D} be a failure detector that can be used to solve Consensus in environment \mathcal{E} . Given an arbitrary failure pattern $F \in \mathcal{E}$, and failure detector history $H_{\mathcal{D}} \in \mathcal{D}(F)$, let G be any infinite DAG with the following properties:

1. The vertices of G are of the form $[p, d, k]$ where $p \in \Pi$, $d \in \mathcal{R}_{\mathcal{D}}$ and $k \in \mathbb{N}$. Each vertex is labeled with a time (an element in \mathcal{T}) such that:
 - (a) If vertex $[p, d, k]$ is labeled with t , then $p \notin F(t)$ and $d = H_{\mathcal{D}}(p, t)$ (i.e., at time t , p has not crashed and the value of p ’s failure detector module is d).

- (b) If vertices v_1 and v_2 are labeled with t_1 and t_2 , respectively, and $v_1 \rightarrow v_2$ is an edge of G , then $t_1 < t_2$.
- 2. If $[p, d, k], [p, d', k']$ are vertices of G and $k < k'$, then $[p, d, k] \rightarrow [p, d', k']$ is an edge of G .
- 3. G is transitively closed.
- 4. Let V be any finite subset of vertices of G and p be any correct process. There is a $d \in \mathcal{R}_{\mathcal{D}}$ and a $k \in \mathbf{N}$ such that for every vertex $v \in V$, $v \rightarrow [p, d, k]$ is an edge of G .

Note that G contains only a “sampling” of the failure detector values that occur in $H_{\mathcal{D}}$, and only a subset of the temporal relationships that relate them. In other words, we do not require that G contain all the values that occur in $H_{\mathcal{D}}$, or that it relate (with an edge) all its values according to the time at which they occur in $H_{\mathcal{D}}$. However, Property 4 implies that G contains infinitely many “samplings” of the failure detector module of each *correct* process.

Let $g = [q_1, d_1, k_1], [q_2, d_2, k_2], \dots$ be any (finite or infinite) path of G . A schedule S is *compatible with g* if it has the same length as g , and $S = (q_1, m_1, d_1), (q_2, m_2, d_2), \dots$, for some (possibly null) messages m_1, m_2, \dots . We say that S is *compatible with G* if it is compatible with some path of G .

Let I be any initial configuration of $\text{Consensus}_{\mathcal{D}}$. Consider a schedule S that is compatible with G and applicable to I . Intuitively, S is the schedule of a possible run of $\text{Consensus}_{\mathcal{D}}$ with initial configuration I , failure pattern F , and failure detector history $H_{\mathcal{D}}$.

We can represent all the schedules that are compatible with G and applicable to I as a tree. This is called the *simulation tree Υ_G^I induced by G and I* and is defined as follows. The set of vertices of Υ_G^I is the set of *finite* schedules S that are compatible with G and are applicable to I . The *root* of Υ_G^I is the empty schedule S_{\perp} . There is an edge from vertex S to vertex S' if and only if $S' = S \cdot e$ for a step e ;⁷ this edge is labeled e . With each (finite or infinite) path in Υ_G^I , we associate the unique schedule $S = e_1, e_2, \dots, e_k, \dots$ consisting of the sequence of labels of the edges on that path. Note that if a path starts from the root of Υ_G^I and it is finite, the schedule S associated with it is also the last vertex of that path.

The following two lemmata make precise the connection between paths of Υ_G^I and runs of $\text{Consensus}_{\mathcal{D}}$. The proofs, which follow directly from the definitions, are included in the Appendix.

Lemma 5: Let S be a schedule associated with a *finite* path of Υ_G^I that starts from the root of Υ_G^I . There is a sequence of times T such that $\langle F, H_{\mathcal{D}}, I, S, T \rangle$ is a partial run of $\text{Consensus}_{\mathcal{D}}$.

Lemma 6: Let S be a schedule associated with an *infinite* path of Υ_G^I that starts from the root of Υ_G^I . If in S every correct process takes an infinite number of steps and every

⁷If u, w are sequences and u is finite then $u \cdot w$ denotes the concatenation of the two sequences.

message sent to a correct process is eventually received, there is a sequence of times T such that $\langle F, H_{\mathcal{D}}, I, S, T \rangle$ is a run of $\text{Consensus}_{\mathcal{D}}$.

The following lemmata state some “richness” properties of the simulation trees induced by G (their proofs are in the Appendix).

Lemma 7: For any two initial configurations I and I' , if S is a vertex of Υ_G^I and is applicable to I' then S is also a vertex of $\Upsilon_G^{I'}$.

Lemma 8: Let S be any vertex of Υ_G^I and p be any correct process. Let m be a message in the message buffer of $S(I)$ addressed to p or the null message. For some d , S has a child $S \cdot (p, m, d)$ in Υ_G^I .

Lemma 9: Let S be any vertex of Υ_G^I and p be any process. Let m be a message in the message buffer of $S(I)$ addressed to p or the null message. Let S' be a descendent of S such that, for some d , $S' \cdot (p, m, d)$ is in Υ_G^I . For each vertex S'' on the path from S to S' (inclusive), $S'' \cdot (p, m, d)$ is also in Υ_G^I .

Lemma 10: Let S, S_0 , and S_1 be any vertices of Υ_G^I . There is a finite schedule E containing only steps of correct processes such that:

1. $S \cdot E$ is a vertex of Υ_G^I and all correct processes have decided in $S \cdot E(I)$.
2. For $i = 0, 1$, if E is applicable to $S_i(I)$ then $S_i \cdot E$ is a vertex of Υ_G^I .

Let I^i , $0 \leq i \leq n$, denote the initial configuration of $\text{Consensus}_{\mathcal{D}}$ in which the initial values of $p_1 \dots p_i$ are 1, and the initial values of $p_{i+1} \dots p_n$ are 0. The *simulation forest induced by G* is the set $\{\Upsilon_G^{I^0}, \Upsilon_G^{I^1}, \dots, \Upsilon_G^{I^n}\}$ of simulation trees induced by G and initial configurations I^0, I^1, \dots, I^n .

6.3 Tagging the simulation forest

We assign a set of *tags* to each vertex of every tree in the simulation forest induced by G . Vertex S of tree Υ_G^I gets tag k if and only if it has a descendent S' (possibly $S' = S$) such that some *correct* process has decided k in $S'(I)$. Hereafter, Υ^i denotes the tagged tree $\Upsilon_G^{I^i}$, and Υ denotes the tagged simulation forest $\{\Upsilon^0, \Upsilon^1, \dots, \Upsilon^n\}$.

Lemma 11: Every vertex of Υ^i has at least one tag.

PROOF: From Lemma 10, every vertex S of Υ^i has a descendent $S' = S \cdot E$ (for some E) such that all correct processes have decided in $S'(I^i)$. \square

A vertex of Υ^i is *monovalent* if it has only one tag, and *bivalent* if it has both tags, 0 and 1. A vertex is *0-valent* if it is monovalent and is tagged 0; *1-valent* is similarly defined.

Lemma 12: Every vertex of Υ^i is either 0-valent, 1-valent, or bivalent.

PROOF: Immediate from Lemma 11. \square

Lemma 13: The ancestors of a bivalent vertex are bivalent. The descendants of a k -valent vertex are k -valent.

PROOF: Immediate from the definitions. \square

Lemma 14: If vertex S of Υ^i has tag k , then no correct process has decided $1 - k$ in $S(I^i)$.

PROOF: Since S has tag k , it has a descendent S' such that a correct process p has decided k in $S'(I^i)$. From Lemma 5, there is a T such that $R = \langle F, H_{\mathcal{D}}, I^i, S', T \rangle$ is a partial run of $Consensus_{\mathcal{D}}$. Since p has decided k in $S'(I^i)$, from the agreement requirement of Consensus, no correct process has decided $1 - k$ in $S'(I^i)$. Since S' is a descendent of S , from the irrevocability requirement of Consensus, no correct process could have decided $1 - k$ in $S(I^i)$. \square

Lemma 15: If vertex S of Υ^i is bivalent, then no correct process has decided in $S(I^i)$.

PROOF: Immediate from Lemma 14. \square

Recall that in I^0 all processes have initial value 0, and in I^n they all have initial value 1.

Lemma 16: The root of Υ^0 is 0-valent; the root of Υ^n is 1-valent.

PROOF: We first show that the root of Υ^0 is 0-valent. Suppose, for contradiction, that the root of Υ^0 has tag 1. There must be a vertex S of Υ^0 such that some correct process has decided 1 in $S(I^0)$. From Lemma 5, there is a T such that $R = \langle F, H_{\mathcal{D}}, I^0, S, T \rangle$ is a partial run of $Consensus_{\mathcal{D}}$. R violates the validity requirement of Consensus—a contradiction. Thus the root of Υ^0 cannot have a tag of 1. From Lemma 11, the root of Υ^0 has at least one tag: thus it is 0-valent.

By a symmetric argument, the root of Υ^n is 1-valent. \square

Index i is *critical* if the root of Υ^i is bivalent, or if the root of Υ^{i-1} is 0-valent while the root of Υ^i is 1-valent. In the first case, we say that index i is *bivalent critical*; in the second case, we say that i is *monovalent critical*.

Lemma 17: There is a critical index i , $0 < i \leq n$.

PROOF: Apply Lemmata 12 and 16 to the roots of $\Upsilon^0, \Upsilon^1, \dots, \Upsilon^n$. \square

The critical index i is the key to extracting the identity of a correct process. In fact, if i is monovalent critical, we shall prove that p_i must be correct (Lemma 19). If i is bivalent critical, the correct process will be found by focusing on the tree Υ^i , as explained in the following section.

6.4 Of hooks and forks

We describe two forms of finite subtrees of Υ^i referred to as *decision gadgets of Υ^i* . Each type of decision gadget is rooted at the root S_\perp of Υ^i and has exactly two leaves: one 0-valent and one 1-valent. The least common ancestor of these leaves is called the *pivot*. The pivot is clearly bivalent.

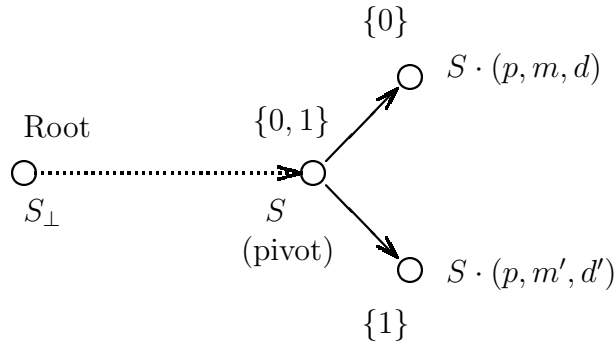


Figure 2: A fork— p is the deciding process

The first type of decision gadget, called a *fork*, is shown in Figure 2. The two leaves are children of the pivot, obtained by applying different steps of the *same* process p . Process p is called the *deciding process of the fork*, because its step after the pivot determines the decision of correct processes.

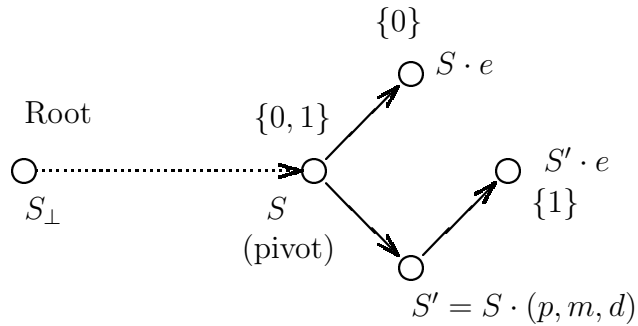


Figure 3: A hook— p is the deciding process

The second type of decision gadget, called a *hook*, is shown in Figure 3. Let S be the pivot of the hook. There is a step e such that $S \cdot e$ is one leaf, and the other leaf is $S \cdot (p, m, d) \cdot e$ for some p, m, d . Process p is called the *deciding process of the hook*, because the decision of correct processes is determined by whether p takes the step (p, m, d) before e . (The tags of the leaves $S \cdot e$ and $S' \cdot e$ may be reversed in a hook relative to the tags shown in Figure 3: it may be that 1 is the only tag of $S \cdot e$ and 0 is

the only tag of $S' \cdot e$.)⁸

We shall prove that the deciding process p of a decision gadget, whether a fork or a hook, must be correct (Lemma 21). Intuitively, this is because if p crashes, then no process can figure out whether p has taken the step that determines the decision value; indeed, this is so even though processes can consult the failure detector \mathcal{D} . Thus, if p crashes, then no process can decide — contradicting the correctness of $\text{Consensus}_{\mathcal{D}}$.

$S \leftarrow S_{\perp}$ $\{S_{\perp} \text{ is the bivalent root of } \Upsilon^i\}$
repeat forever
 Let p be the next correct process in round-robin order
 Let m be the oldest message addressed to p in the message buffer of $S(I^i)$
 (if no such message exists, $m = \lambda$)
 if S has a descendent S' (possibly $S' = S$) such that
 for some d , $S' \cdot (p, m, d)$ is a bivalent vertex in Υ^i
 then $S \leftarrow S' \cdot (p, m, d)$ $\{S \text{ is bivalent}\}$
 else exit

Figure 4: Generating path π in Υ^i

Lemma 18: If index i is bivalent critical then Υ^i has at least one decision gadget (and hence a deciding process).

PROOF: Starting from the bivalent root of Υ^i , we generate a path π in Υ^i , all the vertices of which are bivalent, as follows. We consider all correct processes in round-robin fashion. Suppose we have generated path S so far, and it is the turn of process p . Let m be the oldest message addressed to p that is in the message buffer of $S(I^i)$.⁹ (If no such message exists, we take m to be the null message.) We try to extend the path S so that the last edge in the extension corresponds to p receiving m and the target of that edge is a bivalent vertex. The path construction ends if and when such an extension is no longer possible. This construction is shown in Figure 4. Each iteration of the loop extends the path by at least one edge. Let π be the path generated by these iterations; π is finite or infinite depending on whether the loop terminates.

Claim 1: π is finite.

PROOF: Suppose, for contradiction, that π is infinite. Let S be the schedule associated with π . By construction, in S every correct process takes an infinite number of steps and every message sent to a correct process is eventually received. By Lemma 6, there is a T such that $R = \langle F, H_{\mathcal{D}}, I^i, S, T \rangle$ is a run of $\text{Consensus}_{\mathcal{D}}$. By construction, all vertices in π are bivalent. By Lemma 15, no correct process decides in R , thus violating the termination requirement of Consensus—a contradiction. \square **claim 1**

⁸A fork may be a subgraph of a hook.

⁹By a slight abuse of notation we identify a finite path from the root of Υ^i and its associated schedule.

Let S be the last vertex of π (clearly, S is bivalent). Let p be the next correct process in round-robin order when the loop in Figure 4 terminates. Let m be the oldest message addressed to p in the message buffer of $S(I^i)$ (if no such message exists, m is the null message). The loop exit condition is:

$$\begin{aligned} &\text{For all descendents } S' \text{ of } S \text{ (including } S' = S) \text{ and all } d, \\ &S' \cdot (p, m, d) \text{ is not a bivalent vertex of } \Upsilon^i. \end{aligned} \quad (*)$$

From Lemma 8, for some d , S has a child $S \cdot (p, m, d)$ in Υ^i . By $(*)$ and Lemma 12, $S \cdot (p, m, d)$ is monovalent. Without loss of generality, assume it is 0-valent.

Claim 2: For some d' there is a descendent S' of S such that $S' \cdot (p, m, d')$ is a 1-valent vertex of Υ^i , and the path from S to S' contains no edge labeled $(p, m, -)$.

PROOF: Since S is bivalent, it has a descendent S^* such that some correct process has decided 1 in $S^*(I^i)$. From Lemmata 11 and 14, S^* is 1-valent. There are two cases:

1. The path from S to S^* does not have an edge labeled $(p, m, -)$. Suppose $m \neq \lambda$. Since m is in the message buffer of $S(I^i)$ and p does not receive m in the path from S to S^* , m is still in the message buffer of $S^*(I^i)$. From Lemma 8 (which also applies if $m = \lambda$), for some d' , $S^* \cdot (p, m, d')$ is in Υ^i . Since S^* is 1-valent, by Lemma 13, $S^* \cdot (p, m, d')$ is also 1-valent. In this case, the required S' is S^* .
2. The path from S to S^* has an edge labeled $(p, m, -)$. Let (p, m, d') be the first such edge on that path. Let S' be the source of this edge. By $(*)$ and Lemma 12, $S' \cdot (p, m, d')$ is monovalent. Since $S' \cdot (p, m, d')$ has a 1-valent descendent S^* , by Lemma 13, $S' \cdot (p, m, d')$ is 1-valent. □ **claim 2**

Consider the vertex S' and edge (p, m, d') of Claim 2. By Lemma 9, for each vertex S'' on the path from S to S' (inclusive), $S'' \cdot (p, m, d')$ is also in Υ^i . By $(*)$ and Lemma 12, all such vertices $S'' \cdot (p, m, d')$ are monovalent. In particular, $S \cdot (p, m, d')$ is monovalent. There are two cases (see Figure 5):

1. $S \cdot (p, m, d')$ is 1-valent. Since $S \cdot (p, m, d)$ is 0-valent, Υ^i has a fork with pivot S .
2. $S \cdot (p, m, d')$ is 0-valent. Recall that $S' \cdot (p, m, d')$ is 1-valent and for each vertex S'' between S and S' , $S'' \cdot (p, m, d')$ is monovalent. Thus, the path from S to S' must have two vertices S_0 and S_1 such that S_0 is the parent of S_1 , $S_0 \cdot (p, m, d')$ is 0-valent and $S_1 \cdot (p, m, d')$ is 1-valent. Hence, Υ^i has a hook with pivot S_0 . □

6.5 Extracting the correct process

By Lemma 17, there is a critical index i . If i is monovalent critical, Lemma 19 below shows how to extract a correct process. If i is bivalent critical, a correct process can be found by applying Lemmata 18 and 21.

Lemma 19: If index i is monovalent critical then p_i is correct.

PROOF: Suppose, for contradiction, that p_i crashes. By Lemma 10(1) (applied to the root $S = S_\perp$ of Υ^i), there is a finite schedule E that contains only steps of correct

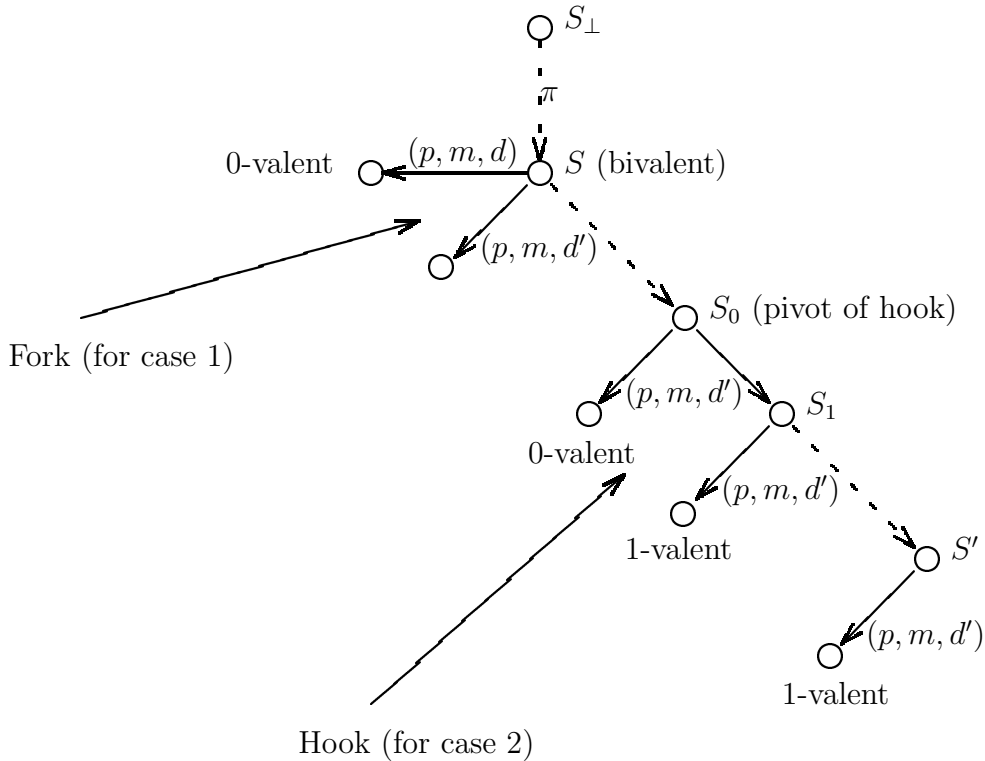


Figure 5: The decision gadgets in Υ^i if i is bivalent critical

processes (and hence no step of p_i) such that all correct processes have decided in $E(I^i)$. Since index i is monovalent critical, the root S_\perp of Υ^i is 1-valent. Hence all correct processes must have decided 1 in $E(I^i)$.

I^i and I^{i-1} only differ in the state of p_i . Since E is applicable to I^i and does not contain any steps of p_i , an easy induction on the number of steps in E shows that: (a) E is also applicable to I^{i-1} , and (b) the state of all processes other than p_i are the same in $E(I^i)$ and $E(I^{i-1})$. Using Lemma 7, (a) implies that E is also a vertex of Υ^{i-1} . By (b), all correct processes have decided 1 in $E(I^{i-1})$. Thus the root of Υ^{i-1} has tag 1. Since i is monovalent critical, the root of Υ^{i-1} is 0-valent—a contradiction. \square

Lemma 20: Let S be any bivalent vertex of Υ^i , and S_0, S_1 be any 0-valent and 1-valent descendents of S . If there is a process p such that the paths from S to S_0 and from S to S_1 contain only steps of the form $(p, -, -)$, then p is correct.

PROOF: Suppose, for contradiction, that p crashes. From Lemma 10, there is a schedule E containing only steps of correct processes (and hence no step of p) such that:

- i. $S \cdot E$ is a vertex of Υ^i and all correct processes have decided in $S \cdot E(I^i)$.
- ii. For $k = 0, 1$, if $S_k \cdot E$ is applicable to I^i then $S_k \cdot E$ is a vertex of Υ^i .

Without loss of generality assume that all correct processes decided 0 in $S \cdot E(I^i)$.

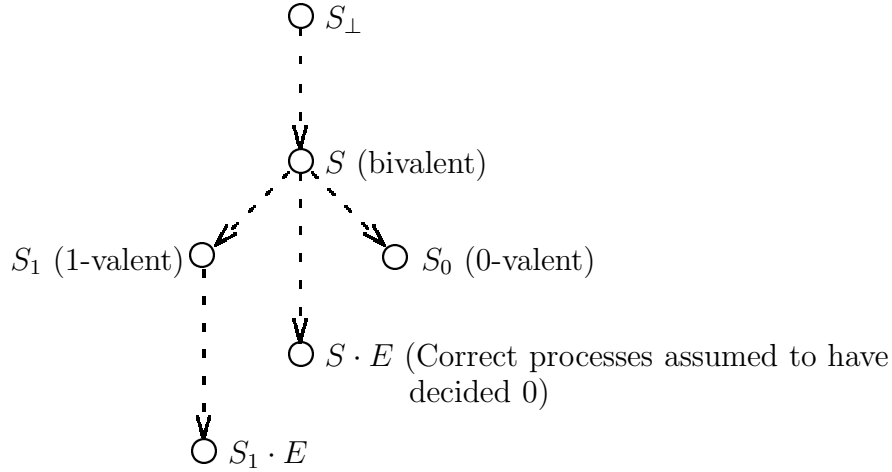


Figure 6: Lemma 20

Refer to Figure 6. Since all steps in the path from S to S_1 are steps of p , the state of every process other than p is the same in $S(I^i)$ and in $S_1(I^i)$. Furthermore, any message addressed to a process other than p that is in the message buffer in $S(I^i)$ is still in the message buffer in $S_1(I^i)$. Since E is applicable to $S(I^i)$ and does not contain any steps of p , an easy induction on the number of steps in E shows that: (a) E is also applicable to $S_1(I^i)$, and (b) the state of every process other than p is the same in $S \cdot E(I^i)$ and $S_1 \cdot E(I^i)$. By (ii), (a) implies that $S_1 \cdot E(I^i)$ is a vertex in Υ^i . By (b), all correct processes decide 0 in $S_1 \cdot E(I^i)$. So S_1 , has tag 0. But S_1 is 1-valent—a contradiction. \square

Lemma 21: The deciding process of a decision gadget is correct.

PROOF: Let γ be any decision gadget of Υ^i . There are two cases to consider:

1. γ is a fork. By Lemma 20, the deciding process of γ is correct.
2. γ is a hook. Assume (without loss of generality) that S is the pivot of γ , $S_0 = S \cdot (p', m', d')$ is the 0-valent leaf of γ and $S_1 = S \cdot (p, m, d) \cdot (p', m', d')$ is the 1-valent leaf of γ (see Figure 7). There are two cases:
 - (a) $p = p'$. By Lemma 20, p is correct.
 - (b) $p \neq p'$. Suppose, for contradiction, that p crashes. By Lemma 10, there is a schedule E containing only steps of correct processes (and hence no step of p) such that:
 - i. $S_0 \cdot E$ is a vertex of Υ^i and all correct processes have decided in $S_0 \cdot E(I^i)$. Since S_0 is 0-valent, all correct processes must have decided 0 in $S_0 \cdot E(I^i)$.
 - ii. If E is applicable to $S_1(I^i)$ then $S_1 \cdot E$ is a vertex of Υ^i .

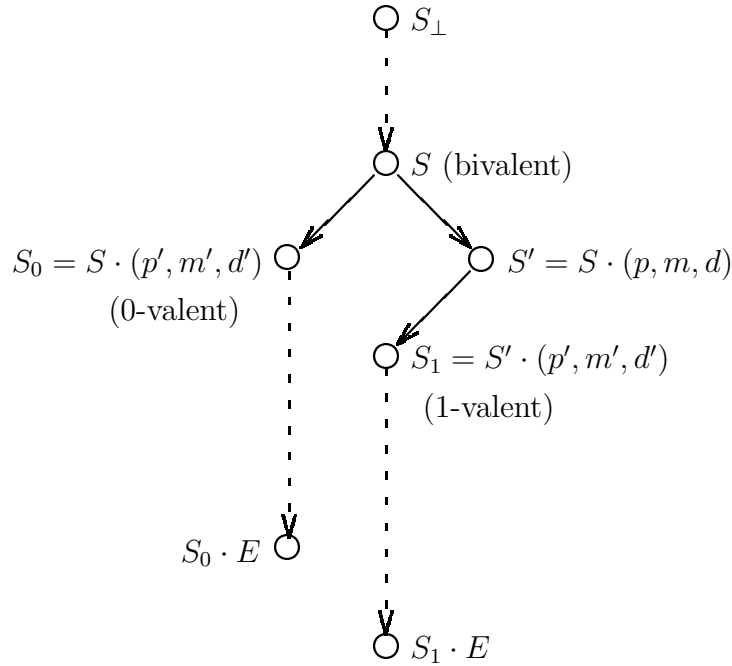


Figure 7: Lemma 21

Let $S' = S \cdot (p, m, d)$ be the parent of S_1 . The state of every process other than p is the same in $S(I^i)$ and $S'(I^i)$. Furthermore, any message addressed to a process other than p that is in the message buffer in $S(I^i)$ is still in the message buffer in $S'(I^i)$. Therefore, since $S_0 = S \cdot (p', m', d')$ and $S_1 = S' \cdot (p', m', d')$, the state of every process other than p is the same in $S_0(I^i)$ and $S_1(I^i)$. In addition, any message addressed to a process other than p that is in the message buffer in $S_0(I^i)$ is also in the message buffer in $S_1(I^i)$. Since E is applicable to $S_0(I^i)$ and does not contain any steps of p , an easy induction on the number of steps in E shows that: (a) E is also applicable to $S_1(I^i)$, and (b) the state of every process other than p is the same in $S_0 \cdot E(I^i)$ and $S_1 \cdot E(I^i)$. By (ii), (a) implies that $S_1 \cdot E$ is a vertex of Υ^i . By (b), all correct processes decide 0 in $S_1 \cdot E(I^i)$. Thus S_1 , receives a tag of 0. But S_1 is 1-valent—a contradiction. \square

There may be several critical indices and several decision gadgets in the simulation forest. Thus, Lemmata 19 and 21 may identify many correct processes. Our selection rule will choose *one* of these, as the failure detector Ω requires, as follows. It first determines the smallest critical index i . If i is monovalent critical, it selects p_i . If, on the other hand, i is bivalent critical, it chooses the “smallest” decision gadget in Υ^i according to some encoding of gadgets, and selects the corresponding deciding process. It is easy to encode finite graphs as natural numbers. Since a decision gadget is just a finite graph, the selection rule can use any such encoding. The whole method of selecting a correct

process is shown in Figure 8 (recall that G is any directed acyclic graph that satisfies Properties 1–4 of Section 6.2 with respect to the given failure pattern F).

```

{Build and tag simulation forest  $\Upsilon$  induced by  $G$ }
for  $i \leftarrow 0, 1, \dots, n$ :
     $\Upsilon^i \leftarrow$  simulation tree induced by  $G$  and  $I^i$ 
    for every vertex  $S$  of  $\Upsilon^i$ 
        if  $S$  has a descendent  $S'$  such that a correct process has decided  $k$  in  $S'(I^i)$ 
            then add tag  $k$  to  $S$ 

{Select a process from tagged simulation forest  $\Upsilon$ }
 $i \leftarrow$  smallest critical index (1)
if  $i$  is monovalent critical then return  $p_i$  (2)
else return deciding process of the smallest decision gadget in  $\Upsilon^i$  (3)

```

Figure 8: Selecting a correct process

Theorem 22: The algorithm in Figure 8 selects a correct process.

PROOF: By Lemma 17, there is a critical index $i, 0 < i \leq n$. If i is monovalent critical, Line (2) returns p_i which, by Lemma 19, is correct. If i is bivalent critical, by Lemma 18, Υ^i contains at least one decision gadget. Let γ be the decision gadget in Υ^i with the smallest encoding. By Lemma 21, the deciding process of γ is correct in F . Thus, Line (3) returns the identity of a process that is correct. \square

6.6 The reduction algorithm $T_{\mathcal{D} \rightarrow \Omega}$

The selection of a correct process described in Figure 8 is not yet the distributed algorithm $T_{\mathcal{D} \rightarrow \Omega}$ that we are seeking: it involves an infinite simulation forest and is “centralized”. To turn it into a distributed algorithm, we will modify it as follows. Each process will cooperate with other processes to construct ever increasing finite approximations of the same simulation forest. Such approximations will eventually contain the decision gadget and the other tagging information necessary to extract the identity of the *same* correct process chosen by the selection method in Figure 8.

Note that the selection method in Figure 8 involves three stages: The construction of G , a DAG representing samples of failure detector values and some of their temporal relationship; the construction and tagging of the simulation forest induced by G ; and, finally, the selection of a correct process using this forest.

Algorithm $T_{\mathcal{D} \rightarrow \Omega}$ consists of two components. In the first component, each process repeatedly queries its failure detector module and sends the failure detector values it sees to the other processes. This component enables correct processes to construct ever increasing finite approximations of the *same* DAG G . Since all inter-process communication occurs in this component, we call it the *communication component* of $T_{\mathcal{D} \rightarrow \Omega}$.

In the second component, each process repeatedly (a) constructs and tags the simulation forest induced by its current approximation of graph G , and (b) selects the identity of a process using its current simulation forest. Since this component does not require any communication, we call it the *computation component* of $T_{\mathcal{D} \rightarrow \Omega}$.

6.6.1 The communication component

In this component processes cooperate to construct ever increasing approximations of the same graph G . The communication component of $T_{\mathcal{D} \rightarrow \Omega}$ for p , shown in Figure 9, works as follows. Let G_p denote p 's current approximation of G . Every process p repeatedly performs the following three tasks. (i) If p receives G_q for some q , it incorporates this information by replacing G_p with the union of G_p and G_q . (ii) Process p queries its own failure detector module. If this is the k -th query of p to its failure detector module, and d is the value p sees in this query, then p adds $[p, d, k]$ to G_p . Let $[p', d', k']$ be any vertex that was in G_p just before this insertion. Clearly p saw d in its k -th query after p' saw d' in its k' -th query. Thus p adds an edge from every such $[p', d', k']$ in G_p to $[p, d, k]$. (iii) Process p sends its updated G_p to all processes.

Note that the body of the **repeat forever** loop in Figure 9 corresponds to a *single step* of p 's algorithm, and that each step is subdivided in three phases: receive, failure detector query, and send. This conforms to our model, as defined in Section 2. Recall that in this model the three phases of a step occur atomically at some discrete time t . If p takes a step at time t , we denote by $G_p(t)$ the value of G_p *at the end* of that step.

```

{Build the directed acyclic graph  $G_p$ }
 $G_p \leftarrow$  empty graph
 $k_p \leftarrow 0$ 
repeat forever
  RECEIVE PHASE:
     $p$  receives  $m$ 
  FAILURE DETECTOR QUERY PHASE:
     $d_p \leftarrow$  query failure detector  $\mathcal{D}$ 
     $k_p \leftarrow k_p + 1$ 
  SEND PHASE:
    if  $m$  is of the form  $(q, G_q, p)$  then  $G_p \leftarrow G_p \cup G_q$  (1)
    add  $[p, d_p, k_p]$  to  $G_p$  and edges from all other vertices of  $G_p$  to  $[p, d_p, k_p]$  (2)
     $output_p \leftarrow$  computation component {Figure 10} (3)
     $p$  sends  $(p, G_p, q)$  to all  $q \in \Pi$  (4)

```

Figure 9: Process p 's communication component

To show that the local graphs constructed by the communication component are ever increasing finite approximations of the same infinite limit graph, we first prove:

LEMMA 6.6.1.1 *For any correct process p and any time t :*

1. $G_p(t)$ is a subgraph of $G_p(t')$, for all $t' \geq t$.
2. For every correct process q , there is a time $t' \geq t$ such that $G_p(t)$ is a subgraph of $G_q(t')$.

PROOF: The first part of the lemma follows from the fact that no vertex or edge of G_p is ever removed. Since p is correct, at some time $t' \geq t$ it sends its local graph $G_p(t')$ to all processes, including q . Since q is correct, it eventually receives $G_p(t')$, and then replaces G_q with $G_q \cup G_p(t')$, say at time t'' . By the first part of the lemma, $G_p(t)$ is a subgraph of $G_p(t')$. Thus, $G_p(t)$ is a subgraph of $G_q(t'')$. □

Lemma 6.6.1.1(1) allows us to define $G_p^\infty = \bigcup_{t \in \mathcal{T}} G_p(t)$.

LEMMA 6.6.1.2 *For any correct processes p and q , $G_p^\infty = G_q^\infty$.*

PROOF: Let o be any vertex or edge of G_p^∞ , i.e., there is a time t at which o is in $G_p(t)$. From Lemma 6.6.1.1(2), there is a time t' such that $G_p(t)$ is a subgraph of $G_q(t')$. Thus o is in G_q^∞ . Thus G_p^∞ is a subgraph of G_q^∞ . By a symmetric argument, G_q^∞ is a subgraph of G_p^∞ , hence $G_p^\infty = G_q^\infty$. □

Lemma 6.6.1.2 allows us to define the *limit graph* G to be G_p^∞ for any correct process p . We will show (Corollary 6.6.1.5) that G is an infinite DAG that has the four properties defined in Section 6.2. To do so we first prove a technical lemma.

LEMMA 6.6.1.3 *Let $v = [p, d, k]$ be a vertex of a local graph during the execution of the communication component, and t be the earliest time when v appears in any local graph.*

1. The first graph that contains v is $G_p(t)$ (i.e., v is in $G_p(t)$, but not in $G_q(t')$, for any process q and any time $t' < t$). Moreover, at time t , p sees d , k is the value of k_p , and p inserts v in G_p .
2. If edge $u \rightarrow v$ is in some local graph during the execution of the communication component, then $u \rightarrow v$ is also in $G_p(t)$.
3. Any local graph that contains v also contains $G_p(t)$.

PROOF: 1. Let q be the process that inserts v in its local graph G_q at time t . This insertion occurs in Line (1) or (2) of q 's communication component. If it occurs in Line (1), q must have received at time t a message with a graph that contains v . The process that sent that message must have had v in its local graph before time t , contradicting the definition of t . Thus, q inserts v in its local graph G_q by executing Line (2). Since $v = [p, d, k]$, it is clear that $q = p$, and the result immediately follows from p 's algorithm.

2. Let t' be the earliest time when edge $u \rightarrow v$ appears in any local graph, and q be the process that adds $u \rightarrow v$ to its local graph G_q at time t' . By definition of t , $t' \geq t$. If $t' > t$, it must be that at time t' process q receives a message that contains a graph with the edge $u \rightarrow v$. The sender of that message had a local graph that contained the edge $u \rightarrow v$ some time before t' , contradicting the definition of t' . Therefore, $t' = t$. Then, by Part (1), $q = p$, and so $u \rightarrow v$ is in $G_p(t)$, as wanted.

3. Suppose, for contradiction, that some local graph contains v , but does not contain $G_p(t)$. Let t' be the earliest time when such a local graph is formed, and say that this occurs at process q . So $G_q(t')$ contains v but not $G_p(t)$. By definition of t , $t' \geq t$. By Lemma 6.6.1.1(1), $G_p(t')$ contains $G_p(t)$, and so $q \neq p$. Therefore, at time t' process q receives a message with a graph that contains v but not $G_p(t)$. The sender of that message must have had a local graph that contains v but not $G_p(t)$ some time before t' , contradicting the definition of t' . \square

Recall that we are considering a fixed run of $T_{\mathcal{D} \rightarrow \Omega}$, with failure pattern F , and failure detector history $H_{\mathcal{D}} \in \mathcal{D}(F)$. We now prove that the local graphs constructed by the communication component of $T_{\mathcal{D} \rightarrow \Omega}$ satisfy four properties that are very similar to those of the DAGs defined in Section 6.2. To state these properties, we define the *label* of a vertex v of a local graph to be the earliest time when v appears in any local graph.

LEMMA 6.6.1.4 *For any correct process p and any time t :*

1. *The vertices of $G_p(t)$ are of the form $[p', d', k']$ where $p' \in \Pi$, $d' \in \mathcal{R}_{\mathcal{D}}$ and $k' \in \mathbf{N}$. The labels of the vertices of $G_p(t)$ are such that:

 - (a) *If vertex $[p', d', k']$ is labeled with t' , then $p' \notin F(t')$ and $d' = H_{\mathcal{D}}(p', t')$.*
 - (b) *If vertices v_1 and v_2 are labeled with t_1 and t_2 , respectively, and $v_1 \rightarrow v_2$ is an edge of $G_p(t)$, then $t_1 < t_2$.¹⁰**
2. *If $[p', d', k']$ and $[p', d'', k'']$ are vertices of $G_p(t)$, and $k' < k''$, then $[p', d', k'] \rightarrow [p', d'', k'']$ is an edge of $G_p(t)$.*
3. *$G_p(t)$ is transitively closed.*
4. *There is a time $t' \geq t$, a $d \in \mathcal{R}_{\mathcal{D}}$ and a $k \in \mathbf{N}$ such that for every vertex v of $G_p(t)$, $v \rightarrow [p, d, k]$ is an edge of $G_p(t')$.*

PROOF:

Property 1 : From Figure 9, it is clear that all the vertices of $G_p(t)$ have the required form. Consider a vertex $[p', d', k']$ labeled with t' . By Lemma 6.6.1.3(1), p' saw d' at time t' . Thus $p' \notin F(t')$ (otherwise p' would not have taken a step at time t' and would not have seen d'), and $d' = H_{\mathcal{D}}(p', t')$, proving 1a.

¹⁰This immediately implies that $G_p(t)$ is acyclic.

Now consider vertices v_1 and v_2 , labeled t_1 and t_2 , respectively, such that $v_1 \rightarrow v_2$ is an edge in $G_p(t)$. Let $v_2 = [p_2, d_2, k_2]$. By Lemma 6.6.1.3(1), p_2 inserted v_2 in G_{p_2} at time t_2 . By Lemma 6.6.1.3(2), $v_1 \rightarrow v_2$ is an edge in $G_{p_2}(t_2)$. Therefore vertex v_1 was in G_{q_2} before time t_2 . Hence, by definition of t_1 , $t_1 < t_2$, proving 1b.

Property 2 : Consider vertices $[p', d', k']$ and $[p', d'', k'']$ of $G_p(t)$, such that $k' < k''$. Let the labels of these vertices be t' and t'' , respectively. By Lemma 6.6.1.3(1), the first graphs that contain these vertices are $G_{p'}(t')$ and $G_{p'}(t'')$, respectively. Furthermore, the values of $k_{p'}$ at times t' and t'' are k' and k'' , respectively. Since the value of $k_{p'}$ never decreases, and $k' < k''$, it follows that $t' < t''$. By Lemma 6.6.1.1(1), $G_{p'}(t')$ is a subgraph of $G_{p'}(t'')$. Thus, when p' inserts $[p', d'', k'']$ in $G_{p'}$ at time t'' , $G_{p'}$ already contains $[p', d', k']$. In Line (2) of the algorithm, p' adds the edge $[p', d', k'] \rightarrow [p', d'', k'']$ to $G_{p'}(t'')$. By Lemma 6.6.1.3(3), since $G_p(t)$ contains $[p', d'', k'']$, it also contains $G_{p'}(t'')$. Thus, $G_p(t)$ contains the edge $[p', d', k'] \rightarrow [p', d'', k'']$.

Property 3 : Let $[q_1, d_1, k_1] \rightarrow \dots \rightarrow [q_\ell, d_\ell, k_\ell]$ ($\ell > 1$) be a path in $G_p(t)$. We must show that there is an edge $[q_1, d_1, k_1] \rightarrow [q_\ell, d_\ell, k_\ell]$.

For all $1 \leq i \leq \ell$, let t_i be the label of $[q_i, d_i, k_i]$. By Lemma 6.6.1.3(1), t_i is the time when q_i inserted $[q_i, d_i, k_i]$ into G_{q_i} . We first show by induction on i that $[q_1, d_1, k_1] \rightarrow \dots \rightarrow [q_i, d_i, k_i]$ is a path in $G_{q_i}(t_i)$. The basis, $i = 1$, is trivial. For the induction step, suppose that $[q_1, d_1, k_1] \rightarrow \dots \rightarrow [q_{i-1}, d_{i-1}, k_{i-1}]$ is a path in $G_{q_{i-1}}(t_{i-1})$. By Lemma 6.6.1.3(2), since $[q_{i-1}, d_{i-1}, k_{i-1}] \rightarrow [q_i, d_i, k_i]$ is an edge in $G_p(t)$, it is also an edge in $G_{q_i}(t_i)$. By Lemma 6.6.1.3(3), since $G_{q_i}(t_i)$ contains vertex $[q_{i-1}, d_{i-1}, k_{i-1}]$, $G_{q_i}(t_i)$ also contains $G_{q_{i-1}}(t_{i-1})$. In particular, $G_{q_i}(t_i)$ contains the path $[q_1, d_1, k_1] \rightarrow \dots \rightarrow [q_{i-1}, d_{i-1}, k_{i-1}]$. Thus, $[q_1, d_1, k_1] \rightarrow \dots \rightarrow [q_i, d_i, k_i]$ is a path in $G_{q_i}(t_i)$, as wanted.

Therefore, vertices $[q_1, d_1, k_1], \dots, [q_\ell, d_\ell, k_\ell]$ are all in $G_{q_\ell}(t_\ell)$. When q_ℓ inserts $[q_\ell, d_\ell, k_\ell]$ in G_{q_ℓ} (at time t_ℓ), G_{q_ℓ} already contains $[q_1, d_1, k_1]$. Thus, q_ℓ also adds the edge $[q_1, d_1, k_1] \rightarrow [q_\ell, d_\ell, k_\ell]$ in $G_{q_\ell}(t_\ell)$. By Lemma 6.6.1.3(3), since $G_p(t)$ contains $[q_\ell, d_\ell, k_\ell]$, it also contains $G_{q_\ell}(t_\ell)$. Therefore, edge $[q_1, d_1, k_1] \rightarrow [q_\ell, d_\ell, k_\ell]$ is in $G_p(t)$.

Property 4 : Since p is correct, it takes a step at some time t' after t . In the failure detector query phase of this step, p queries its failure detector module and obtains a value, say $d \in \mathcal{R}_D$. In Line (2) of this step, p inserts vertex $[p, d, k]$ in G_p (where k is the current value of k_p) and an edge from all other vertices of $G_p(t')$ to $[p, d, k]$. By Lemma 6.6.1.1(1), $G_p(t)$ is a subgraph of $G_p(t')$, hence the result follows. \square

Lemma 6.6.1.4 and the definition of the limit graph G immediately imply:

COROLLARY 6.6.1.5 *The limit graph G satisfies the four properties of the DAGs defined in Section 6.2.*

As before, Υ^i denotes the tagged simulation tree induced by the limit graph G and initial configuration I^i , and Υ denotes the tagged simulation forest $\{\Upsilon^0, \Upsilon^1, \dots, \Upsilon^n\}$.

```

{Build and tag simulation forest  $\Upsilon_p$  induced by  $G_p$ }
for  $i \leftarrow 0, 1, \dots, n$ :
   $\Upsilon_p^i \leftarrow$  simulation tree induced by  $G_p$  and  $I^i$ 
  for every vertex  $S$  of  $\Upsilon_p^i$ 
    if  $S$  has a descendent  $S'$  such that  $p$  has decided  $k$  in  $S'(I^i)$ 
      then add tag  $k$  to  $S$ 

{Select a process from tagged simulation forest  $\Upsilon_p$ }
if there is no critical index then return  $p$ 
else
   $i \leftarrow$  smallest critical index (1)
  if  $i$  is monovalent critical then return  $p_i$  (2)
  else if  $\Upsilon_p^i$  has no decision gadgets then return  $p$ 
    else return deciding process of the smallest decision gadget in  $\Upsilon_p^i$  (3)

```

Figure 10: Process p 's computation component

6.6.2 The computation component

Since the limit graph G has the four properties of Section 6.2 (Corollary 6.6.1.5), we can apply the “centralized” selection method of Figure 8 to identify a correct process. This method involved:

- Constructing and tagging the infinite simulation forest Υ induced by G .
- Applying a rule to Υ to select a particular correct process p^* .

In the computation component of $T_{\mathcal{D} \rightarrow \Omega}$, each process p approximates the above method by repeatedly:

- Constructing and tagging the *finite* simulation forest Υ_p induced by G_p , its present finite approximation of G .
- Applying the same rule to Υ_p to select a particular process.

Since the limit of Υ_p over time is Υ , and the information necessary to select p^* is in a finite subgraph of Υ , we can show that *eventually* p will keep selecting the correct process p^* , forever.

Actually, p cannot quite use the tagging method of Figure 8: that method requires knowing which processes are correct! Instead, p assigns tag k to a vertex S in Υ_p^i if and only if S has a descendent S' such that p itself has decided k in $S'(I^i)$. If p is correct, this is eventually equivalent to the tagging method of Figure 8. If p crashes, we do not care how it tags its forest. Also, p cannot use exactly the same selection method as that of Figure 8: its current simulation forest Υ_p may not *yet* have a critical index or contain

any decision gadget (although it eventually will!). In that case, p temporizes by just selecting itself. The computation component of $T_{\mathcal{D} \rightarrow \Omega}$ is shown in Figure 10 (compare it with the selection method of Figure 8).

We first show that Υ_p , the simulation forest that p constructs, is indeed an increasingly accurate approximation of Υ (Lemma 23). We then show that the tags that p gives to any vertex S in Υ_p are eventually the same ones that the tagging rule of Figure 8 gives to S in Υ (Lemma 24). Let $\Upsilon_p(t)$ denote Υ_p at time t , i.e., $\Upsilon_p(t)$ is the finite simulation forest induced by $G_p(t)$.

Lemma 23: For any correct p and any time t :

1. $\Upsilon_p(t)$ is a subgraph¹¹ of Υ .
2. $\Upsilon_p(t)$ is a subgraph of $\Upsilon_p(t')$, for all $t' \geq t$.
3. $\bigcup_{t \in \mathcal{T}} \Upsilon_p(t) = \Upsilon$.

PROOF:

Property 1 : Let S be any vertex of tree $\Upsilon_p^i(t)$ (for some i , $0 \leq i \leq n$). From the definition of $\Upsilon_p^i(t)$, S is compatible with some path g of $G_p(t)$ and applicable to I^i . Since $G_p(t)$ is a subgraph of G , g is also a path of G . Thus, S is compatible with G ; since it is also applicable to I^i , it is a vertex of Υ^i .

Similarly, let $S \rightarrow S'$ be an edge e of $\Upsilon_p^i(t)$. Since S and S' are also vertices of Υ^i , and $S' = S \cdot e$, $S \rightarrow S'$ is also an edge of Υ^i .

Property 2 : Follows from Lemma 6.6.1.1(1).

Property 3 : We first show that Υ is a subgraph of $\bigcup_{t \in \mathcal{T}} \Upsilon_p(t)$. Let S be any vertex of any tree Υ^i of Υ . From the definition of Υ^i , S is compatible with some finite path g of G and is applicable to I^i . Since $G = \bigcup_{t \in \mathcal{T}} G_p(t)$ and g is a finite path of G , there is a time t such that g is also a path of $G_p(t)$. Since S is compatible with g of $G_p(t)$ and is applicable to I^i , S is a vertex of $\Upsilon_p^i(t)$.

Let $S \rightarrow S'$ be any edge e of Υ^i . By the argument above, there is a time t after which both S and S' are vertices of Υ_p^i . Since $S' = S \cdot e$, after time t the edge e is also in Υ_p^i . Thus, every vertex and every edge of Υ is also in $\bigcup_{t \in \mathcal{T}} \Upsilon_p(t)$, i.e., Υ is a subgraph of $\bigcup_{t \in \mathcal{T}} \Upsilon_p(t)$. By Property 1, $\bigcup_{t \in \mathcal{T}} \Upsilon_p(t) = \Upsilon$. \square

Lemma 24: Let p be any correct process, and S be any vertex of Υ_p . There is a time after which the tags of S in Υ_p are the same as the tags of S in Υ .

PROOF: Suppose that at some time t , p assigns tag k to vertex S of tree Υ_p^i . Thus S has a descendent S' in $\Upsilon_p^i(t)$ such that p has decided k in $S'(I^i)$. By Lemma 23(1), S' is also a descendent of S in Υ^i , and since p is correct, S has tag k in Υ^i as well.

¹¹The subgraph and graph equality relations ignore the tags.

Conversely, suppose a vertex S of a tree Υ^i of Υ has tag k . We show that, eventually, p also assigns tag k to S in Υ_p^i . Since S has tag k in Υ^i , S has a descendent S' in Υ^i such that some correct process has decided k in $S'(I^i)$ (cf. tagging rule in Figure 8). By Lemma 10(1), there is a descendent S'' of S' in Υ^i , such that *all* correct processes, including p , have decided in $S''(I^i)$. By Lemma 5, $S''(I^i)$ is a configuration of a partial run of Consensus_D . By the Agreement property of Consensus, p must have decided k in $S''(I^i)$. Consider the path that starts from the root of Υ^i and goes to vertex S and then to S'' . By Lemma 23(3), there is a time t after which this path is also in Υ_p^i . Therefore, when p executes the tagging rule of Figure 10 after time t , p assigns tag k to S in Υ_p^i (because p has decided k in $S''(I^i)$, and S'' is a descendent of S in Υ_p^i). \square

Recall that p^* is the correct process obtained by applying the selection rule of Figure 8 to the infinite simulation forest Υ . We now show that there is a time after which any correct p always selects p^* when it applies the corresponding selection rule of Figure 10 to its own finite approximation of the simulation forest Υ_p . Roughly speaking, the reason is as follows. By Lemma 24, there is a time t after which the tags of all the roots in p 's forest Υ_p are the same as in the infinite forest Υ . Since these tags determine the sets of monovalent and bivalent critical indices, after time t these sets according to p are the same as in Υ . Let i be the minimum critical index in these sets, and consider the situation after time t . If i is monovalent critical, the selection rule of Figure 10 selects p_i , which is what p^* is in this case. If i is bivalent critical, then p selects the deciding process of its current minimum decision gadget of Υ_p^i (if it has one). This case is examined below.

Let γ^* be the minimum decision gadget of Υ^i (so, p^* is the deciding process of γ^*). For a while, γ^* may not be the minimum decision gadget of Υ_p^i . This may be because γ^* (and its tags) is not yet in Υ_p^i . However, by Lemmata 23(3) and 24, γ^* (including its tags) will eventually be in Υ_p^i . Alternatively, it may be because Υ_p^i contains a subgraph γ whose encoding is smaller than γ^* 's, and for a while γ looks like a decision gadget according to its *present* tags. However, by Lemma 24, p will eventually determine *all* the tags of γ , and discover that γ is not really a decision gadget. Since there are only *finitely* many graphs whose encoding is smaller than γ^* 's, p will eventually discard all the “fake” decision gadgets (like γ) that are smaller than γ^* , and then select γ^* as its minimum decision gadget. After that time, p always selects the deciding process of γ^* — which is precisely p^* , in this case.

Theorem 25: For all correct processes p , there is a time after which $\text{output}_p = p^*$, forever.

PROOF: Let i^* denote the critical index selected by Line (1) of Figure 8 applied to Υ . By Lemma 24, there is a time t_{init} after which every root of Υ_p has the same tags as the corresponding root of Υ . Thus after time t_{init} , p always sets $i = i^*$ in Line (1) of Figure 10. We now show that there is a time after which the computation component of p (Figure 10) always returns p^* . There are two cases:

1. i^* is monovalent critical. In this case, p^* is process p_{i^*} (by Line (2) of the selection rule Figure 8). Similarly, after time t_{init} : (a) p always sets i to i^* (Line (1) of Figure 10); (b) p always returns p_{i^*} (Line (2) of Figure 10).

2. i^* is bivalent critical. Let γ^* denote the smallest decision gadget of Υ^{i^*} . In this case, p^* is the deciding process of γ^* . Since γ^* is a finite subgraph of Υ^{i^*} , by Lemma 23(3), there is a time after which γ^* is also a subgraph of Υ_p^i . By Lemma 24, there is a time t_{γ^*} after which all the (finitely many) vertices of γ^* receive the same tags in Υ^{i^*} and $\Upsilon_p^{i^*}$. Thus after time t_{γ^*} , γ^* is also a decision gadget of $\Upsilon_p^{i^*}$.

Since each graph is encoded as a unique natural number, there are finitely many graphs with a smaller encoding than γ^* . Let \mathcal{G} denote the set of graphs with a smaller encoding than γ^* , and γ be any graph in \mathcal{G} . We show that there is a time after which γ is not a decision gadget of $\Upsilon_p^{i^*}$. There are two cases:

- (a) γ is not a subgraph of Υ^{i^*} . In this case, by Lemma 23(1), γ is never a subgraph of $\Upsilon_p^{i^*}$.
- (b) γ is a subgraph of Υ^{i^*} . Since γ^* is the smallest decision gadget of Υ^{i^*} and γ is smaller than γ^* , γ is not a decision gadget of Υ^{i^*} . By Lemma 24, there is a time t_γ after which all the (finitely many) vertices of γ have the same tags in Υ^{i^*} and $\Upsilon_p^{i^*}$. Thus after time t_γ , γ is not a decision gadget of $\Upsilon_p^{i^*}$.

Since \mathcal{G} is finite, there is a time $t_{\mathcal{G}}$ after which no graph in \mathcal{G} is a decision gadget of Υ_p^i .

Consider the process that is returned by the computation component of p (Figure 10) at any time $t > \max(t_{init}, t_{\gamma^*}, t_{\mathcal{G}})$. Since $t > t_{init}$, p always sets i to i^* in Line (1). Since $t > t_{\gamma^*}$, γ^* is a decision gadget of $\Upsilon_p^i(t)$. Finally, since $t > t_{\mathcal{G}}$, γ^* is the smallest decision gadget of $\Upsilon_p^i(t)$. Thus, since i^* is bivalent, at any time after $\max(t_{init}, t_{\gamma^*}, t_{\mathcal{G}})$, Line (3) of Figure 10 returns the deciding process of γ^* . Therefore, after time $\max(t_{init}, t_{\gamma^*}, t_{\mathcal{G}})$, the computation component of p always returns p^* .

From the above, there is a time after which p sets $output_p \leftarrow p^*$, forever, in Line (3) of Figure 9. \square

We now have all the pieces needed to prove our main result, Theorem 2 in Section 5:

Theorem 2: For all environments \mathcal{E} , if a failure detector \mathcal{D} can be used to solve Consensus in \mathcal{E} , then $\mathcal{D} \succeq_{\mathcal{E}} \Omega$.

PROOF: Consider the execution of algorithm $T_{\mathcal{D} \rightarrow \Omega}$ in any environment \mathcal{E} . By Theorem 25, there is a time after which all correct processes set $output_p = p^*$, forever. By Theorem 22, p^* is a correct process. Thus, $T_{\mathcal{D} \rightarrow \Omega}$ is a reduction algorithm that transforms \mathcal{D} into Ω . In other words, Ω is reducible to \mathcal{D} . \square

7 Discussion

7.1 Granularity of atomic actions

Our model incorporates very strong assumptions about the atomicity of steps. First, the three phases of each step are assumed to occur indivisibly, and at a single time. In

particular, the failure of a process cannot happen in the “middle of a step”. This allows us to associate a single time t with a step and think of the step as occurring at that time. Second, in the send phase of a step a message is sent to *all* processes. Given that the entire step is indivisible, this means that either all or none of the correct processes eventually receive the message sent in a step. Finally, no two steps can occur at the same time.¹² These assumptions are convenient because they make the formal model simpler to describe. Also, they are consistent with those made in the model of [FLP85] that provided the impetus for this work.

On the other hand, in [CT91] a model with weaker properties is used. There, the three phases of a step need not occur indivisibly, and may occur at different times. Even within the send phase, the messages sent to the different processes may be sent at different times. Thus, a failure may occur in the middle of the send phase, resulting in some correct processes eventually receiving the messages sent to them in that step while others never do. Also, actions of *different* processes may take place simultaneously, subject to the restriction that a message can only be received strictly *after* it was sent. Since [CT91] is mainly concerned with showing how to use various types of failure detectors to achieve Consensus, the use of a weaker model strengthens the results. (In fact, the negative results of [CT91] hold even in the model of this paper, with the stronger atomicity assumptions.)

The question naturally arises whether our result also applies to this weaker model. In other words, if a failure detector \mathcal{D} can be used to solve Consensus in the weak model, is it true that we can transform \mathcal{D} to \mathcal{W} *in that model*? It turns out that the answer is affirmative. To see this, first note that if \mathcal{D} solves Consensus in the weak model then it surely solves Consensus in the strong model. By our result, \mathcal{D} can be transformed to \mathcal{W} in the strong model. It remains to show that \mathcal{D} can be transformed to \mathcal{W} in the weak model. This is not obvious, since it is conceivable that the extra properties of the strong model are crucial in the transformation of \mathcal{D} to \mathcal{W} . Fortunately, the transformation presented in this paper actually works even in the weak model!

To see this, it is sufficient to ensure that the communication component of the transformation (cf. Figure 9 in Section 6.6.1) constructs graphs that satisfy the properties listed in Lemma 6.6.1.4, *even* if we run it in the weak model. It is not difficult to verify that this is indeed so. The proof is virtually the same, except for the fact that we must distinguish the time t in which a process p queries its failure detector and the time t' in which p adds the value it saw into G_p . In our proof we assume that $t = t'$; in the weak model we would have $t \leq t'$. Similar comments apply to all actions within a step that are no longer assumed to occur at the same instant of time. These changes make the proofs slightly more cumbersome, since we must introduce notation for all the different times in which relevant actions within a step take place, but the reasoning remains essentially the same.¹³

¹²This is reflected in our formal model by the fact that the list of times in a run (which indicate when the events in the run’s schedule occur) is *increasing*.

¹³Another problem that must be confronted is that in the proofs of Lemmata 6.6.1.3 and 6.6.1.4 we often refer to the “first graph” in which a vertex or edge is present. In the strong model there is no difficulty with this, since processes cannot execute steps simultaneously. In the weak model, we have to justify that it makes sense to speak of the “first” graph to contain a vertex or edge, in spite of the fact

Thus, our result is not merely a fortuitous consequence of some whimsical choice of model. We view the robustness of the result across different models of asynchrony as further testimony to the significance of the failure detector \mathcal{W} .

7.2 Failure detection and partial synchrony

The fundamental reason why Consensus cannot be solved in completely asynchronous systems is the fact that, in such systems, it is impossible to reliably distinguish a process that has crashed from one that is merely very slow. In other words, Consensus is unsolvable because accurate failure detection is impossible. On the other hand, it is well-known that Consensus is solvable (deterministically) in completely synchronous systems — that is, systems where all processes take steps at the same rate and each message arrives at its destination a fixed and known amount of time after it is sent. In such a system we can use timeouts to implement a “perfect” failure detector — i.e., one in which no process is ever wrongly suspected, and every faulty process is eventually suspected. Thus the ability to solve Consensus in a given system is intimately related to the failure detection capabilities of that system. This realization led to the extension of the asynchronous model of computation with failure detectors in [CT91]. In that paper Consensus is shown to be solvable even with very weak failure detectors that could make an infinite number of “mistakes”.

A different tack on circumventing the unsolvability of Consensus is pursued in [DDS87] and [DLS88]. The approach of those papers is based on the observation that between the completely synchronous and completely asynchronous models of distributed systems there lie a variety of intermediate “partially synchronous” models. For instance, in one model of partial synchrony, processes take steps at the same rate, but message delays are unbounded (albeit finite). Alternatively, it may be known that message delays are bounded, but the actual bound may be unknown. In yet another variation, the *eventual* maximum message delay is known, but during some initial period of finite but unknown duration some messages may experience longer delays. These and many other models of partial synchrony are studied in [DDS87] and [DLS88], and the question of solvability of Consensus in each of them is answered either positively or negatively.

In particular, [DDS87] defines a space of 32 models by considering five key parameters, each of which admits a “favorable” and an “unfavorable” setting. For instance, one of the parameters is whether the maximum message delay is known (favorable setting) or not (unfavorable setting). Each of the 32 models corresponds to a particular setting of the 5 parameters. [DDS87] identifies four “minimal” models in which Consensus is solvable. These are minimal in the sense that the weakening of any parameter from favorable to unfavorable would yield a model of partial synchrony where Consensus is unsolvable. Thus, within the space of the models considered, [DDS87] and [DLS88] delineate precisely the boundary between solvability and unsolvability of Consensus, and provide an answer to the question “What is the least amount of synchrony sufficient to solve Consensus?”.

that certain actions can be executed at the same time. The fact that a message can be received only *after* it was sent is needed here.

Failure detectors can be viewed as a more abstract and modular way of incorporating partial synchrony assumptions into the model of computation. Instead of focusing on the *operational features* of partial synchrony (such as the five parameters considered in [DDS87]), we can consider the *axiomatic properties* that failure detectors must have in order to solve Consensus. The problem of implementing a given failure detector in a specific model of partial synchrony becomes a separate issue; this separation affords greater modularity.

To see the connection between partial synchrony and failure detectors, it is useful to examine how one might go about implementing a failure detector. By the impossibility result of [FLP85], a failure detector that can be used to solve Consensus cannot be implemented in a completely asynchronous system. Now consider partially synchronous systems in which correct processes have accurate timers (i.e., they can measure elapsed time). If in such a system message delays are bounded and the maximum delay is known, we can use timeouts to implement the “perfect” failure detector described above. In a weaker system where message delays are bounded but the maximum delay is *not* known, we can implement a failure detector satisfying a weaker property: *eventually* no correct process is suspected. This can be done by using timeouts of increasing length; once the timeout period has been increased sufficiently to exceed the unknown maximum delay, no correct process will be suspected. A failure detector with the same property can also be implemented in a distributed system where the *eventual* maximum message delay is known, but messages may be delayed for longer during some initial period of finite but unknown duration. With these remarks we illustrate two points: First, that stronger failure detectors correspond to stronger models of partial synchrony; and second, that the same failure detector can be implemented in different models of partial synchrony.

Studying failure detectors rather than various models of partial synchrony has several advantages. By determining whether Consensus is solvable using some specific failure detector we thereby determine whether Consensus is solvable in *all* systems in which that failure detector can be implemented. An algorithm that relies on the axiomatic properties of a given failure detector is more general, more modular, and simpler to understand than one that relies directly on some specific operational features of partial synchrony (that can be used to implement the given failure detector).

From this more abstract point of view, the question “What is the least amount of synchrony sufficient to solve Consensus?” translates to “What is the weakest failure detector sufficient to solve Consensus?”. In contrast to [DDS87], which identified a *set* of *minimal* models of partial synchrony in which Consensus is solvable, we are able to exhibit a *single minimum* failure detector that can be used to solve Consensus. The technical device that made this possible is the notion of *reduction* between failure detectors. We suspect that a corresponding notion of reduction between models of partial synchrony, although possible, would be more complex. This is because there are models which are not comparable in general (in the sense that there are tasks that are possible in one but not in the other and vice versa), although they are comparable *as far as failure detection is concerned* — which is all that matters for solving Consensus! In this connection, it is useful to recall our earlier observation, that the same failure detector can be implemented in different (indeed, incomparable) models of partial synchrony.

7.3 Weak Consensus

[FLP85] actually showed that even the *Weak Consensus* problem cannot be solved (deterministically) in an asynchronous system. Weak Consensus is like Consensus except that the validity property is replaced by the following, weaker, property:

Non-triviality: There is a run of the protocol in which correct processes decide 0, and a run in which correct processes decide 1.

Unlike validity, this property does not explicitly prescribe conditions under which the correct processes must decide 0 or 1 — it merely states that it is possible for them to reach each of these decisions. It is natural to ask whether our result holds for this weaker problem as well. That is, we would like to know if the following holds:

Theorem: For all environments \mathcal{E} , if a failure detector \mathcal{D} can be used to solve *Weak Consensus* in \mathcal{E} , then $\mathcal{D} \succeq_{\mathcal{E}} \Omega$.

Under the above definition of non-triviality, this is not quite right. But as we shall argue, the problem really lies with the definition! Under a slightly stronger definition, which is more appropriate for our model that incorporates failure detectors, the theorem is actually true.

Intuitively, the problem with the above definition of non-triviality in our model of failure detectors is that it is possible for the decision of correct processes to depend entirely on the the values returned by the failure detector. Consider, for example, a failure detector \mathcal{D} so that for each failure pattern F , $\mathcal{D}(F) = \{H_0, H_1\}$, where for all processes p and times t , and for all $i \in \{0, 1\}$, $H_i(p, t) = i$. In other words, in any given run, this failure detector returns the same binary value to all processes at all times, independent of the run’s failure pattern. It is trivial to use this failure detector to solve Weak Consensus: A process merely queries its failure detector and decides the value returned! It is easy to see that $\mathcal{D} \not\succeq_{\mathcal{E}} \Omega$, for any environment \mathcal{E} : \mathcal{D} provides absolutely no information about which processes are correct or faulty.¹⁴

At this point, the reader may justifiably object that \mathcal{D} is “cheating” — it is really not a failure detector, but a mechanism that non-deterministically chooses the decision value. One possible way of fixing this problem would be to make our definition of failure detector less general than it presently is. We could then try to prove the theorem for this restricted definition of failure detectors. This approach, however, is fraught with the danger of restricting the definition too much and ruling out legitimate failure detectors in addition to bogus ones, like \mathcal{D} . Intuitively, the failure detector is supposed to provide some information about faulty processes. As this information may be encoded in a complex way, we should not arbitrarily rule out such encodings because, in doing so, we may be inadvertently ruling out useful failure detectors.

Instead of modifying our definition of failure detector, we strengthen the non-triviality property to require that the failure detector values seen by the processes do not, by themselves, determine the decision value. To formalize this, let R be a run of a Consensus algorithm, and $(p_1, m_1, d_1), (p_2, m_2, d_2), \dots$, be the schedule of R . We denote by $\mathbf{fd}(R)$

¹⁴In fact, \mathcal{D} cannot be used to solve Consensus.

the sequence $[p_1, d_1], [p_2, d_2], \dots$, i.e., the sequence of failure detector values seen by the processes in R . Consider the relation \equiv on runs, where $R \equiv R'$ if and only if $\mathbf{fd}(R) = \mathbf{fd}(R')$. It is immediate that \equiv is an equivalence relation. We now redefine the non-triviality property in our model (where processes have access to failure detectors) as follows:

Non-triviality: In every equivalence class of the relation \equiv , there is a run of the protocol in which correct processes decide 0, and a run in which correct processes decide 1.

This captures the idea that the decision value cannot be ascertained merely on the basis of the failure detector values seen by the processes. It must also depend on other aspects of the run (such as the initial values, the particular messages sent, or other features).

If we define Weak Consensus using *this* version of non-triviality, then the Theorem stated above is, in fact, true. We briefly sketch the modifications of our proof needed to obtain this strengthening of Theorem 2. The only use of the validity property is in the proof of Lemma 16 which states that the root of Υ^0 is 0-valent and the root of Υ^n is 1-valent. This, in turn, is used in the proof of Lemma 17, which states that a critical index exists.

To prove the stronger theorem, we concentrate on the forest induced by *all* initial configurations — not just I^0, \dots, I^n . Thus, the forest now will have 2^n trees, rather than only $n+1$. Consider the n initial values of processes in an initial configuration as an n -bit vector, and fix any n -bit Gray code.¹⁵ Let I^0, \dots, I^{2^n-1} be the initial configurations listed in the order specified by the Gray code, and Υ^i be the tree $\Upsilon_G^{I^i}$, for all $i \in \{0, \dots, 2^n - 1\}$. We use the same definition for a critical index as we had before: Index $i \in \{0, \dots, 2^n - 1\}$ is *critical* if the root of Υ^i is bivalent or the root of Υ^i is 1-valent while the root of Υ^{i-1} is 0-valent. The only difference is that we now take subtraction to be modulo 2^n , so that when $i = 0$, $i - 1 = -1 = 2^n - 1$. We can now prove an analogue to Lemma 17.

Lemma: There is a critical index i , $0 \leq i \leq 2^n - 1$.

PROOF: First, we claim that the forest contains both nodes tagged 0 and nodes tagged 1. To see this, let S be a node in some tree of the forest. By Lemma 11, S has a tag; without loss of generality, assume that S has tag 0. Consider an infinite path that extends S . By Lemma 6 and the fact that S is tagged 0, there is a run R of the Weak Consensus algorithm in which a correct process decides 0. By non-triviality, there is a run $R' \equiv R$, so that correct processes decide 1 in R' . Let S' be the infinite schedule and I^ℓ be the initial configuration of run R' . Using the definition of the \equiv relation and the construction of the induced forest, it is easy to show that every finite prefix of S' is a node of Υ^ℓ . Since correct processes decide 1 in R' , all these nodes are tagged 1.

Since there are both nodes tagged 0 and nodes tagged 1, by Lemma 13, there are both roots tagged 0 and roots tagged 1. If the root of some Υ^i is tagged both 0 and 1, it is bivalent and we are done. Otherwise, the roots of all trees are monovalent, and there are both 0- and 1-valent roots. Thus, there exist $0 \leq i, j \leq 2^n - 1$ so that the root of Υ^i

¹⁵An n -bit Gray code is a sequence of all possible n -bit vectors where successive vectors, as well as the first and last vectors, differ only in the value of one position. It is well-known that such codes exist for all $n \geq 1$.

is 0-valent and the root of Υ^j is 1-valent. By considering the sequence $\Upsilon^i, \Upsilon^{i+1}, \dots, \Upsilon^j$, (where addition is modulo 2^n) it is easy to see that the root of some Υ^k , $k \neq i$, that appears in that sequence is 1-valent, while the root of Υ^{k-1} is 0-valent. By definition, k is a critical index. \square

The rest of the proof remains unchanged.

7.4 Failure detectors with infinite range of output values

The failure detectors in [RB91, CT91] only output lists of processes suspected to have crashed. Since the set of processes is finite, the range of possible output values of these failure detectors is also finite. In this paper our model allows for failure detectors with arbitrary ranges of output values, including the possibility of infinite ranges! We illustrate the significance of this generality by describing a natural class of failure detectors whose range of output values is infinite (though each value output is finite).

Example: One apparent weakness with our formulation of failure detection is that a brief change in the value output by a failure detector module may go unnoticed. For example, process p 's module of the given failure detector, \mathcal{D} , may output d_1 at time t_1 , d_2 at a later time t_2 and d_1 again at time t_3 after t_2 . If, due to the asynchrony of the system, p does not take a step between time t_1 and t_3 , p may never notice that its failure detector module briefly output d_2 . A natural way of overcoming this problem is to replace \mathcal{D} with failure detector \mathcal{D}' that has the following property: \mathcal{D}' maintains the same list of suspects as \mathcal{D} but when queried, \mathcal{D}' returns the entire history of its list of suspects up to the present time. In this manner, correct processes are guaranteed to notice every change in \mathcal{D}' 's list of suspects. As the system continues executing, the values output by \mathcal{D}' grow in size. Thus \mathcal{D}' has an infinite range of output values.

However, since \mathcal{D} is a function of F , the failure pattern encountered, \mathcal{D}' is also a function of F , and can be described by our model. Thus, the result in this paper applies to \mathcal{D}' , a natural failure detector with an infinite range of output values.

Appendix

In this appendix, G is a DAG that satisfies Properties 1–4 listed in Section 6.2. To give the proofs of the lemmata in Section 6.2, we first show two auxiliary results.

Lemma 26: Let V be any finite subset of vertices in G . G has an infinite path g such that:

- There is an edge from every vertex of V to the first vertex of g .
- If $[p, d, k]$ is a vertex of g then p is correct; for each correct p , there are infinitely many vertices of the form $[p, -, -]$ in g .

PROOF: By repeated application of Property 4 of the DAG G . \square

Lemma 27: S is a schedule associated with a path of Υ_G^I that starts from the root of Υ_G^I if and only if S is a schedule compatible with G and applicable to I .

PROOF: The lemma obviously holds if S is a *finite* schedule (this is immediate from the definitions). Now let $S = e_1, e_2, \dots, e_i, \dots$ be an *infinite* schedule, where $e_i = (q_i, m_i, d_i)$. We define $S_0 = S_{\perp}$, $S_1 = e_1$, $S_2 = S_1 \cdot e_2$, and in general $S_i = S_{i-1} \cdot e_i$ for all $i \geq 1$.

Assume that S is compatible with G and applicable to I . We must show that S is a schedule associated with a path of Υ_G^I that starts from the root of Υ_G^I . To see this, note that for all $i \geq 0$, S_i is a finite schedule that is also compatible with G and applicable to I . Thus, all the schedules $S_0, S_1, S_2, \dots, S_{i-1}, S_i, \dots$ are vertices of Υ_G^I . Since $S_i = S_{i-1} \cdot e_i$, the tree Υ_G^I has an edge from S_{i-1} to S_i which is labeled e_i , for all $i \geq 1$. Thus, $S = e_1, e_2, \dots, e_i, \dots$ is the schedule associated with the infinite path $S_0 \rightarrow S_1 \rightarrow S_2 \rightarrow \dots \rightarrow S_{i-1} \rightarrow S_i \rightarrow \dots$ of Υ_G^I ; this path starts from the root $S_0 = S_{\perp}$ of Υ_G^I .

Assume that S is a schedule associated with an infinite path of Υ_G^I that starts from the root of Υ_G^I . We must show that S is compatible with G and is applicable to I . First note that for all i , S_i is a vertex in Υ_G^I , thus S_i is applicable to I and compatible with G . From the definition of applicability, since S_i is applicable to I for all $i \geq 1$, S is applicable to I . It remains to show that S is compatible with G .

For all $i \geq 1$, S_i is compatible with G , so S_i is compatible with a finite path π_i in G . Since $S_i = e_1, e_2, \dots, e_i$ where $e_j = (q_j, m_j, d_j)$,

$$\pi_i = [q_1, d_1, k_1^i], [q_2, d_2, k_2^i], \dots, [q_i, d_i, k_i^i] \quad (1)$$

for some $k_1^i, k_2^i, \dots, k_i^i \in \mathbf{N}$. We now show that S is compatible with some infinite path π in G .

For each positive integer i , consider the set $\{k_i^i, k_i^{i+1}, k_i^{i+2}, \dots\}$. This set has a minimum element (because it is a non-empty set of natural numbers). Let $\mu(i)$ be a superscript of that minimum element, i.e., $k_i^{\mu(i)} \leq k_i^j$ for all $j \geq i$. Note that $\mu(i) \geq i$. We claim that

$$\pi = [q_1, d_1, k_1^{\mu(1)}], [q_2, d_2, k_2^{\mu(2)}], \dots, [q_i, d_i, k_i^{\mu(i)}], [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}], \dots$$

is an infinite path in G and that S is compatible with π .

The fact that S is compatible with π is immediate from the definitions of S , π , and compatibility. To prove that π is an infinite path in G , we show that for all $i \geq 1$, $[q_i, d_i, k_i^{\mu(i)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}]$ is an edge in G . By (1), and the fact that $\mu(i+1) \geq i+1$, path $\pi_{\mu(i+1)}$ contains the edge $[q_i, d_i, k_i^{\mu(i+1)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}]$. Since $\pi_{\mu(i+1)}$ is a path in G ,

$$\text{edge } [q_i, d_i, k_i^{\mu(i+1)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}] \text{ is in } G. \quad (2)$$

By definition of $\mu(i)$, $k_i^{\mu(i)} \leq k_i^{\mu(i+1)}$. There are two possible cases:

1. $k_i^{\mu(i)} = k_i^{\mu(i+1)}$. In this case, the vertices $[q_i, d_i, k_i^{\mu(i)}]$ and $[q_i, d_i, k_i^{\mu(i+1)}]$ coincide, and by (2), edge $[q_i, d_i, k_i^{\mu(i)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}]$ is in G .
2. $k_i^{\mu(i)} < k_i^{\mu(i+1)}$. In this case, by Property 2 of G , $[q_i, d_i, k_i^{\mu(i)}] \rightarrow [q_i, d_i, k_i^{\mu(i+1)}]$ is an edge in G . By (2) and Property 3 (transitivity) of G , edge $[q_i, d_i, k_i^{\mu(i)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}]$ is in G .

Thus, in both cases, edge $[q_i, d_i, k_i^{\mu(i)}] \rightarrow [q_{i+1}, d_{i+1}, k_{i+1}^{\mu(i+1)}]$ is in G . \square

Lemma 5: Let S be a schedule associated with a *finite* path of Υ_G^I that starts from the root of Υ_G^I . There is a sequence of times T such that $\langle F, H_{\mathcal{D}}, I, S, T \rangle$ is a partial run of *Consensus $_{\mathcal{D}}$* .

PROOF: By Lemma 27, S is applicable to I and compatible with G . Thus S is compatible with some finite path $g = [q_1, d_1, k_1], [q_2, d_2, k_2], \dots, [q_i, d_i, k_i], \dots, [q_\ell, d_\ell, k_\ell]$ of G . For all i , $1 \leq i \leq \ell$, let t_i be the label of vertex $[q_i, d_i, k_i]$. From Property 1a of G (applied to every vertex of the path g), for all i , $1 \leq i \leq \ell$, $d_i = H_{\mathcal{D}}(q_i, t_i)$ and $q_i \notin F(t_i)$. From Property 1b of G (applied to every edge of the path g), for all i , $1 \leq i < \ell$, $t_i < t_{i+1}$. Thus T is a sequence of increasing times, and, by definition, $\langle F, H_{\mathcal{D}}, I, S, T \rangle$ is a partial run of *Consensus $_{\mathcal{D}}$* . \square

Lemma 6: Let S be a schedule associated with an *infinite* path of Υ_G^I that starts from the root Υ_G^I . If in S every correct process takes an infinite number of steps and every message sent to a correct process is eventually received, there is a sequence of times T such that $\langle F, H_{\mathcal{D}}, I, S, T \rangle$ is a run of *Consensus $_{\mathcal{D}}$* .

PROOF: Similar to Lemma 5. \square

Lemma 7: For any two initial configurations I and I' , if S is a vertex of Υ_G^I and is applicable to I' then S is also a vertex of $\Upsilon_G^{I'}$.

PROOF: Follows directly from the definitions. \square

Lemma 8: Let S be any vertex of Υ_G^I and p be any correct process. Let m be a message in the message buffer of $S(I)$ addressed to p or the null message. For some d , S has a child $S \cdot (p, m, d)$ in Υ_G^I .

PROOF: From the definition of Υ_G^I , S is compatible with some finite path g of G and applicable to I . Let v denote the last vertex of g . By Property 4, there is a d and a k such that $v \rightarrow [p, d, k]$ is an edge of G . Therefore, $g \cdot [p, d, k]$ is a path of G , and $S \cdot (p, m, d)$ is compatible with G .

It remains to show that $S \cdot (p, m, d)$ is applicable to I . Since S is applicable to I , it suffices to show that (p, m, d) is applicable to $S(I)$. But this is true since, by hypothesis, m is in the message buffer of $S(I)$ and addressed to p , or the null message. \square

Lemma 9: Let S be any vertex of Υ_G^I and p be any process. Let m be a message in the message buffer of $S(I)$ addressed to p or the null message. Let S' be a descendent of S such that, for some d , $S' \cdot (p, m, d)$ is in Υ_G^I . For each vertex S'' on the path from S to S' (inclusive), $S'' \cdot (p, m, d)$ is also in Υ_G^I .

PROOF: Since they are vertices of Υ_G^I , S , S'' and $S' \cdot (p, m, d)$ are compatible with some finite paths g , $g \cdot g''$ and $g \cdot g'' \cdot g' \cdot [p, d, k]$ of G , respectively. From Property 3 (transitive closure) of G , $g \cdot g'' \cdot [p, d, k]$ is also a path of G . So $S'' \cdot (p, m, d)$ is compatible with this path of G . We now show that $S'' \cdot (p, m, d)$ is also applicable to I , and therefore it is a vertex of Υ_G^I .

Since S'' is a vertex of Υ_G^I , S'' is applicable to I . If $m = \lambda$, then (p, m, d) is obviously applicable to $S''(I)$. Now suppose $m \neq \lambda$. Since $S' \cdot (p, m, d)$ is a vertex of Υ_G^I , (p, m, d) is applicable to $S'(I)$, and thus m is in the message buffer of $S'(I)$. Since each message is sent at most once and m is in the message buffers of $S(I)$ and $S'(I)$, there is no edge of the form $(p, m, -)$ on the path from S to S' . So m is also in the message buffer of $S''(I)$, and (p, m, d) is applicable to $S''(I)$. \square

```

j ← 0
S0 ← S                                {S0 is compatible with g and applicable to I}
repeat forever
  j ← j + 1
  Let [qj, dj, kj] be the j-th vertex of path g∞
  Let mj be the oldest message addressed to qj in the message buffer of Sj-1(I)
  (if no such message exists, mj = λ)
  ej ← (qj, mj, dj)
  Sj ← Sj-1 · ej                    {Sj is compatible with g · [q1, d1, k1] · ... · [qj, dj, kj]}
  {and applicable to I}

```

Figure 11: Generating schedule $S \cdot E^\infty$, compatible with path $g \cdot g_\infty$, in Υ_G^I

Lemma 10: Let S, S_0 , and S_1 be any vertices of Υ_G^I . There is a finite schedule E containing only steps of correct processes such that:

1. $S \cdot E$ is a vertex of Υ_G^I and all correct processes have decided in $S \cdot E(I)$.
2. For $i = 0, 1$, if E is applicable to $S_i(I)$ then $S_i \cdot E$ is a vertex of Υ_G^I .

PROOF: Since S is a vertex of Υ_G^I , S is compatible with some finite path g of G and is applicable to I . Similarly, S_0 and S_1 are compatible with some finite paths g_0 and g_1 , respectively, of G . From Lemma 26 (applied to the last vertices of g, g_0 and g_1), G has an infinite path $g_\infty = [q_1, d_1, k_1], [q_2, d_2, k_2], \dots, [q_j, d_j, k_j], \dots$ with the following two properties:

1. There is an edge from the last vertex of g, g_0 and g_1 to the first vertex of g_∞ . (Thus, $g \cdot g_\infty, g_0 \cdot g_\infty$, and $g_1 \cdot g_\infty$ are infinite paths in G .)
2. If $[p, d, k]$ is a vertex of g_∞ then p is correct; for each correct p , there are infinitely many vertices of the form $[p, -, -]$ in g_∞ .

We now show how to construct the required schedule E . Consider the infinite sequence of schedules $S^0, S^1, S^2, \dots, S^j, \dots$ constructed by the algorithm in Figure 11. An easy induction shows that for all $j > 0$, S^j is applicable to I and is compatible with $g \cdot [q_1, d_1, k_1] \cdot \dots \cdot [q_j, d_j, k_j]$, a prefix of the path $g \cdot g_\infty$ in G . So, for all $j > 0$, S^j is a vertex

of Υ_G^I . Consider the infinite path of Υ_G^I that starts from the root of Υ_G^I then goes to $S^0 = S$, and then to $S^1, S^2, \dots, S^j, \dots$. The infinite schedule associated with that path is $S^\infty = S \cdot e_1 \cdot e_2 \cdot \dots \cdot e_j \dots$. Note that schedule $E^\infty = e_1 \cdot e_2 \cdot \dots \cdot e_j \dots$ is compatible with path g_∞ of G . By Property (2) of path g_∞ , every correct process p takes an infinite number of steps in E^∞ (and thus also in $S^\infty = S \cdot E^\infty$). Since in each one of these steps p receives the oldest message that is addressed to it, every message sent to p (in S^∞) is eventually received. By Lemma 6, there is a T such that $R = \langle F, H_{\mathcal{D}}, I, S^\infty, T \rangle$ is a run of *Consensus $_{\mathcal{D}}$* .

From the termination requirement of *Consensus*, S^∞ has a finite prefix S^d such that all correct processes have decided in $S^d(I)$. There are two cases:

- S^d is a prefix of S . Since decisions are irrevocable, all correct processes remain decided in $S(I)$. Thus S_\perp , the empty schedule, is the required E .
- S is a prefix of S^d . Thus, $S^d = S \cdot E$ where E is a finite prefix of E^∞ . Since E^∞ is compatible with g_∞ , E is compatible with a prefix of g_∞ . Now consider S_0 (the following argument also applies to S_1). Since S_0 is compatible with g_0 , $S_0 \cdot E$ is compatible with a prefix of $g_0 \cdot g_\infty$, a path in G . So, $S_0 \cdot E$ is compatible with G . If $S_0 \cdot E$ is also applicable to I , then, by the definition of Υ_G^I , it is a vertex of Υ_G^I . The same argument holds for S_1 . It remains to show that E contains only steps of correct processes. This is immediate from Property (2) of g_∞ and from the fact that E is compatible with a prefix of g_∞ . \square

Acknowledgements

We would like to thank Prasad Jayanti for his valuable comments on various versions of this paper. We would also like to thank Cynthia Dwork, Dexter Kozen and the distributed systems group at Cornell for many helpful discussions, and the anonymous referees.

References

- [ABD⁺87] Hagit Attiya, Amotz Bar-Noy, Danny Dolev, Daphne Koller, David Peleg, and Rüdiger Reischuk. Achievable cases in an asynchronous environment. In *Proceedings of the Twenty-Eighth Symposium on Foundations of Computer Science*, pages 337–346. IEEE Computer Society Press, October 1987.
- [BMZ88] Ofer Biran, Shlomo Moran, and Shmuel Zaks. A combinatorial characterization of the distributed tasks that are solvable in the presence of one faulty processor. In *Proceedings of the Seventh ACM Symposium on Principles of Distributed Computing*, pages 263–275. ACM Press, August 1988.
- [BW87] Michael Bridgland and Ronald Watro. Fault-tolerant decision making in totally asynchronous distributed systems. In *Proceedings of the Sixth ACM*

Symposium on Principles of Distributed Computing, pages 52–63. ACM Press, August 1987.

- [CD89] Benny Chor and Cynthia Dwork. Randomization in byzantine agreement. *Advances in Computer Research*, 5:443–497, 1989.
- [CT91] Tushar Deepak Chandra and Sam Toueg. Unreliable failure detectors for asynchronous systems. In *Proceedings of the Tenth ACM Symposium on Principles of Distributed Computing*, pages 325–340. ACM Press, August 1991. To appear in the *Journal of the ACM*. An extended and revised version is also available by anonymous ftp from `ftp.cs.cornell.edu` in `pub/sam/failure.detectors.algorithms.ps`.
- [DDS87] Danny Dolev, Cynthia Dwork, and Larry Stockmeyer. On the minimal synchronism needed for distributed consensus. *Journal of the ACM*, 34(1):77–97, January 1987.
- [DLP⁺86] Danny Dolev, Nancy A. Lynch, Shlomit S. Pinter, Eugene W. Stark, and William E. Weihl. Reaching approximate agreement in the presence of faults. *Journal of the ACM*, 33(3):499–516, July 1986.
- [DLS88] Cynthia Dwork, Nancy A. Lynch, and Larry Stockmeyer. Consensus in the presence of partial synchrony. *Journal of the ACM*, 35(2):288–323, April 1988.
- [FLP85] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, April 1985.
- [RB91] Aleta Ricciardi and Kenneth P. Birman. Using process groups to implement failure detection in asynchronous environments. In *Proceedings of the Tenth ACM Symposium on Principles of Distributed Computing*, pages 341–351. ACM Press, August 1991.