

Minimal System Conditions to Implement Unreliable Failure Detectors *

Antonio Fernández[†] Ernesto Jiménez[‡] Sergio Arévalo[†]

Abstract

In this paper we explore the minimal system requirements to implement unreliable failure detectors. We first consider systems formed by lossy asynchronous and eventually timely links. On these systems we define two properties, the Weak Property and the Strong Property, depending on whether all correct processes can be reached with links that are not lossy asynchronous from one or from all correct processes, respectively. We present necessary conditions based on these properties. We show that there is no algorithm that implements $\diamond\mathcal{S}$, Ω , nor \mathcal{S} (resp. $\diamond\mathcal{P}$ nor \mathcal{P}) if we allow one single failure in a system that, when all processes are correct, does not satisfy the Weak (resp. Strong) Property. Then, we propose an algorithm that implements $\diamond\mathcal{P}$ if the Strong Property is satisfied, and $\diamond\mathcal{S}$ (and Ω with an additional assumption) if only the Weak Property is satisfied. For systems formed by synchronous and lossy asynchronous links only, we propose another algorithm that implements detector class \mathcal{P}_4 if the Strong Property is satisfied, and implements a new detector class \mathcal{S}' (and Ω with an additional assumption) if only the Weak Property is satisfied.

1. Introduction

Unreliable failure detectors were proposed by Chandra and Toueg [6] as devices that can be used to solve consensus in asynchronous systems with crash failures. These devices allow to circumvent the seminal result from Fischer et al [13] that shows the impossibility of solving consensus in these systems even with one single failure. Beyond consensus, it is also known that unreliable failure detectors are useful to solve other fundamental problems in distributed computing. For these reasons there is a growing interest in deriving practical efficient algorithms to implement failure detectors [4, 3, 2, 1, 7, 18, 17, 19].

In a tutorial at PODC 2002, Keidar and Rajsbaum [15] asked, among other questions, for the weakest requirements on systems that allow to implement the different classes of failure detectors. In their recent works Aguilera et al. [2, 3] have partially answered this question for the class of failure detectors Ω [5]. As usual, Aguilera et al. consider that all the processes are connected by directed links. They consider several classes of links. A *timely* link is a reliable link with a known upper bound on the message delivery time. An *eventually timely* link is a link that behaves timely (but the bound may not be known) after some unknown stabilization time. Before this time messages can be lost. A *fair lossy* link is one in which if infinite messages are sent, infinite messages are received (but infinite messages can also be lost). Finally, a *lossy asynchronous* link is one in which messages can be lost and there is no bound on the message delivery time. In [1] they present an algorithm that implements Ω in a system with all links lossy asynchronous except the (input and output) links of an (unknown) correct process. Later in [2] they show that a failure detector of class Ω can be implemented in a system with lossy asynchronous links as long as there is one (unknown) correct process whose *output* links are eventually timely. However, they prove that in this case any algorithm will have runs in which $\Omega(n^2)$ links are permanently carrying messages. They also show that if, additionally, some (unknown) correct process has all its (input and output) links fair lossy, Ω can be implemented efficiently, i.e. eventually only one process (the leader) sends messages. In [14], it is shown that Ω can be implemented as long as all correct processes can be reached from a (unknown) correct process with eventually timely paths (paths of eventually timely links and correct processes), even if initially processes only know their own identity. It is also shown there that none of the classes originally proposed in [6] can be implemented if the identity of (at least) one process is unknown to the rest of processes.

In [3], it is shown that Ω can be implemented in a system with fair lossy links, in which at most t processes can crash and t is known, if some correct process has t eventually timely output links. If additionally $n > 2t$, then consensus can be solved. If the links are reliable, they also have an algorithm in which eventually only t links carry messages. They also show that even if all processes have $t - 1$

*A preliminary version of this paper was presented as brief announcement at PODC 2005. Partially supported by the Spanish MEC under grants TIN2005-09198-C02-01, TIN2004-07474-C02-02, and TIN2004-07474-C02-01, and the Comunidad de Madrid under grant S-0505/TIC/0285.

[†]LADyR, GSyC, Universidad Rey Juan Carlos, 28933 Móstoles, Spain

[‡]EUI, Universidad Politécnica de Madrid, 28031 Madrid, Spain

timely output links and all links are reliable neither Ω can be implemented nor consensus can be solved. In [12] it is shown that if initially processes only know their own identity and t , Ω can be implemented if all correct processes are connected via fair lossy paths (paths of fair lossy links and correct processes) and there is one process that can reach $t - f$ (f is the actual number of failures in the run) other correct processes via eventually timely paths. Additionally, if all correct processes are fully connected with fair lossy links and one process has all output links eventually timely, Ω can be implemented efficiently.

Regarding other classes of failure detectors, Larrea et al. [18] consider the implementation of the eight original classes proposed by Chandra and Toueg (namely, the *perpetual* classes \mathcal{P} , \mathcal{S} , \mathcal{Q} , and \mathcal{W} , and the *eventual* classes $\diamond\mathcal{P}$, $\diamond\mathcal{S}$, $\diamond\mathcal{Q}$, and $\diamond\mathcal{W}$) in systems with all links eventually timely. They show that no perpetual detector can be implemented in these systems even if only one process can fail. Then, they show that all eventual detectors can be implemented in these systems. Implicitly, all it is required for this possibility result is that all the links in a ring formed by the correct processes are eventually timely. Only these links carry messages forever. Another related work is [17], in which an algorithm that implements $\diamond\mathcal{S}$ and Ω in a system with all links eventually timely is proposed. With this algorithm, eventually only the output links of one process (the correct process with smallest identifier) carry messages. In fact, the synchrony requirement they impose could be relaxed so that only these links are required to be eventually timely.

A different approach to limit the system requirements to implement failure detectors has been proposed by Mostefaoui et al. [20, 21]. In these works the condition on the system has to do with its behavior. They show that their algorithms implement a $\diamond\mathcal{S}$ detector if messages are received by the processes in the appropriate order. They do a probabilistic analysis and show that these requirements are met with high probability for one single failure.

Contributions. In this paper we continue the exploration of the system limits to implement unreliable failure detectors. We study five of the traditional classes of failure detectors: Ω , \mathcal{P} , \mathcal{S} , $\diamond\mathcal{P}$, and $\diamond\mathcal{S}$, and two additional (perpetual) classes, \mathcal{P}_4 and \mathcal{S}' , which are weak versions (they have weaker accuracy) of \mathcal{P} and \mathcal{S} , respectively. As far as we know, the class \mathcal{S}' has never been previously proposed and seems an interesting class to study further. The class \mathcal{P}_4 was first proposed in [16].

We consider systems formed either by lossy asynchronous and timely links (class Ψ), or lossy asynchronous and eventually timely links (class \mathcal{E}). In the complete directed graph formed by the processes and the links, we call a path that does not contain lossy asynchronous links nor

faulty processes a *good* path. We are interested on the set R of correct processes that can reach all correct processes via good paths. We explore the implementability of failure detectors depending on whether the set R has at least one correct process (*Weak Property*) or R contains all correct processes (*Strong Property*). Additionally, we say that we have the *Min Property* if R contains the correct processor with smallest identifier. Observe that the Strong Property implies the Min Property, and the latter implies the Weak Property.

First, we show that there is no algorithm that implements $\diamond\mathcal{P}$ (and hence \mathcal{P} and \mathcal{P}_4) with single failures in a system that, when all processes are correct, does not satisfy the Strong Property. Similarly, we show that $\diamond\mathcal{S}$ (and hence Ω , \mathcal{S} , and \mathcal{S}') cannot be implemented if we allow one single failure in a system that, when all processes are correct, does not satisfy the Weak Property.

Then, we present algorithms that work on minimal system conditions. First, we propose an algorithm for the systems in Ψ that implements \mathcal{P}_4 if the Strong Property is satisfied and implements \mathcal{S}' if only the Weak Property is satisfied. We propose a second algorithm for the systems in \mathcal{E} that implements $\diamond\mathcal{P}$ if the Strong Property is satisfied and $\diamond\mathcal{S}$ if only the Weak Property is satisfied. If the Min Property is satisfied, both algorithms implement Ω .

For the systems in \mathcal{E} , the Weak Property is the minimal requirement imposed above to implement Ω (with f unknown) [14]. The Min Property is also strictly weaker than the requirements in [17]. However it cannot be compared with those in [2], since they have no requirement regarding the correct process with smallest identifier. Note as well that the Strong Property is weaker than the property for $\diamond\mathcal{P}$ in [18], which requires all links connecting the correct processes in a ring to be eventually timely.

The rest of the paper is structured as follows. In Section 2 we present the model and notations we use in the rest of the paper. In Section 3 we show necessary conditions to implement the classes of failure detectors we consider here. Section 4 presents an algorithm that implements perpetual failure detectors in systems with weak synchrony. Similarly, Section 5 presents an algorithm that implements eventual failure detectors in systems with weak synchrony. Finally, Section 7 contains some concluding remarks.

2. The model

We consider systems formed by a finite set Π of $n > 1$ processes. Processes have unique and totally ordered identifiers, known to all the processes. We assume that processes can communicate with each other only by sending and receiving messages, and that every pair of processes is connected by a pair of directed links (with opposite directions). Let $\Lambda = \{(p, q) : p, q \in \Pi; p \neq q\}$ denote the set of di-

rected links of a system. Clearly, if we see the system as a graph $G = (\Pi, \Lambda)$, G is a complete directed graph.

Failure-prone processes. A process can fail by permanently crashing. We say that a process is correct if it does not fail. We denote by *correct* the set of correct processes. The complementary set $\Pi - \text{correct}$ is denoted by *crashed*. We assume that the algorithms have no a priori knowledge of the number of failures that can occur.

To add generality, for the impossibility results (necessary conditions) of Section 3 we assume that the execution of processes is synchronous, and that their clocks are synchronized. However, we do not need such strong assumptions in our algorithms. In them, we only assume that each line of the algorithm takes no more than σ time to be executed (σ can be a very loose bound). For the *perpetual detectors algorithm* (Section 4) we additionally assume that σ is known to the algorithm. In the algorithms we also assume the availability of timers at each process. We use $\tau_p(T)$ to denote the time it takes a timer of process p started with a value T to expire. We first require that, for all p , $\tau_p(T)$ is finite if T is finite. Then, for the *perpetual detectors algorithm* we also require that for all p and all T , $\tau_p(T) \geq T$, while for the *eventual detectors algorithm* (Section 5) we only require that $\tau_p(T)$ is non-decreasing with T and $\lim_{T \rightarrow \infty} \tau_p(T) = \infty$. Note that processes do not use global clocks.

Unreliable failure detector classes. Chandra and Toueg defined several classes of unreliable failure detectors by specifying their corresponding *completeness* and *accuracy*. These properties are defined on the lists of suspected processes maintained by (the failure detector modules of) the processes. The completeness property requires that every process that actually crashes is eventually suspected, while the accuracy property restricts the mistakes (i.e., the false suspicions) that a failure detector can make. Chandra and Toueg defined two completeness and four accuracy properties in [6]. We only consider here one completeness property:

- *Strong Completeness*: Eventually every process that crashes is permanently suspected by every correct process.

Regarding accuracy, we consider the four properties proposed by Chandra and Toueg:

- *(Perpetual) Strong Accuracy*: No process is suspected before it crashes.
- *(Perpetual) Weak Accuracy*: Some correct process is never suspected.

- *Eventual Strong Accuracy*: There is a time after which no correct process is ever suspected by any correct process.
- *Eventual Weak Accuracy*: There is a time after which some correct process is never suspected by any correct process.

We are going to consider here two more accuracy properties, which are slightly weaker than Strong Accuracy and Weak Accuracy, respectively, as defined above. The first property has been stated in [9, 16]. The second has never been stated as far as we know.

- *(Perpetual) Quasi-Strong Accuracy*. No correct process is ever suspected by any correct process.
- *(Perpetual) Quasi-Weak Accuracy*. Some correct process is never suspected by any correct process.

Failure detectors with eventual accuracy may suspect every process at one time or another, while failure detectors with perpetual accuracy require that at least one correct process is never suspected. Combining the completeness property with one of the accuracy properties we obtain one class of failure detectors. We consider here six different classes, which are presented in Figure 1. Quasi-Strong Accuracy combined with Strong Completeness yields a weak version of the class \mathcal{P} , which has been denoted in [16] by \mathcal{P}_4 . Quasi-Weak Accuracy combined with Strong Completeness yields a weak version of the class \mathcal{S} , which we are going to denote here by \mathcal{S}' .

Finally, we define the Ω failure detector. In an Ω failure detector, each process chooses some process in the system as *leader*. The detector must guarantee that all correct processes eventually agree on a single correct process as leader. More formally, Ω failure detectors must satisfy the following property.

Property 1 *There is a time after which every process $p \in \text{correct}$ permanently has the same process $l \in \text{correct}$ as leader.*

We will consider in this paper that an Ω failure detector also implements $\diamond\mathcal{S}$, because the set $\Pi \setminus \{\text{leader}\}$ guarantees the properties of $\diamond\mathcal{S}$.

Types of links. We consider the following three types of links¹.

- Lossy asynchronous (LA)*: A message sent across a lossy asynchronous link can be lost or arbitrarily delayed; however, a message that is not lost will eventually be delivered. (Note that all messages sent using a lossy asynchronous link may be lost.)

¹Observe that the type of a link depends on its observed behavior. This behavior can be satisfied in all or in a specific set of runs of the system.

Completeness	Accuracy					
	Strong	Weak	Quasi-Strong	Quasi-Weak	Eventual Strong	Eventual Weak
Strong	\mathcal{P}	\mathcal{S}	\mathcal{P}_4	\mathcal{S}'	$\diamond\mathcal{P}$	$\diamond\mathcal{S}$

Figure 1. Classes of failure detectors we consider, defined in terms of completeness and accuracy.

- (ii) *Timely (T)*: The link is reliable and there is a *known* bound Δ on the maximum message delay. (Hence, a message that is sent at time t is received by time $t+\Delta$.)
- (iii) *Eventually timely (ET)*: There is a *possibly unknown* global stabilization time GST such that until GST the link behaves like lossy asynchronous; after GST the link is reliable and there is a *possibly unknown* bound Δ on the maximum message delay. (Hence, a message that is sent at time $t > GST$ is received by time $t+\Delta$.)

We assume that links do not modify the messages they carry nor they generate spontaneous messages. In order to make the negative results stronger, we assume for the impossibility results that links do not replicate messages and that they deliver them in order. However, in our algorithms we allow messages to be replicated and received out of order. In fact, our algorithms explicitly resend messages, creating replicated messages. We finally assume that messages are unique, in the sense that an algorithm can determine whether a message received is a duplicate (either generated by the link or resent by some other process) of a previously received message. This can be easily implemented using a message identifier formed by the unique identifier of the sending process and a sequence number, unique for that process.

Classes of systems. In this paper we consider two large classes of systems. A system belongs to the class Ψ if each of its links² is either lossy asynchronous or timely. A system belongs to the class \mathcal{E} if each of its links is either lossy asynchronous or eventually timely. Since timely links are special cases of eventually timely links (with $GST = 0$ and known Δ), we have that $\Psi \subset \mathcal{E}$. From now on, when we assume a system in \mathcal{E} we also include those in Ψ .

Then, we consider a system $S \in \mathcal{E}$ characterized by the pair $(L_S, correct_S)$, where $L_S : \Lambda \rightarrow \{LA, T, ET\}$, and $correct_S \subseteq \Pi^3$. We denote this by $S = (L_S, correct_S)$. L_S is a function that determines for each link in S its type. Of course L_S must be consistent with the class to which S belongs. Then if $S \in \Psi$, L_S can take the values LA and T only. The set $correct_S$ is the set of correct processes in S . When convenient we will use the complementary set

²I.e., the behavior of its links in each run.

³Intuitively, a system as defined is the set of all runs in which all processes in $correct_S$ are correct and the links behave as implied by L_S .

$crashed_S$, and when clear from the context we will remove the subindex.

For any system $S \in \mathcal{E}$ we are going to define an associated graph $G(S)$ which is derived from the attributes of S . We mentioned above that S can be seen as the complete directed graph (Π, Λ) . Then, $G(S)$ is the directed subgraph induced in this graph by the set $correct_S$, from which all the lossy asynchronous links (links with $L_S(p, q) = LA$) have been removed. Then, $G(S)$ only contains correct processes as vertices and directed timely links (if $S \in \Psi$) or directed eventually timely links (if $S \in \mathcal{E}$).

Given two correct processes p and q in S , we say that q can be reached from p if either $p = q$ or there is a directed path from p to q in $G(S)$. The set of processes that can be reached from a process p is denoted by $reach(p)$. It can be trivially observed that

Observation 1 *If $q \in reach(p)$, then $reach(q) \subseteq reach(p)$.*

We define now the following properties that can be satisfied by a system $S \in \mathcal{E}$.

Property 2 (Weak) *There is some process $p \in correct$ such that $reach(p) = correct$ in $G(S)$.*

Property 3 (Min) $reach(\min(correct)) = correct$ in $G(S)$.

Property 4 (Strong) *For all process $p \in correct$, $reach(p) = correct$ in $G(S)$.*

Observe that the Strong Property implies the Min Property, and the Min Property implies the Weak Property.

Finally, we will use the following notation. Given a system $S = (L_S, correct_S) \in \mathcal{E}$, we denote by $S(p)$, for $p \in correct_S$, the system obtained from S by removing p from the set $correct_S$. That is, $S(p) = (L_S, correct_S \setminus \{p\})$. Then, we denote by $\Phi(S)$ the set that contains S and all systems $S(p)$, $p \in correct_S$, i.e. $\Phi(S) = \{S\} \cup \{S(p) : p \in correct_S\}$.

Algorithms. We study here algorithms that implement unreliable failure detectors of the above classes. These algorithms are implemented as one local module for each process of the system. A module exists as long as its local process has not crashed. Modules exchange messages among

each other to provide the required completeness and accuracy properties. They also, upon request by their local process, provide it with the current list of suspected processes or the current leader.

An algorithm \mathcal{A} implements a failure detector of a given class C for a set of systems Σ . If it is claimed that \mathcal{A} implements the detector for a system set Σ , and $S \in \Sigma$, then every run of \mathcal{A} in S must guarantee that the implemented detector belongs to C , independently of when the processes in $crashed_S$ fail and the behavior of the links (as long as they are consistent with their type). The sets of systems we consider here are subsets of the class \mathcal{E} . We will make clear the subset of systems for which an algorithm implements the detector in each case.

3. Necessary conditions to implement failure detectors

In this section we obtain minimal conditions that systems in \mathcal{E} must satisfy to be able to implement a failure detector of the classes of interest.

3.1. Conditions for detectors with weak accuracy

We consider first the detector classes $\diamond\mathcal{S}$, Ω , \mathcal{S}' , and \mathcal{S} . We show that it is not possible to implement a $\diamond\mathcal{S}$ detector (and hence Ω , \mathcal{S} , and \mathcal{S}' detectors) for a set of systems that contains one system S that does not satisfy the Weak Property and those obtained from it with one more failure (denoted $S(p)$). Recall that we denote this set by $\Phi(S)$. The following theorem shows this.

Theorem 1 *Let $S \in \mathcal{E}$ be a system that does not satisfy the Weak Property, $S(p)$ be the system obtained from S by removing process $p \in correct_S$ from the set of correct processes, and $\Phi(S) = \{S\} \cup \{S(p) : p \in correct_S\}$. There is no algorithm that implements a $\diamond\mathcal{S}$ failure detector for the set of systems $\Phi(S)$.*

Proof: For the sake of contradiction, let us assume there is such an algorithm \mathcal{A} . Let us consider a run R of \mathcal{A} in S in which all the messages sent across lossy asynchronous links are lost. By Eventual Weak Accuracy, there will be a process $p \in correct_S$ and a time t such that p is never suspected by any correct process after t . Now, note that the set of processes $Q = correct_S \setminus reach(p)$ is not empty, since the Weak Property does not hold. Furthermore, from Observation 1, no message from the processes in $reach(p)$ ever reaches any process in Q .

Let us consider now system $S(p)$. By assumption, \mathcal{A} should also implement a $\diamond\mathcal{S}$ detector in $S(p)$. Let us consider a run R' of \mathcal{A} in $S(p)$ in which all the messages sent across lossy asynchronous links are lost and all

processes behave as closely as possible to their behavior in R . Note that the processes in Q do not have any way of distinguishing R' from R , since like before they receive no message from $reach(p)$, which are the processes that noticed the failure of p . Hence, they can in fact behave exactly like in R , and they will never suspect p after time t . This violates the Strong Completeness property, and hence \mathcal{A} cannot exist. ■

Clearly, if there is no algorithm for a given set, there is no algorithm for any of its supersets. Also, since all detectors in Ω , \mathcal{S} , and \mathcal{S}' are also in $\diamond\mathcal{S}$, we have the following corollary.

Corollary 1 *Let $S \in \mathcal{E}$ be a system that does not satisfy the Weak Property, and $\Sigma \supseteq \Phi(S)$ be a set of systems. There is no algorithm that implements a $\diamond\mathcal{S}$, Ω , \mathcal{S}' , or \mathcal{S} failure detector for the set of systems Σ .*

Note that this holds even if S has no failures.

Corollary 2 *Let $S \in \mathcal{E}$ be a system without failures that does not satisfy the Weak Property, and Σ be a set of systems that include S and the systems obtained from S with one single failure (i.e., $\Sigma \supseteq \Phi(S)$). There is no algorithm that implements a $\diamond\mathcal{S}$, Ω , \mathcal{S}' , or \mathcal{S} failure detector for the set of systems Σ .*

3.2. Conditions for detectors with strong accuracy

Let us now look at detectors with some form of strong accuracy. Like before, we will show first that there is no algorithm to implement a detector in $\diamond\mathcal{P}$ for a set of systems that contains a system that does not satisfy the Strong Property and those obtained from it with one additional failure. The following theorem shows this.

Theorem 2 *Let $S \in \mathcal{E}$ be a system that does not satisfy the Strong Property, $S(p)$ be the system obtained from S by removing process $p \in correct_S$ from the set of correct processes, and $\Phi(S) = \{S\} \cup \{S(p) : p \in correct_S\}$. There is no algorithm that implements a $\diamond\mathcal{P}$ failure detector for the set of systems $\Phi(S)$.*

Proof: For the sake of contradiction, let us assume there is such an algorithm \mathcal{A} . Let us consider a run R of \mathcal{A} in S in which all the messages sent across lossy asynchronous links are lost. By Eventual Strong Accuracy, there will be a time t such that no correct process is ever suspected by any correct process after t . Since S does not satisfy the Strong Property, there is some process $p \in correct_S$ such that $reach(p) \neq correct_S$. Then, the set of processes $Q = correct_S \setminus reach(p)$ is not empty (note that $reach(p)$

only contains correct processes). Furthermore, from Observation 1, no message from the processes in $reach(p)$ ever reaches any process in Q .

Let us consider now system $S(p)$. By assumption, \mathcal{A} should also implement a $\diamond\mathcal{P}$ detector in $S(p)$. Let us consider a run R' of \mathcal{A} in $S(p)$ in which all the messages sent across lossy asynchronous links are lost and all processes behave as closely as possible to their behavior in R . Note that the processes in Q do not have any way of distinguishing R' from R , since like before they receive no message from $reach(p)$, which are the processes that noticed the failure of p . Hence, they can in fact behave exactly like in R , and they will never suspect p after time t . This violates the Strong Completeness property, and hence \mathcal{A} cannot exist. ■

Again, if there is no algorithm for a given set, there is no algorithm for any of its supersets. Also, since the detectors in \mathcal{P} and \mathcal{P}_4 are also in $\diamond\mathcal{P}$, we have the following corollary.

Corollary 3 *Let $S \in \mathcal{E}$ be a system that does not satisfy the Strong Property, and $\Sigma \supseteq \Phi(S)$ be a set of systems. There is no algorithm that implements a $\diamond\mathcal{P}$, \mathcal{P}_4 , or \mathcal{P} failure detector for the set of systems Σ .*

Corollary 4 *Let $S \in \mathcal{E}$ be a system without failures that does not satisfy the Strong Property, and Σ be a set of systems that include S and the systems obtained from S with one single failure (i.e., $\Sigma \supseteq \Phi(S)$). There is no algorithm that implements a $\diamond\mathcal{P}$, \mathcal{P}_4 , or \mathcal{P} failure detector for the set of systems Σ .*

4. Algorithm for perpetual failure detectors

In this section we present an algorithm that implements a failure detector for systems in the class Ψ . For all the systems in Ψ that satisfy the Weak Property the algorithm implements a detector of class \mathcal{S}' . If additionally they satisfy the Min Property the algorithm implements a detector of class Ω . For all the systems that satisfy the Strong Property the algorithm implements a detector of class \mathcal{P}_4 . Figure 2 presents the algorithm in detail. For all $p \in \Pi$, the sets $suspected_p$ provide the required completeness and accuracy properties for \mathcal{S}' and \mathcal{P}_4 , while the values $leader_p$ satisfy Property 1 for Ω .

The proof of the following theorem has been removed due to space limitations. It can be found at [11].

Theorem 3 *Let $S \in \Psi$ be a system in which the Algorithm Perpetual (Figure 2) is executed. Then,*
(i) *if S satisfies the Weak Property, Algorithm Perpetual implements a failure detector of class \mathcal{S}' ,*
(ii) *if, additionally, S satisfies the Min Property, Algorithm*

Algorithm Perpetual

init:

- (1) $suspected_p \leftarrow \emptyset$
- (2) $leader_p \leftarrow \min(\Pi)$
- (3) $Timeout_p \leftarrow \eta + (n - 1)(\Delta + 4\sigma)$
- (4) reset $timer_p(q)$ to $Timeout_p$, for each process $q \neq p$
- (5) **start tasks 1 and 2**

Task 1:

- (6) **loop forever**
- (7) send ($ALIVE, p$) to every process except p every η time

Task 2:

- (8) **upon reception of message ($ALIVE, q$) do**
- (9) **if** [message ($ALIVE, q$) was not previously received] **then**
- (10) reset $timer_p(q)$ to $Timeout_p$
- (11) resend message ($ALIVE, q$) to every process except p
- (12) **upon expiration of $timer_p(q)$ do**
- (13) $suspected_p \leftarrow suspected_p \cup \{q\}$
- (14) $leader_p \leftarrow \min(\Pi \setminus suspected_p)$

Figure 2. Algorithm to implement perpetual failure detectors in systems of class Ψ . The code is for process p .

Perpetual implements a failure detector of class Ω , and (ii) if, additionally, S satisfies the Strong Property, Algorithm Perpetual implements a failure detector of class \mathcal{P}_4 .

5. Algorithm for eventual failure detectors

In this section we present an algorithm that implements failure detectors for systems in class \mathcal{E} . We show that, for all systems in \mathcal{E} that satisfy the Weak Property the algorithm implements $\diamond\mathcal{S}$, if they satisfy the Min Property it implements Ω , and if they satisfy the Strong Property it implements $\diamond\mathcal{P}$. Figure 3 presents the algorithm in detail.

We have removed the proof of the following theorem due to space limitation. It can be found at [11].

Theorem 4 *Let $S \in \mathcal{E}$ be a system in which the Algorithm Eventual (Figure 3) is executed. Then,*

- (i) *if S satisfies the Weak Property, Algorithm Eventual implements a failure detector of class $\diamond\mathcal{S}$,*
- (ii) *if, additionally, S satisfies the Min Property, Algorithm Eventual implements a failure detector of class Ω , and*
- (ii) *if, additionally, S satisfies the Strong Property, Algorithm Eventual implements a failure detector of class $\diamond\mathcal{P}$.*

6. The Failure Detector Class \mathcal{S}'

In this section we explore the new class \mathcal{S}' of unreliable failure detectors proposed in this paper. We first show that this class is strictly weaker than the class \mathcal{S} . To do so, we show that while any detector in \mathcal{S} can be used to solve uniform consensus, the same is not so for any detector in \mathcal{S}' .

Algorithm Eventual**init:**

- (1) $suspected_p \leftarrow \emptyset$
- (2) $leader_p \leftarrow \min(\Pi)$
- (3) $Timeout_p(q) \leftarrow \eta + 1$, for each process $q \neq p$
- (4) reset $timer_p(q)$ to $Timeout_p(q)$, for each process $q \neq p$
- (5) **start tasks 1 and 2**

Task 1:

- (6) **loop forever**
- (7) send $(ALIVE, p)$ to every process except p every η time

Task 2:

- (8) **upon reception of message** $(ALIVE, q)$ **do**
- (9) **if** [message $(ALIVE, q)$ was not previously received] **then**
- (10) $suspected_p \leftarrow suspected_p \setminus \{q\}$
- (11) $leader_p \leftarrow \min(\Pi \setminus suspected_p)$
- (12) reset $timer_p(q)$ to $Timeout_p(q)$
- (13) resend message $(ALIVE, q)$ to every process except p
- (14) **upon expiration of** $timer_p(q)$ **do**
- (15) $Timeout_p(q) \leftarrow Timeout_p(q) + 1$
- (16) $suspected_p \leftarrow suspected_p \cup \{q\}$
- (17) $leader_p \leftarrow \min(\Pi \setminus suspected_p)$

Figure 3. Algorithm to implement eventual failure detectors in systems of class \mathcal{E} . The code is for process p .

Then, we show that any detector in \mathcal{S}' can be used to solve non-uniform consensus in an asynchronous system with any number of failures. We show this by showing how to transform a failure detector in \mathcal{S}' into a failure detector of class (Ω, Σ') , which can be used to solve nonuniform consensus [10].

Theorem 5 *It is not enough to have a failure detector of class \mathcal{S}' to solve uniform consensus in a crash-prone asynchronous system with single crashes.*

Proof: (Sketch) By way of contradiction, let \mathcal{A} be an algorithm that solves uniform consensus with any detector in \mathcal{S}' . Consider a system with two processes p and q and a detector $D \in \mathcal{S}'$. Let R_0 be a run of \mathcal{A} in which both processes propose 0 and p fails before sending any message. D always returns $\{p\}$ to q as the set of suspected processes. Then, q must decide 0 at some time t_0 . Similarly, let run R_1 be a run of \mathcal{A} in which both processes propose 1 and q fails before sending any message; D always returns $\{q\}$ to p as the list of suspected processes; and p decides 1 at time t_1 . Finally, consider a run R of \mathcal{A} in which all messages sent are delayed until a time $t > \max(t_0, t_1)$, and p fails at this time. D always returns $\{p\}$ to q and $\{q\}$ to p as their respective lists of suspected processes. Until time t process p behaves exactly like in R_1 , and hence decides 1 at time t_1 . Similarly, process q behaves exactly like in R_0 until time t , and hence decides 0 at time t_0 . Since they decide different values, \mathcal{A} does not solve uniform consensus. ■

Now we show that any failure detector $D \in \mathcal{S}'$ can be transformed into a failure detector of class (Ω, Σ') . A failure detector of class Ω can be obtained from D by using, for instance, the algorithm proposed by Chu [8] for $\diamond\mathcal{W}^4$. A detector of class Σ' is trivially obtained from D by returning as quorum at each process the complement of the list of suspected processes returned by D . Clearly, all the quorums returned at the correct processes contain at least the correct process that is never suspected by them. Eventually these quorums only contain correct processes by the strong completeness of \mathcal{S}' .

7. Conclusions

In this paper we explore the minimal system synchrony to implement unreliable failure detectors. We present algorithms that implement detectors in systems with weak synchrony requirements and show that these requirements are in fact needed.

There are still a number of open problems related with this work. For instance, we present algorithms that implement detectors of classes \mathcal{P}_4 and \mathcal{S}' , which are weaker than \mathcal{P} and \mathcal{S} , in systems with limited synchrony. It would be nice to have algorithms that implement \mathcal{P} and \mathcal{S} in systems with the same synchrony.

Finally, observe that our algorithms have a quadratic number of links carrying messages forever in the worse case. We believe this is the best we can hope for $\diamond\mathcal{S}$ (since there are similar bounds for Ω in a system with one correct process whose output links are eventually timely [2]). However, we would like to have a proof of that.

References

- [1] M. Aguilera, C. Delporte-Gallet, H. Fauconnier, and S. Toueg. Stable leader election. In *Proceedings of the 15th International Symposium on Distributed Computing (DISC'2001)*, pages 108–122, Lisbon, Portugal, October 2001. LNCS 2180, Springer-Verlag.
- [2] M. Aguilera, C. Delporte-Gallet, H. Fauconnier, and S. Toueg. On implementing Ω with weak reliability and synchrony assumptions. In *Proceedings of the 22nd Annual Symposium on Principles of Distributed Computing (PODC'2003)*, Boston, Massachusetts, July 2003.
- [3] M. Aguilera, C. Delporte-Gallet, H. Fauconnier, and S. Toueg. Communication-efficient leader election and consensus with limited link synchrony. In *Proceedings of the 23rd Annual Symposium on Principles of Distributed Computing (PODC'2004)*, 2004.
- [4] M. Bertier, O. Marin, and P. Sens. Implementation and performance evaluation of an adaptable failure detector. In *Proceedings of the 2002 International Conference on Dependable Systems and Networks*, June 2002.

⁴Any detector in \mathcal{S}' is trivially in $\diamond\mathcal{W}$.

- [5] T. D. Chandra, V. Hadzilacos, and S. Toueg. The weakest failure detector for solving consensus. *Journal of the ACM*, 43(4):685–722, July 1996.
- [6] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225–267, March 1996.
- [7] W. Chen, S. Toueg, and M. K. Aguilera. On the quality of service of failure detectors. *IEEE Transactions on Computers*, 51(1):13–32, January 2002.
- [8] F. Chu. Reducing Ω to $\diamond W$. *Information Processing Letters*, 67:289–293, 1998.
- [9] X. Défago, A. Schiper, and P. Urbán. Total order broadcast and multicast algorithms: Taxonomy and survey. *ACM Comput. Surv.*, 36(4):372–421, 2004.
- [10] J. Eisler, V. Hadzilacos, and S. Toueg. The weakest failure detector to solve nonuniform consensus. In *Proceedings of the 24th Annual Symposium on Principles of Distributed Computing (PODC'2005)*, pages 189–196, Las Vegas, NV, July 2005.
- [11] A. Fernández, E. Jiménez, and S. Arévalo. Minimal system conditions to implement unreliable failure detectors. Technical report, Reports on Systems and Communications, Grupo de Sistemas y Comunicaciones, 2005.
- [12] A. Fernández, E. Jiménez, and M. Raynal. Eventual leader election with weak assumptions on initial knowledge, communication reliability, and synchrony. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN-2006)*, Philadelphia, PA, USA, June 2006.
- [13] M. Fischer, N. Lynch, and M. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, April 1985.
- [14] E. Jiménez, S. Arévalo, and A. Fernández. Implementing unreliable failure detectors with unknown membership. *Information Processing Letters*, 100(2):60–63, 2006.
- [15] I. Keidar and S. Rajsbaum. On the cost of fault-tolerant consensus when there are no faults. In *Tutorial on the 21st Annual Symposium on Principles of Distributed Computing (PODC'2002)*, Monterey, California, July 2002.
- [16] M. Larrea. Brief announcement: Understanding perfect failure detectors. In *Proceedings of the 21st Annual ACM Symposium on Principles of Distributed Computing, PODC 2002*, page 257, 2002.
- [17] M. Larrea, A. Fernández, and S. Arévalo. Optimal implementation of the weakest failure detector for solving consensus. In *Proceedings of the 19th IEEE Symposium on Reliable Distributed Systems (SRDS'2000)*, pages 52–59, Nuremberg, Germany, October 2000.
- [18] M. Larrea, A. Fernández, and S. Arévalo. On the implementation of unreliable failure detectors in partially synchronous systems. *IEEE Transactions on Computers*, 53(7):815–828, July 2004.
- [19] M. Larrea, A. Fernández, and S. Arévalo. Eventually consistent failure detectors. *Journal of Parallel and Distributed Computing*, 65(3):361–373, March 2005.
- [20] A. Mostéfaoui, E. Mourgaya, and M. Raynal. Asynchronous implementation of failure detectors. In *2003 International Conference on Dependable Systems and Networks (DSN 2003)*, pages 351–360, 2003.
- [21] A. Mostéfaoui, D. Powell, and M. Raynal. A hybrid approach for building eventually accurate failure detectors. In *10th IEEE Pacific Rim International Symposium on Dependable Computing (PRDC 2004)*, pages 57–65, 2004.