

On the Computability Power and the Robustness of Set Agreement-oriented Failure Detector Classes

Achour Mostefaoui, Sergio Rajsbaum, Michel Raynal, Corentin Travers

► **To cite this version:**

Achour Mostefaoui, Sergio Rajsbaum, Michel Raynal, Corentin Travers. On the Computability Power and the Robustness of Set Agreement-oriented Failure Detector Classes. [Research Report] PI 1819, 2006, pp.31. <inria-00107220>

HAL Id: inria-00107220

<https://hal.inria.fr/inria-00107220>

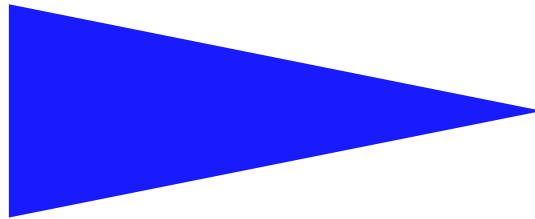
Submitted on 17 Oct 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

IRISA
INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

PUBLICATION
INTERNE
N° 1819



ON THE COMPUTABILITY POWER AND THE
ROBUSTNESS OF SET AGREEMENT-ORIENTED
FAILURE DETECTOR CLASSES

A. MOSTEFAUI S. RAJSBAUM M. RAYNAL C. TRAVERS

On the Computability Power and the Robustness of Set Agreement-oriented Failure Detector Classes*

A. Mostefaui** S. Rajsbaum*** M. Raynal**** C. Travers*****

Systèmes communicants

Publication interne n° 1819 — Octobre 2006 — 31 pages

Abstract: Solving agreement problems, such as consensus and k -set agreement, in asynchronous distributed systems prone to process failures has been shown to be impossible. To circumvent this impossibility, unreliable failure detectors have been widely studied. These are oracles that provide information on failures. The exact nature of such information is defined by a set of abstract properties that a particular class of failure detectors satisfy. The weakest class of failure detectors that allow to solve consensus is Ω .

This paper considers the failure detector classes that have been considered in the literature to solve k -set agreement, and studies their relative power. It shows that the family of failure detector classes $\diamond\mathcal{S}_x$ ($0 \leq x \leq n$), and $\diamond\psi^y$ ($0 \leq y \leq n$), can be “added” to provide a failure detector of the class Ω^z (a generalization of Ω). It also characterizes the power of such an “addition”, namely, $\diamond\mathcal{S}_x + \diamond\psi^y \rightsquigarrow \Omega^z \Leftrightarrow x + y + z > t + 1$, where t is the maximum number of processes that can crash in a run. As an example, the paper shows that, while $\diamond\mathcal{S}_t$ allows solving 2-set agreement (and not consensus) and $\diamond\psi^1$ allows solving t -set agreement (but not $(t - 1)$ -set agreement), a system with failure detectors of both classes can solve consensus. More generally, the paper studies the failure detector classes $\diamond\mathcal{S}_x$, $\diamond\psi^y$ and Ω^z , and shows which reductions among these classes are possible and which are not.

The paper presents also a message-passing Ω^k -based k -set agreement protocol. In that sense, it can be seen as a step toward the characterization of the weakest failure detector class that allows solving the k -set agreement problem.

Key-words: Asynchronous system, Distributed algorithm, Eventual leader, Fault-tolerance, Limited scope accuracy, Process crash, Message-passing system, Reduction algorithm, Robustness, Scalability, Set agreement, Unreliable failure detector.

(Résumé : *tsvp*)

* An extended abstract of this paper has appeared in the proceedings of PODC 2006 [21]. This work has been supported by a grant from LAFMI (Franco-Mexican Lab in Computer Science) and PAPIIT-UNAM.

** IRISA, Université de Rennes 1, Campus de Beaulieu, 35042 Rennes Cedex, France achour@irisa.fr

*** Instituto de Matemáticas, UNAM, D. F. 04510, Mexico rajsbaum@matem.unam.mx

**** IRISA, Université de Rennes 1, Campus de Beaulieu, 35042 Rennes Cedex, France, raynal@irisa.fr

***** IRISA, Université de Rennes 1, Campus de Beaulieu, 35042 Rennes Cedex, France, travers@irisa.fr



Sur la puissance et la robustesse de classes de détecteurs de fautes

Résumé : Ce rapport étudie la puissance de calcul et la robustesse de classes de détecteur de fautes.

Mots clés : Systèmes répartis asynchrones, Tolérance aux fautes, Crash de processus, Détecteur de fautes.

1 Introduction

Context of the work: failure detectors for agreement problems *Consensus* is one of the most fundamental problems in fault-tolerant distributed computing: each process proposes a value, and every non-faulty process must decide a value (termination) such that no two different values are decided (agreement) and the decided value is a proposed value (validity). Despite the simplicity of its definition and its use as a basic building block to solve distributed agreement problems, consensus cannot be solved in asynchronous systems where even a single process can crash [9].

Several approaches have been investigated to circumvent this impossibility result. One of them is the failure detector approach [4, 27]. It consists in equipping the underlying system with a distributed oracle that provides each process with (possibly inaccurate) hints on process failures. According to the type and the quality of the hints, several classes of failure detectors can be defined. As far as consensus is concerned, two classes are particularly important.

- The class of *leader* failure detectors [3], denoted Ω . This class includes all the failure detectors that continuously output at each process the identity of a process such that, after some time, all the correct processes are provided with the same identity that is the identity of a correct process (eventual leadership). Before that time, different processes can be provided with distinct leaders (that can also change over time), and there is no way for the processes to know when this anarchy period is over. Ω -based asynchronous consensus protocols can be found in [11, 17, 25]¹.

- The class of *eventually strong* failure detectors [4], denoted $\diamond\mathcal{S}$. A failure detector of that class provides each process with a set of suspected processes such that this set eventually includes all the crashed processes (strong completeness) and there is a correct process p and a time after which no set contains the identity of p (eventual strong accuracy). $\diamond\mathcal{S}$ -based asynchronous consensus protocols can be found in [4, 11, 22, 29].

Two important results are associated with Ω and $\diamond\mathcal{S}$. First, they are equivalent (which means that it is possible, from any failure detector of any of these classes, to build a failure detector of the other class) [3, 6, 20]. Second, as far as information on failures is concerned, they are the weakest classes of failure detectors that allow solving consensus in asynchronous systems where a majority of processes are correct [3].

The *k-set agreement* problem relaxes the consensus requirement to allow up to k different values to be decided [5] (consensus is 1-set agreement). This problem is solvable in asynchronous system despite up to $k - 1$ process crash failures, but has been shown to be impossible to solve as soon as k or more processes can crash [1, 16, 28].

A weakened form of the failure detector class $\diamond\mathcal{S}$ has been first proposed in [12] and investigated to solve consensus in [23]. It has then been considered in [24, 30] with the k -set agreement problem in mind. While the scope of the accuracy property of $\diamond\mathcal{S}$ spans the whole system (there is a correct process that, after some time, is not suspected by any process), the class $\diamond\mathcal{S}_x$ is defined by the same completeness property and a limited scope accuracy property, namely, there is a correct process that, after some time, is not suspected by x processes. It is easy to see that $\diamond\mathcal{S}_n$ (where n is the total number of processes) is $\diamond\mathcal{S}$, while $\diamond\mathcal{S}_1$ provides no information on failures. Moreover, $\diamond\mathcal{S}_{x+1} \subseteq \diamond\mathcal{S}_x$. It has been shown that, when we consider the family $(\diamond\mathcal{S}_x)_{1 \leq x \leq n}$ of failure detectors, $\diamond\mathcal{S}_x$ is the weakest class that allows solving k -set agreement in asynchronous systems for $k = t - x + 2$ (where t is an upper bound on the number of processes that can crash) [14] (message-passing systems must also satisfy the additional constraint of a majority of correct processes, $t < n/2$). The class \mathcal{S}_x of failure detectors is a subset of $\diamond\mathcal{S}_x$. It has the same completeness property but a stronger accuracy property: it requires that, from the very beginning, there is a subset of x processes that never suspect one correct process.

A family of failure detectors, denoted $(\phi^y)_{0 \leq y \leq n}$, has recently been introduced in [19] where it is used in conjunction with conditions [18] to solve set agreement problems². A failure detector of the class ϕ^y provides the processes with a query primitive that has as parameter a set X of processes, and returns a boolean answer. When $|X|$ is too small or too big, the invocation $\text{QUERY}(X)$ by a process returns systematically *true*

¹It is important to notice that the first version of the leader-based Paxos protocol dates back to 1989, i.e., before the Ω formalism was introduced.

²A *condition* is a restriction on the possible inputs to a distributed problem. When a distributed problem is not solvable in a given system, conditions that allow to solve it are considered.

(resp., *false*). Otherwise, namely, when $t - y < |X| \leq t$, $\text{QUERY}(X)$ returns *true* only if all the processes in X have crashed; moreover, if all the processes of X have crashed and a process repeatedly issues $\text{QUERY}(X)$, it eventually obtains the answer *true*. We have $\phi^{y+1} \subseteq \phi^y$. Moreover, ϕ^0 provides no information on failures, while, $\forall y \geq t$, ϕ^y is equivalent to a perfect failure detector (one that never does a mistake [4]). The class $\diamond\phi^y$ has been introduced in [21]. A failure detector of that class eventually satisfies the properties defining the class ϕ^y . It is shown in [21] that, when we consider the family $(\diamond\phi^y)_{0 \leq y \leq t}$, $\diamond\phi^y$ is the weakest class for solving the asynchronous k -set agreement problem where $k = t - y + 1$.

The family of failure detector classes $(\Omega^z)_{1 \leq z \leq n}$ [26] has been introduced to augment the synchronization power of object types in the wait-free hierarchy. A failure detector of the class Ω^z outputs at each process a set of at most z process identities such that, after some time, the same set including the identity of at least one correct process is output at all correct processes. Clearly, Ω^1 is Ω . Moreover, $\Omega^z \subseteq \Omega^{z+1}$.

Motivation and results Given that we know of three families of failure detectors $(\diamond\mathcal{S}_x)_{1 \leq x \leq n}$, $(\diamond\phi^y)_{0 \leq y < n}$, and $(\Omega^z)_{1 \leq z \leq n}$, we are interested in studying their relative power. We have that k -set agreement can be solved with

- $\diamond\mathcal{S}_x$, $k = t - x + 2$,
- $\diamond\phi^y$, $k = t - y + 1$, and
- Ω^z , $k = z$ as we show in this paper.

Thus, natural questions are the following:

Are the classes $\diamond\mathcal{S}_x$, $\diamond\phi^y$ and Ω^z that solve k -set agreement, equivalent?

Is the hierarchy represented by these three families of failure detectors, robust, or is it possible to use two of them that cannot solve k -set agreement and together solve it?

If so, which failure detector class do they produce? Etc.

In their seminal work on failure detectors, Chandra, Hadzilacos and Toueg [3, 4] define the output of a failure detector query according to the failure pattern of the corresponding run and the invocation time of that query. Differently, the output of a query of ϕ^y or $\diamond\phi^y$ depends also on a parameter provided by the invoking process (the set of processes that the invoking process inquiries about). In that sense, the definition of this family $(\diamond\phi^y)_{0 \leq y \leq t}$ does not fit the Chandra and Toueg's failure detector definition framework [4]. We start with the following.

- **CONTRIBUTION #1:** The classes $(\psi^y)_{0 \leq y \leq n}$ and $(\diamond\psi^y)_{0 \leq y \leq n}$.

The paper introduces two new classes of failure detectors (denoted ψ^y and $\diamond\psi^y$) that are defined in the Chandra and Toueg's failure detector framework [4], i.e., the output of a failure detector query depends only on the failure pattern and the time at which the failure detector is queried. These classes are rather "natural" as they output an integer that approximates the number of crashed processes.

More precisely, a query to a failure detector of the class ψ^y returns an integer that is always comprised between $t - y$ and the number of processes that crash during the run. Let f^τ be the number of processes that have crashed at time τ . For any τ there is a time $\tau' \geq \tau$ from which the outputs returned by the queries issued after τ' are $\geq f^\tau$. The class $\diamond\psi^y$ allows the properties defining ψ^y to be satisfied only eventually which means that during an arbitrary (but finite) period, the integers returned by the queries can be arbitrary.

A first result of the paper shows that the classes ψ^y and $\diamond\psi^y$ are equivalent to ϕ^y and $\diamond\phi^y$, respectively. "Equivalent" means that, given any failure detector of one class (e.g., $\diamond\phi^y$), it is possible to build a failure detector of the other class (e.g., $\diamond\psi^y$); both provide the same information on failures.

In addition to the previous one, the paper has the three following contributions. In the following, the notation $A + B \rightsquigarrow C$ means that, given as inputs a failure detector of the class A and a failure detector of the class B , there is an algorithm that constructs a failure detector of the class C . The notation $A + B \not\rightsquigarrow C$ means that there is no such transformation algorithm. The notations $A \rightsquigarrow C$ and $A \not\rightsquigarrow C$ have the same meaning considering a single failure detector class as input.

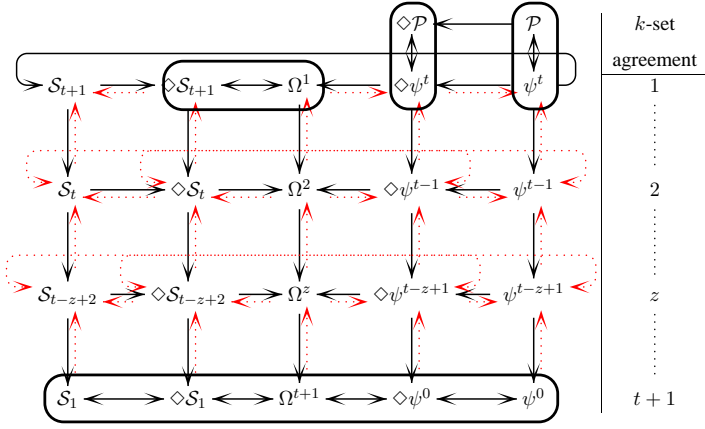


Figure 1: Grid of failure detector classes

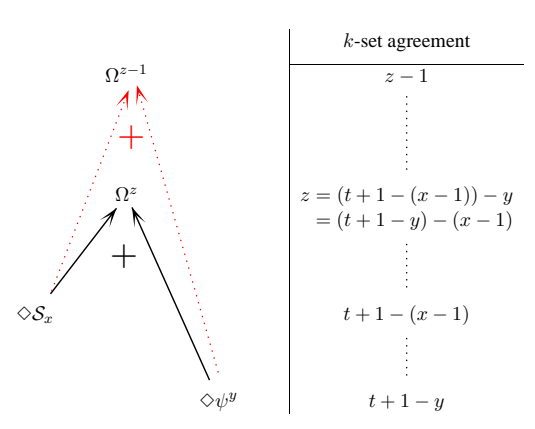


Figure 2: Additivity of $\diamond S_x$ and $\diamond \psi^y$

• CONTRIBUTION #2: Reducibility, Irreducibility and Minimality.

- Relations linking $\psi^y / \diamond \psi^y$ and $S_x / \diamond S_x$:
 - Let $1 \leq x \leq t+1$ and $1 \leq y \leq t$. $S_x \not\rightsquigarrow \diamond \psi^y$. (Theorem 10.)
 - Let $0 \leq y < t$ and $1 < x \leq t+1$. $\psi^y \not\rightsquigarrow \diamond S_x$. (Theorem 11.)
- Relations linking $\psi^y / \diamond \psi^y$ and Ω^z :
 - $\diamond \psi^y \rightsquigarrow \Omega^z$ iff $y + z > t$. (Corollary 5.)
 - Let $1 \leq z \leq t+1$ and $1 \leq y \leq t$. $\Omega^z \not\rightsquigarrow \diamond \psi^y$. (Theorem 12.)
- Relations linking $\diamond S_x$ and Ω^z :
 - $\diamond S_x \rightsquigarrow \Omega^z$ iff $x + z > t+1$. (Corollary 6.)
 - Let $1 < x, z \leq t$. $\forall z : \Omega^z \not\rightsquigarrow \diamond S_x$. (Theorem 13.)

All these relations are depicted in Figure 1 where the bold arrows mean reducibility, and the dotted arrows mean irreducibility. The class S_x is the subclass of $\diamond S_x$ where the accuracy is perpetual (namely, there is a correct process that is not suspected by x processes from the very beginning). \mathcal{P} is the class of *perfect* failure detectors [4] (the ones that never do a mistake). The column at the right of the figure concerns k -set agreement: all the failure detector classes in the z th line allow solving z -set agreement. Moreover, in the family of failure detectors defined by a column, the class on the line “ z ” is the weakest for solving k -set agreement; and given a line “ z ” of the figure, Ω^z is the weakest failure detector class of that line that allows solving k -set agreement. It is important to notice that, for $1 \leq z \leq t$, we have (1) $\diamond S_{t-z+2}$ and $\diamond \psi^{t-z+1}$ cannot be compared, and (2) both are stronger than Ω^z .

• CONTRIBUTION #3: Additivity. The paper addresses the question of adding failure detectors of distinct classes. This is an important issue as “additivity” is a crucial concept as soon as modularity and scalability of distributed systems are concerned.

As an example, assuming $t > 1$, let us consider the class $\diamond S_t$ that allows solving 2-set agreement (but not consensus), and the class $\diamond \psi^1$ that allows solving t -set agreement (but not $(t-1)$ -set agreement). What about a system with a failure detector in $\diamond S_t$ and one in $\diamond \psi^1$? Which type of information on failures is provided by their combination? The paper shows that $\diamond S_t + \diamond \psi^{t-1}$ allows solving the consensus problem. More generally, with respect to the grid described in the previous figure, the paper characterizes which classes can be added and which cannot. More explicitly, it shows the following result (see also Figure 2): $\diamond S_x + \diamond \psi^y \rightsquigarrow \Omega^z \Leftrightarrow x + y + z > t + 1$. To that end, the paper presents a construction algorithm (sufficiency part, figures 7 and 8), and an impossibility proof (necessity part, Theorem 9).

Intuitively, this shows that $\diamond\mathcal{S}_x$ and $\diamond\psi^y$ provide different types of information on failures to build Ω^z . To see the gain provided by such an addition, let us analyze it as follows:

- As $\diamond\mathcal{S}_x \rightsquigarrow \Omega^{t-x+2}$, the previous addition shows that adding $\diamond\psi^y$ allows strengthening Ω^{t-x+2} to obtain Ω^z with $z = (t - x + 2) - y$.
- Similarly, as $\diamond\psi^y \rightsquigarrow \Omega^{t-y+1}$, the previous addition shows that adding $\diamond\mathcal{S}_x$ allows strengthening Ω^{t-y+1} to Ω^z with $z = (t - y + 1) - (x - 1)$.

It is remarkable that the previous addition of failure detectors (Figure 2) shows that, when we consider both of them, the failure detector classes $\diamond\mathcal{S}_x$ and $\diamond\psi^y$ are not robust: adding them allows solving a problem (the $(t + 2 - (x + y))$ -set agreement problem), that none of them taken alone can solve ($\diamond\mathcal{S}_x$ can solve only $(t + 2 - x)$ -set agreement, and $\diamond\psi^y$ can solve only $(t + 1 - y)$ -set agreement).

- **CONTRIBUTION #4: Asynchronous Ω^k -based k -set agreement.** This paper proposes such an algorithm. To our knowledge, no previous work has addressed the design of Ω^z -based k -set agreement algorithms. The proposed algorithm (Figure 5) is very simple. The paper also establishes that, when one is interested in solving the k -set agreement problem in an asynchronous message-passing system equipped with a failure detector of the $(\Omega^z)_{1 \leq z \leq n}$ family, the bounds $t < n/2$ and $z \leq k$ are tight (Theorem 6). Consequently, among all the classes described in Figure 1, Ω^k is the weakest class for solving asynchronous k -set agreement (hence, the algorithm is optimal in that respect). This constitutes a step towards the characterization of the weakest failure detector class that allows solving the k -set agreement problem.

Roadmap The paper is made up of 7 sections plus an appendix. Section 2 describes the asynchronous computing model and the classes of failure detectors we are interested in. Section 3 shows that the failure detector classes ψ^y and ϕ^y (resp., $\diamond\psi^y$ and $\diamond\phi^y$) are equivalent. Section 4 presents the asynchronous Ω^k -based k -set agreement algorithm. Then, Section 5 presents an algorithm that builds a failure detector of the class Ω^z from a pair of underlying failure detectors, one of the class $\diamond\psi^y$, the other of the class $\diamond\mathcal{S}$. Section 6 shows that $x + y + z > t + 1$ is a necessary requirement for the previous construction, and establishes the irreducibility relations depicted by the grid of Figure 1. Finally, Section 7 provides concluding remarks. From a methodology point of view, as much as possible the paper uses reductions (striving not to reinvent the wheel).

2 Computation Model

2.1 Asynchronous System with Process Crash Failures

We consider a system consisting of a finite set Π of $n \geq 2$ processes, namely, $\Pi = \{p_1, p_2, \dots, p_n\}$. When it is not ambiguous we also use Π to denote the set of the identities $1, \dots, n$ of the processes. A process can fail by *crashing*, i.e., by prematurely halting. It behaves correctly (i.e., according to its specification) until it (possibly) crashes. By definition, a process is *correct* in a run if it does not crash in that run; otherwise it is faulty. As previously indicated, t denotes the maximum number of processes that can crash in a run ($1 \leq t < n$). The identity of the process p_i is i , and each process knows all the identities.

Processes communicate and synchronize by sending and receiving messages through channels. Every pair of processes is connected by a channel. Channels are assumed to be reliable: they do not create, alter or lose messages. In particular, if p_i sends a message to p_j , then eventually p_j receives that message unless it fails. There is no assumption about the relative speed of processes or message transfer delays (let us observe that channels are not required to be FIFO).

Broadcast(m) is an abbreviation for “**for_each** $p_j \in \Pi$ **do** *send*(m) **to** p_j **end_for**”. Moreover, we assume (without loss of generality) that the communication system provides the processes with a *reliable broadcast* abstraction [13]. Such an abstraction is made up of two primitives *Broadcast*() and *Deliver*() that allow a process to broadcast and deliver messages (we say accordingly that a message is **R_broadcast** or **R_delivered** by a process) and satisfy the following properties:

- **Validity.** If a process **R_delivers** m , then some process has **R_broadcast** m . (No spurious messages.)

- **Integrity.** A process $R_delivers$ a message m at most once. (No duplication.)
- **Termination.** If a correct process $R_broadcasts$ or $R_delivers$ a message m , then all the correct processes $R_deliver$ m . (No message $R_broadcast$ or $R_delivered$ by a correct process is missed by a correct process.)

As we can see, the messages sent (resp., $R_broadcast$) by a process are not necessarily received (resp., $R_delivered$) in their sending order. Moreover, different processes can $R_deliver$ messages in different order. There is no assumption on message transfer delays. The communication system is consequently reliable and asynchronous.

2.2 The Failure Detector Classes $(\mathcal{S}_x)_{1 \leq x \leq n}$ and $(\diamond \mathcal{S}_x)_{1 \leq x \leq n}$

As indicated in the Introduction, the failure detector classes \mathcal{S}_x and $\diamond \mathcal{S}_x$ have been introduced and used in [12, 23, 24, 30]. A failure detector of the class \mathcal{S}_x or $\diamond \mathcal{S}_x$ consists of a set of modules, each one attached to a process: the module attached to p_i maintains a set (named *suspected_i*) of processes it currently suspects to have crashed. As in other papers devoted to failure detectors, we say “process p_i suspects process p_j at some time τ ”, if $p_j \in \text{suspected}_i$ at that time. Moreover, (by definition) a crashed process suspects no process.

The failure detector $\diamond \mathcal{S}_x$ class generalizes the class $\diamond \mathcal{S}$ defined in [4] (we have $\diamond \mathcal{S}_n = \diamond \mathcal{S}$). A failure detector belongs to the class $\diamond \mathcal{S}_x$ if it satisfies the following properties:

- **Strong Completeness.** Eventually, every process that crashes is permanently suspected by every correct process.
- **Limited Scope Eventual Weak Accuracy:** There is a time after which there is a set Q of x processes such that Q contains a correct process and that process is never suspected by the processes of Q .

Similarly, the class \mathcal{S}_x generalizes the class \mathcal{S} [4] (and we have $\mathcal{S}_n = \mathcal{S}$). A failure detector of the class \mathcal{S}_x satisfies the previous strong completeness property, plus the following accuracy property:

- **Limited Scope Perpetual Weak Accuracy.** there is a set Q of x processes such that (from the very beginning) Q contains a correct process and that process is never suspected by the processes of Q .

It is easy to see that $\mathcal{S}_{x+1} \subseteq \mathcal{S}_x$, $\diamond \mathcal{S}_{x+1} \subseteq \diamond \mathcal{S}_x$, and $\mathcal{S}_x \subseteq \diamond \mathcal{S}_x$. It is also easy to see that the failure detectors of the classes \mathcal{S}_1 and $\diamond \mathcal{S}_1$ provide no information on failures. It is shown in [14] that $\diamond \mathcal{S}_x$ is the weakest failure detector class of the family $(\diamond \mathcal{S}_x)_{1 \leq x \leq n}$ that allows solving k -set agreement for $k = t - x + 2$, in asynchronous message-passing systems with a majority of correct processes ($t < n/2$).

2.3 The Failure Detector Classes $(\Omega^z)_{1 \leq z \leq n}$

This family of failure detectors has been introduced in [26]. A failure detector of the class Ω^z maintains at each process p_i a set of processes of size at most z (denoted *trusted_i*) that satisfies the following property:

- **Eventual Multiple Leadership.** there is a time after which the sets *trusted_i* of the correct processes contain forever the same set of processes and at least one process of this set is correct.

The family $(\Omega^z)_{1 \leq z \leq n}$ generalizes the class of failure detectors Ω defined in [3], with $\Omega^1 = \Omega$.

Recently, another generalization of Ω has been studied in [8] that considers Ω_S , where S is a predefined subset of the processes of the system. Ω_S requires that all the correct processes of S eventually agree on the same correct leader (it is not required that their eventual common leader belongs to S). Let X be the set of all the pairs of processes. It is shown in [8] that, given all the Ω_x , $x \in X$, it is possible to build Ω .

2.4 The Failure Detector Classes $(\phi^y)_{0 \leq y < n}$ and $(\diamond\phi^y)_{0 \leq y < n}$

These failure detector classes have been introduced in [19] and [21]. As noticed in the Introduction, their definition does not comply with the Chandra and Toueg's failure detector framework that restricts the output of a failure detector to depend only on the failure pattern and the invocation time. Here, differently from the previous classes of failure detectors that provide each process p_i with a set (*suspected_i* or *trusted_i*) that p_i can only read, a failure detector provides the processes with a primitive $\text{QUERY}(X)$, where X is a set of process identities supplied by the invoking process. Such a primitive allows a process p_i to query about the crash of a region X of the system.

The classes $(\phi^y)_{0 \leq y < n}$ A failure detector of the class ϕ^y is defined by the following properties (recall that t is an upper bound on the number of process crashes):

- **Triviality property.** If $|X| \leq t - y$, $\text{QUERY}_y(X)$ returns *true*. If $|X| > t$, $\text{QUERY}(X)$ returns *false*.
- **Safety property.** If $t - y < |X| \leq t$ and at least one process in X has not crashed when $\text{QUERY}(X)$ is invoked, the invocation returns *false*.
- **Liveness property.** Let X be such that $t - y < |X| \leq t$. Let τ be a time such that, at time τ , all the processes in X have crashed. There a finite time $\tau' \geq \tau$ from which all the invocations of $\text{QUERY}(X)$ return *true*.

The triviality property provides the invoking process with a pre-determined output when the set X is too small (because the failure detector is not powerful enough to give an answer) or too big (because the answer is obvious). The safety property states that if the output is *true*, then all the processes in X have crashed. The liveness property states that $\text{QUERY}(X)$ eventually outputs *true* when all the processes in X have crashed. It is shown in [19] that (1) $\phi^{y+1} \subseteq \phi^y$, and (2) ϕ^t and the class \mathcal{P} of perfect failure detectors are equivalent in any system where at most t processes can crash. Moreover, it is easy to see that ϕ^0 provides no information on failures.

The classes $(\diamond\phi^y)_{0 \leq y < n}$ The failure detector class $\diamond\phi^y$ is the “eventual” counterpart of the class ϕ^y . More precisely, a failure detector of the class $\diamond\phi^y$ is defined by the following properties (recall that t is an upper bound on the number of process crashes):

- **Triviality property.** If $|X| \leq t - y$, then $\text{QUERY}_y(X)$ returns *true*. If $|X| > t$, then $\text{QUERY}(X)$ returns *false*.
- **Eventual Safety property.** Let X be such that $t - y < |X| \leq t$. Suppose that at least one correct process belongs to X . There a finite time τ from which all the invocations of $\text{QUERY}(X)$ return *false*.
- **Liveness property.** Let X be such that $t - y < |X| \leq t$. Let τ be a time such that, at time τ , all the processes in X have crashed. There a finite time $\tau' \geq \tau$ from which all the invocations of $\text{QUERY}(X)$ return *true*.

As for the classes $(\phi^y)_{0 \leq y \leq t}$, it follows from these properties that (1) $\diamond\phi^{y+1} \subseteq \diamond\phi^y$, and (2) $\diamond\phi^t$ and the class $\diamond\mathcal{P}$ are equivalent in any system where at most t processes can crash.

2.5 The Failure Detector Classes $(\psi^y)_{0 \leq y < n}$ and $(\diamond\psi^y)_{0 \leq y < n}$

The classes $(\psi^y)_{0 \leq y < n}$ A failure detector of the class ψ^y provides each process with an integer nb_c_i that p_i can only read. The current value of this number is an approximation of the number of processes that have crashed (hence the name nb_c_i).

More precisely, let f denote the number of processes that crash in a given run ($0 \leq f \leq t$), f^τ denote the number of processes that have crashed up to time τ , and $nb_c_i^\tau$ denote the value of the failure detector local variable nb_c_i at time τ .

- Safety property. $\forall \tau: t - y \leq nb_c_i^\tau \leq \max(t - y, f^\tau)$.
- Liveness property. $\exists \tau: \forall \tau' \geq \tau: nb_c_i^{\tau'} = \max(t - y, f)$.

The safety property states that the failure detector outputs a value that is never smaller than $t - y$, and is an underestimate of the current number of crashes as soon as at least $t - y$ processes have crashed. The parameter y allows defining a failure detector instance for the algorithms that have to cope with failures only when there are more than $t - y$ crashes. The liveness property states that eventually each nb_c_i local variable converges towards the number of processes that crash in the considered run.

The classes $(\diamond\psi^y)_{0 \leq y < n}$ That class is the eventual counterpart of $(\psi^y)_{0 \leq y < n}$. It allows the previous safety property not to be satisfied during an arbitrary but finite period. This weakening combined with the liveness property can be combined into the following property, where f denote the number of processes that crash in a given run ($0 \leq f \leq t$). This single property is formulated as follows.

- Eventual convergence property. $\exists \tau: \forall \tau' \geq \tau: nb_c_i^{\tau'} = \max(t - y, f)$.

It is easy to see that, differently from the definitions of $(\phi^y)_{0 \leq y < n}$ and $(\diamond\phi^y)_{0 \leq y < n}$, the definitions of $(\psi^y)_{0 \leq y < n}$ and $(\diamond\psi^y)_{0 \leq y < n}$ do comply with the Chandra and Toueg's failure detector definition framework.

2.6 Notation

Let \mathcal{F} and \mathcal{G} be any two classes among the previous classes of failure detectors. The notation $\mathcal{AS}_{n,t}[\mathcal{F}]$ is used to represent a message-passing asynchronous system made up of n processes, where up to t may crash, equipped with a failure detector of the class \mathcal{F} . Similarly, $\mathcal{AS}_{n,t}[\mathcal{F}, \mathcal{G}]$ denotes a system equipped with a failure detector of the class \mathcal{F} and a failure detector of the class \mathcal{G} . Finally, $\mathcal{AS}_{n,t}[\emptyset]$ denotes a “pure” asynchronous message-passing system (i.e., with no failure detector).

3 The Classes ψ^y ($\diamond\psi^y$) and ϕ^y ($\diamond\phi^y$) are Equivalent

This section shows that the failure detector classes ϕ^y and ψ^y (resp., $\diamond\phi^y$ and $\diamond\psi^y$) have the same computational power as far as the information on failures is concerned.

Once we know that ϕ^y and ψ^y ($\diamond\phi^y$ and $\diamond\psi^y$) are equivalent, it becomes possible to use ϕ^y ($\diamond\phi^y$) instead of ψ^y ($\diamond\psi^y$) to prove lower bounds and (ir)reducibility results (as done in Section 6).

3.1 From ϕ^y ($\diamond\phi^y$) to ψ^y ($\diamond\psi^y$)

This section shows that, for any y , $1 \leq y \leq n$, given any failure detector of the class ψ^y (resp., $\diamond\psi^y$) it is possible to build a failure detector of the class ϕ^y (resp., $\diamond\phi^y$).

A transformation For each α , $t - y + 1 \leq \alpha \leq t$, let $Sets(\alpha)$ be the set including all the subsets of Π that contain α processes. There are y such sets, namely, from $Sets(t - y + 1)$ until $Sets(t)$.

The algorithm described in Figure 3 builds a failure detector of ψ^y (resp., $\diamond\psi^y$) from any failure detector of ϕ^y (resp., $\diamond\phi^y$). At each process p_i , it consists in an infinite loop that repeatedly updates the local variable nb_c_i whose value defines the current output of ψ^y (resp., $\diamond\psi^y$). The primitive ϕ -QUERY(X), where X is a set of processes, allows a process p_i to query its underlying ϕ^y failure detector that returns *true* or *false* according to the current state (alive or crashed) of the processes of X .

The body of the loop for p_i consists in invoking ϕ -QUERY(X) for each possible set X of α processes, with α varying from $t - y + 1$ to t . If ϕ -QUERY(X) answers *true* for the current set X , p_i concludes that the α processes of X have crashed; accordingly, it keeps the current value of α in a set S_i . When it has probed all the possible sets, p_i updates nb_c_i according to the value of S_i . (This algorithm can be improved. We do not do it in order to keep it as simple as possible.)

```

nbci ← t - y;
repeat forever
  Ai ← ∅;
  for_each α ∈ {t - y + 1, ..., t} do
    for_each X ∈ Sets(α) do
      if φ-QUERY(X) then Ai ← Ai ∪ {α} end_if
    end_for
  end_for;
  if Ai ≠ ∅ then nbci ← max(Ai) else nbci ← (t - y) end_if
end_repeat

```

Figure 3: From ϕ^y to ψ^y (resp. From $\diamond\phi^y$ to $\diamond\psi^y$), (code for p_i)

Theorem 1 *Given any failure detector of the class ϕ^y (resp., $\diamond\phi^y$), the algorithm described in Figure 3 builds a failure detector of the class ψ^y (resp., $\diamond\psi^y$).*

Proof The proof addresses simultaneously the case where the underlying failure detector belongs to the class ϕ^y , and the case where it belongs to $\diamond\phi^y$. Taking an arbitrary run, it considers two cases according to the number f of processes that crash in that run ($0 \leq f \leq t$).

- $f < t - y + 1$. In that case, any set X , that belongs to a set $Sets(\alpha)$ for some α ($t - y + 1 \leq \alpha \leq t$), contains at least one correct process. It follows from the safety property of the underlying failure detector that there is a finite time τ ($\tau = 0$ for ϕ^y and $\tau \geq 0$ for $\diamond\phi^y$) after which, for any X as defined previously, ϕ -QUERY(X) returns *false*. Consequently, after time τ , for any process p_i , we always have $A_i = \emptyset$ at the end of the outer **for_each** loop. We conclude from the text of the algorithm that, after τ , each local variable nb_{c_i} remains forever equal to $t - y$.
- $f \geq t - y + 1$. Let E be the set of processes that crash (so, $|E| = f$). Due to the definition of the sets $Sets(t - y + 1), \dots, Sets(t)$, there is a set X in one of these sets such that $E = X$. According to the order in which the processes of E crash, let τ be the time at which the last process of E crashes.

Let us first observe that, when the underlying failure detector belongs to the class ϕ^y , it follows from its safety property that all the ϕ -QUERY(E) invocations issued before τ returns *false*. Differently, if it belongs to $\diamond\phi^y$, a ϕ -QUERY(E) invocation issued before τ can return *true* or *false*. Moreover, it follows from the liveness property of ϕ^y and $\diamond\phi^y$, that there is a time $\tau' \geq \tau$ after which all the invocations ϕ -QUERY(E) return *true*.

- Case 1: The underlying failure detector belongs to $\diamond\phi^y$. There is a time $\tau'' \geq \tau'$ after which any ϕ -QUERY(X) issued by a process p_i and such that $|X| > f$ returns *false* (eventual safety property of $\diamond\phi^y$), and any ϕ -QUERY(E) returns *true* (liveness property of $\diamond\phi^y$). It follows that, after time τ'' , we always have $\max(A_i) = f$ before executing the last **if** statement. Consequently, after τ'' , nb_{c_i} keeps forever the value f . As $f > t - y$, the eventual convergence property of $\diamond\phi^y$ follows.
- Case 2: The underlying failure detector belongs to ϕ^y . During the period during which no more than $t - y$ processes crash, as all the sets X used in the algorithm are such that $|X| > t - y$, it follows that all the invocations ϕ -QUERY(X) issued during that period return *false*. The set A_i remains consequently empty, and $nb_{c_i} = t - y$ during that period.

Let time $\tau(f')$ be a time at which exactly f' ($t - y < f' \leq f$) processes have crashed (i.e., the remaining $f - f'$ processes have not yet crashed). For notational convenience, let $\tau(f + 1) = +\infty$. It follows from the safety property of ϕ^y that any ϕ -QUERY(X) with $|X| > f'$ returns *false* at least until $\tau(f' + 1)$. Consequently, until $\tau(f' + 1)$, the greatest value that A_i can contain is f' , which proves the safety property of ψ^y .

To prove the liveness property of ψ^y , it is sufficient to show that there is a time after which nb_{c_i} keeps forever the value f . There is a finite time $\tau' \geq \tau(f)$ after which ϕ -QUERY(E) returns always *true* (liveness property of ϕ), and ϕ -QUERY(X) with $|X| > f$ always return *false* (safety property

of ϕ). It follows from this observation that, after τ' , we always have $\max(A_i) = f = |E|$ before executing the last **if** statement. Consequently, from τ' , nb_c_i keeps forever the value $f = |E|$.

□*Theorem 1*

3.2 From ψ^y ($\diamond\psi^y$) to ϕ^y ($\diamond\phi^y$)

A transformation The algorithm that builds a failure detector of the class ϕ^y ($\diamond\phi^y$) from a failure detector of the class ψ^y ($\diamond\psi^y$) is described in Figure 4. Let ϕ -QUERY(X) denote the operation of the failure detector of the class ϕ^y ($\diamond\phi^y$). The underlying failure detector of the class ψ^y ($\diamond\psi^y$) provides each process p_i with an integer local variable nb_c_i that p_i can only read.

When p_i invokes ϕ -QUERY(X), it first checks the size of X . If X is too small (resp., too big), the value *true* (resp., *false*) is returned. Otherwise, the size of X is such that $t - y < |X| \leq t$. In that case, p_i saves the current value of nb_c_i in a local variable est_c_i , and sends an INQUIRY(sn_i) message (timestamped with the next sequence number) to every process. It then waits (line 06) until either it has received “enough” corresponding responses (i.e., that carry the sequence number sn_i) or the value of $n - nb_c_i$ has changed. “Enough” means here $n - nb_c_i$ (while it is waiting, p_i checks regularly the condition; each time it checks it, it reads the (possibly new) value of nb_c_i). If the value of nb_c_i has changed, p_i starts a new inquiry (line 04). Otherwise the inquiry timestamped sn_i is successful and p_i collects in rec_i the ids of the processes that sent a response matching the inquiry. Finally, if one process p_j in X is also in rec_i , that process was not crashed when p_i sent the inquiry message. The value *false* is then returned. Otherwise ($rec_i \cap X = \emptyset$), the value *true* is returned.

<pre> operation ϕ-QUERY(X): (0 1) case $X \leq t - y$ then return (<i>true</i>) (0 2) $t < X$ then return (<i>false</i>) (0 3) $t - y < X \leq t$ then (0 4) repeat $sn_i \leftarrow sn_i + 1; est_c_i \leftarrow nb_c_i;$ (0 5) for_each $j \in \{1, \dots, n\}$ do send INQUIRY(sn_i) to p_j end_do; (0 6) wait until ((RESPONSE(sn_i) received from $n - est_c_i$ processes) \vee ($est_c_i \neq nb_c_i$)); (0 7) until $est_c_i = nb_c_i$ end_repeat; (0 8) let $rec_i = \{j \mid \text{RESPONSE}(sn_i) \text{ has been received from } p_j\}$; (0 9) return ($X \cap rec_i = \emptyset$) (10) endcase Background task: when INQUIRY(sn) is received from p_j: send RESPONSE(sn) to p_j </pre>
--

Figure 4: From ψ^y to ϕ^y (resp. From $\diamond\psi^y$ to $\diamond\phi^y$), (code for p_i)

Theorem 2 *Given any failure detector of the class ψ^y (resp., $\diamond\psi^y$), the algorithm described in Figure 4 builds a failure detector of the class ϕ^y (resp., $\diamond\phi^y$).*

Proof Considering an arbitrary run, let f be the number of processes that crash in that run. The proof is decomposed in five parts.

- [Termination] Let us first show that each invocation of ϕ -QUERY(X) by a correct process terminates. If $|X| \leq t - y$ or $|X| > t$, the operation trivially terminates. So, assuming that $t - y < |X| \leq t$, let us consider two cases.
 - Case 1: $f \leq t - y$. In that case, nb_c_i is constant and always equal to $t - y$. Consequently, $n - f \geq n - (t - y) = n - est_c_i$. As there are $n - f$ correct processes, p_i always receive $n - est_c_i$ matching responses to each inquiry message. It follows that the inner **wait until** always terminates. As, in the current case, we always have $est_c_i n - nb_c_i = t - y$, the **repeat** statement always terminates.

- Case 2: $f > t - y$. Let us first consider the **wait until** statement, and let us assume that p_i remains blocked forever. This means that more than est_{c_i} (say x) processes have crashed (otherwise, p_i will receive enough responses to proceed). As p_i is blocked forever, we conclude from the wait condition (line 06) that the predicate $est_{c_i} = nb_{c_i}$ remains true forever, which means that nb_{c_i} remains equal to x forever. But, due to the properties of ψ^y and $\diamond\psi^y$, as more than x processes have crashed there is a time after which we always have $nb_{c_i} > x = est_{c_i}$, which contradicts the fact that p_i blocks forever in the **wait until** statement.

Let us now consider the **repeat** statement. Its termination follows from the liveness of ϕ^y , or the eventual convergence of $\diamond\phi^y$, that states there is a time after which nb_{c_i} remains always equal to f . Consequently, after that time we necessarily always have $est_{c_i} = nb_{c_i} = f$, which proves the termination of the **repeat** statement.

- [Triviality property of ϕ^y and $\diamond\phi^y$] That property is trivially guaranteed by the **case** statement.
- [Liveness property of ϕ^y and $\diamond\phi^y$] Let E , with $|E| > t - y$, be a set of processes that crash. Moreover, let $\tau(E)$ be a time after which all the processes of E have crashed. Due to the liveness property of ψ^y or $\diamond\psi^y$, there is a time τ after which nb_{c_i} remains forever equal to f .

Let $\tau' \geq \max(\tau(E), \tau)$. Any ϕ -QUERY(E) issued after τ' (1) terminates (see above), and (2) does not receive responses from the processes in E as they have crashed before $\tau(E)$. It follows that $rec_i \cap E = \emptyset$, and ϕ -QUERY(E) returns *true*.

- [Safety property of $\diamond\phi^y$] Let X be a set of processes such that $t - y < |X| \leq t$ and at least one process of X does not crash. We have to show that there is a time after which any ϕ -QUERY(X) returns *false*.

- Case 1: $f \leq t - y$. In that case, it follows from the eventual convergence property of $\diamond\psi^y$ that there is a time τ after which we always have $est_{c_i} = nb_{c_i} = t - y$. As there are $n - f$ correct processes and $n - f \geq n - (t - y)$, it follows that after τ , a process p_i receives $n - (t - y)$ matching responses each time it broadcasts an inquiry. As $|X| > t - y$, it follows that $X \cap rec_i \neq \emptyset$, and ϕ -QUERY(X) returns *false*. The eventual safety property of $\diamond\phi^y$ is consequently satisfied.

- Case 2: $f > t - y$. In that case there is a time τ after which all the faulty processes have crashed and we always have $est_{c_i} = nb_{c_i} = f > t - y$. After τ , a process p_i that invokes ϕ -QUERY(X) always receives $n - f$ corresponding responses, one from each correct process. It follows that, after τ , $X \cap rec_i \neq \emptyset$ iff at least one correct process belongs to set X , which proves the eventual safety property of $\diamond\phi^y$.

- [Safety property of ϕ^y] Let X be a set of processes such that $t - y < |X| \leq t$ and at least one process of X has not crashed at time τ . We have to show that any ϕ -QUERY(X) issued before time τ returns *false*.

- Case 1: $f \leq t - y$. In that case, it follows from the safety property of ϕ^y that, from $\tau = 0$, we always have $nb_{c_i} = t - y$, i.e., $est_{c_i} = nb_{c_i} = t - y$. Taking $\tau = 0$, the proof is then the same as the proof of the corresponding case in the proof of the safety property of $\diamond\phi^y$.

- Case 2: $f > t - y$. We claim (claim C) that, when a process p_i terminates the **repeat** loop, it has received a matching response from each process that does not crash before the ϕ -QUERY(X) returns a value.

Let A be the set of processes that have not crashed before ϕ -QUERY(X) terminates. It follows from the claim that $A \subseteq rec_i$. Hence, if X contains a process that has not crashed before ϕ -QUERY(X) terminates, we have $X \cap A \neq \emptyset$, and consequently, $X \cap rec_i \neq \emptyset$. It follows that ϕ -QUERY(X) returns *false*.

Proof of the claim C. Considering the last execution of the repeat loop body of a ϕ -QUERY(X) invocation issued by a process p_i , let sn be the corresponding sequence number, τ_b be the time at which p_i reads the current value x of nb_{c_i} (line 04), and τ_e be the time at which it reads again x from nb_{c_i} (line 07). We have $est_{c_i} = x$ during this loop execution. Let f^{τ_b} and f^{τ_e} be the number of processes that have crashed by time τ_b and τ_e , respectively.

Due to the safety property of ϕ^y , we have $x \leq f^{\tau_b} \leq f^{\tau_e}$. Moreover, (1) no process crashed at time τ_b sends a `RESPONSE(sn)` message; (2) all the processes that are alive at τ_e sent a `RESPONSE(sn)` message to p_i ; (3) the p_i has received $n - x$ `RESPONSE(sn)` messages; and (4) $n - x \geq n - f^{\tau_e}$. It follows from the previous points that p_i received `RESPONSE(sn)` from each process that was alive at time τ_e . The claim follows. *End of the proof of the claim C.*

□*Theorem 2*

A simpler transformation for the class $\diamond\phi^y$ The proof of the safety properties of Theorem 2 relies on a strong synchronization realized by the **repeat** loop and the `estci` `sni` local variables (lines 04-07). This synchronization is used to isolate an inquiry period during which `nbci` remains constant.

Actually, this synchronization is stronger than necessary to ensure the eventual safety property of $\diamond\phi^y$. A much less synchronized transformation algorithm works for this class. More precisely, the local variables `estci` and `sni` can be suppressed, and the **repeat** statement (lines 04-07) can be replaced by the two following lines:

```
for_each  $j \in \{1, \dots, n\}$  do send INQUIRY() to  $p_j$  end_do;  
wait until (RESPONSE() received from  $n - nb_{c_i}$  processes).
```

The proof is left to the reader. (That proof has to consider the fact that there is a time after which all the response messages sent by a crashed process have arrived.)

4 Using Ω^k to Solve k -Set Agreement

This section presents an Ω^k -based k -set agreement algorithm, and lower bounds on when solving k -set agreement with failure detector classes of the family $(\Omega^z)_{1 \leq z \leq n}$ is possible. These lower bounds are $t < n/2$ and $z \leq k$. Interestingly, the proof of these bounds is based on a reduction to a $\diamond\mathcal{S}_x$ -based k -set agreement algorithm and a corresponding lower bound [14].

4.1 A k -Set Agreement Algorithm

The algorithm, described in Figure 5, is a simple adaptation of an Ω -based consensus algorithm described in [11] (which is in turn inspired from a $\diamond\mathcal{S}$ -based consensus algorithm described in [22]); it assumes $t < n/2$. A process p_i invokes `k-SET-AGREEMENT(v_i)`, where v_i is the value it proposes. If it does not crash, it terminates when it executes the statement `return(v)`, where v is then the value it decides.

The function `k-SET-AGREEMENT(v_i)` is made up of two tasks. The task *T2* is used to disseminate a decided value and prevent deadlock: due to the reliable broadcast, as soon as a process decides, all the correct processes decide. In the main task *T1*, the processes proceed in consecutive asynchronous rounds, each round being made up of two phases, each including a communication step. When considering a process p_i , the local variable `esti` is the local estimate of the decision value; r_i is its current round number.

During the first phase of round r , p_i first reads `trustedi` (the set provided by its underlying failure detector module of the class Ω^z), stores its value in `Li`, and sends a `PHASE1(r_i, L_i, est_i)` message to all the processes. Then, p_i waits until it has received such round r messages from $n - t$ processes (i.e., from at least a majority) and it has either received such a message from a process of its `Li` set or the set `trustedi` has changed. Then, if a majority of processes have the same leader set L , and p_i has received an estimate value v_L from a process in this set L , it keeps v_L in `auxi`, otherwise it sets `auxi` to \perp . Let us notice that we can conclude from the previous statements (see the proof) that, at the end of the first phase of each round, the set of the `auxi` local variables contains at most $|L_i| = k$ distinct values different from \perp .

The second phase of a round aims at allowing the processes to decide, while ensuring that no more than k different values are decided, whatever the round during which a process decides. To that end, each process p_i broadcasts a `PHASE2(r_i, aux_i)` message to all the processes, and then waits until it has received such messages from $n - t$ processes. If it receives a non- \perp value v , it adopts v as its new estimate (if there are several such values, it takes one arbitrarily). Moreover, if none of the values it has received is \perp , it decides the estimate value v it has just adopted; this is done by broadcasting v in a reliable way, and then returning that value (in task *T2*).

<p>Function $k\text{-SET_AGREEMENT}(v_i)$: Init: $est_i \leftarrow v_i; r_i \leftarrow 0$</p> <p>Task T1:</p> <pre style="margin: 0;"> (0 1) repeat forever ----- Phase 1 ----- (0 2) $r_i \leftarrow r_i + 1; L_i \leftarrow trusted_i;$ (0 3) <i>Broadcast</i> PHASE1(r_i, L_i, est_i); (0 4) wait until (PHASE1($r_i, _, _$) received from $\geq (n - t)$ processes); (0 5) wait until ((PHASE1($r_i, _, _$) received from a process $\in L_i$) \vee ($L_i \neq trusted_i$)); (0 6) if (($\exists L : \text{PHASE1}(r_i, L, _)$ received from a majority of processes) (0 7) \wedge(PHASE1($r_i, _, v_L$) received from a process $\in L$)) (0 8) then $aux_i \leftarrow v_L$ else $aux_i \leftarrow \perp$ end_if; % Here $\{aux_j : j \in \Pi \wedge aux_j \neq \perp\} \leq L_i = k$ % ----- Phase 2 ----- (0 9) <i>Broadcast</i> PHASE2(r_i, aux_i); (10) wait until (PHASE2($r_i, _$) received from $(n - t)$ processes); (11) let $rec_i = \{ aux : \text{PHASE2}(r_i, aux)$ has been received }; (12) if ($\exists v : v \neq \perp \wedge v \in rec_i$) then $est_i \leftarrow v$ end_if; (13) if ($\perp \notin rec_i$) then <i>R-Broadcast</i> DECISION(est_i); stop T1 end_if (14) end_repeat Task T2: when DECISION(v) is R_delivered: <i>return</i>(v); stop T2 </pre>
--

Figure 5: Ω^k -based k -set agreement algorithm (code for p_i)

4.2 Short Discussion

The notion of *perfection*, *oracle-efficiency* and *zero-degradation* used below are straightforward generalizations of the same notions introduced in [7, 11] in the context of failure detector-based consensus algorithms.

Let a failure detector of the class Ω^k be *perfect* if, from the very beginning, it delivers to the processes the same set of at most k processes including at least one correct process. A set agreement algorithm is *oracle-efficient* if it terminates in two communication steps (a single round) when its underlying failure detector is perfect and there is no crash. It is easy to see that the previous algorithm is oracle-efficient. This algorithm satisfies an even stronger property, namely, it is zero-degrading. A set agreement algorithm is *zero-degrading* if it terminates in two steps when its underlying failure detector is perfect and there are only initial crashes (a crash is *initial* if the corresponding process crashes before the algorithm starts). The reader can easily check that the proposed algorithm is zero-degrading. Zero-degradation is particularly important when a set agreement algorithm is used repeatedly: it means that future executions do not suffer from past process failures as soon as the failure detector behaves perfectly.

4.3 Proof of the Algorithm

The proof is similar to the proof of the Ω -based consensus algorithm described in [11]. It assumes $t < n/2$ and $z \leq k$ (see Theorem 6).

Lemma 1 *No correct process blocks forever in a round.*

Proof Let p_i be a correct process. We have to show, whatever the round number r , that p_i cannot be blocked forever in the **wait** statements (lines 04, 05 and 10) of round r . This follows from (1) the fact that t being the maximum number of faulty processes, (2) the termination and integrity properties of the reliable broadcast primitive, as well as (3) the fact that the leader set eventually permanently contains a correct process. In more detail, we have the following.

If a process decides, then by the termination property of the reliable broadcast of the corresponding DECISION() message, every correct process decides, and consequently no correct process can block forever in a round. Assume by contradiction that no process decides. Let r be the smallest round in which some correct process p_i blocks forever. So p_i blocks at line 04, 05 or 10. Consider the case of line 04. Since no correct process blocks in a round $r' < r$ and no correct process decides, all correct processes broadcast a

PHASE1($r, _, _$) message. As the maximum number of faulty processes is t , it follows from the integrity and termination of the broadcast primitive that p_i eventually delivers $n - t$ such messages. Consequently, p_i cannot block at line 04. The fact that p_i cannot block forever at line 05 follows directly from the fact that its local set $trusted_i$ eventually permanently contains the identity of a correct process and the fact that all the correct processes broadcast a PHASE1($r, _, _$) message. Consider line 10: as we have just shown that no correct process blocks forever in phase 1 of round r , it follows that all correct processes broadcast a PHASE2($r, _, _$) message. Consequently (as in line 04), p_i does not block forever at line 10. $\square_{Lemma\ 1}$

Assuming p_i completes line 08 during round r , let $aux_i[r]$ be the value of aux_i after it has been updated by p_i at line 08. Moreover, let $AUX[r] = \{aux_i[r] \mid p_i \text{ completes phase 1 of } r\}$.

Lemma 2 $\forall r : |\{v : v \in AUX[r] \wedge v \neq \perp\}| \leq k$.

Proof Let p_i be a process that completes phase 1 of round r . Let us observe that p_i sets aux_i to a value $v \neq \perp$ only if it sees that a majority of processes have the same leader set L (lines 06-08). Moreover, v is a value proposed by a process that belongs to L . Let us notice that there is at most one set that is considered leader set by a majority of processes. Consequently, all the values $aux_i \neq \perp$ at the end of the round r are estimate values of processes belonging to the same set L . Since this set is of size k , it follows that $|\{aux_i[r] : aux_i[r] \neq \perp \wedge p_i \text{ completes phase 1 of round } r\}| \leq k$. $\square_{Lemma\ 2}$

Lemma 3 *Suppose that no process decides. $\exists r : \perp \notin AUX[r]$.*

Proof It follows from the eventual multiple leadership of the class Ω^k that there is a time τ after which all the processes have permanently the same leader set L and this set contains a correct process. Let r be a round that starts after that time (i.e., the first process, say p_i , that executes $r_i \leftarrow r$ does so at time $\tau' > \tau$). As no correct process blocks in the round r (Lemma 1), each correct process broadcasts PHASE1($r, _, _$), from which it follows that the condition of the **if** statement of line 06-07 is satisfied for all the processes that complete phase 1 of round r . Consequently, no process p_i sets its aux_i variable to \perp . $\square_{Lemma\ 2}$

Theorem 3 [Validity] *Any decided value is a proposed value.*

Proof The special value \perp cannot be decided (lines 12-13). Moreover, it follows from the integrity and validity of the broadcast primitive that the aux_i and est_i variables can only contain proposed values or \perp . $\square_{Theorem\ 3}$

Theorem 4 [Agreement] *At most k distinct values are decided.*

Proof If no process decides, the theorem is trivially true. So, let us assume that a process decides and let r be the smallest round during which some process decides (“decide v during r ” means “during r , execute line 13 with $\perp \notin rec_i \wedge est_i = v$ ”). We first show that there is a set V of values, $|V| \leq k$, such that any process that decides during r decides a value from V . We then show that any value decided during a subsequent round belongs to V .

Let $V = \{v : v \in AUX[r] \wedge v \neq \perp\}$. Let us first notice that $|V| \leq k$ (Lemma 2). Let p_i be a process that decides during round r . Let $rec_i[r]$ be the value of the set rec_i computed at line 11 of round r . Let us observe that $rec_i[r] \subseteq AUX[r]$ (lines 10-11). Since p_i decides a value $v \neq \perp$ in $rec_i[r]$, we have $v \in V$.

Assuming that some process p_i decides a value $v \in V$ during round r , we now prove that the estimate est_j of any process p_j that progresses to $r + 1$ belongs to V . As there are at least $n - t$ PHASE2($r, _$) messages carrying a value $aux \neq \perp$ (these are the messages that allowed p_i to decide during round r) and $n - t > n/2$, it follows from the integrity and validity properties of the broadcast primitive that p_j has received at least one of these PHASE2 messages. Consequently, when p_j executes line 12, it updates its estimate to a value $aux \neq \perp$. Hence, from the definition of set V we have $est_j \in V$. It follows that estimate est_j of all the processes p_j that start the round $r + 1$ belong to V . $\square_{Theorem\ 4}$

Theorem 5 [Termination] *Every correct process eventually decides.*

Proof The proof is by contradiction. Assume that no correct process decides. By Lemma 1, the correct processes progress from round to round. Hence, due to Lemma 3, there is a round r such that $\perp \notin AUX[r]$. Consequently, any $PHASE2(r, aux)$ that is broadcast is such that $aux \neq \perp$. Due to the integrity and termination properties of the broadcast primitive, we have $\perp \notin rec_i$ for any process p_i executing the second phase of round r . We can then conclude (line 13) that the correct processes decide: a contradiction. $\square_{Theorem 5}$

4.4 A Lower Bound

Considering an asynchronous message-passing system equipped with a failure detector of the class Ω^z , $1 \leq z \leq n$, this section establishes that $t < n/2$ and $z \leq k$ are necessary and sufficient conditions for solving the k -set agreement problem. As already noticed, this result is obtained by a reduction to the problem of the weakest failure detector in the family $(\diamond\mathcal{S}_x)_{1 \leq x \leq n}$ that allows solving k -set agreement.

Theorem 6 *The k -set agreement problem is solvable in $\mathcal{AS}_{n,t}[\Omega^z]$ if and only if $t < n/2$ and $z \leq k$.*

Proof [\Rightarrow part] The proof is by contradiction. let us assume that there is an algorithm \mathcal{A} that solves the k -set agreement problem in $\mathcal{AS}_{n,t}[\Omega^z]$ such that $t \geq n/2$ or $z > k$. Due to Theorem 6, there is an algorithm \mathcal{T} that builds a failure detector of the class Ω^z in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_{t-z+2}]$. Moreover, there are such transformation algorithms (e.g., the one presented in Section 5 with $y = 0$) that are independent of the value of t (i.e., $t < n$). Combining such a transformation \mathcal{T} and the algorithm \mathcal{A} , we obtain an algorithm that solves the k -set agreement problem in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_{t-z+2}]$. It then follows from the lower bound established by Herlihy and Penso [14] for solving the k -set agreement problem in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_{t-z+2}]$ that $t < \min(n/2, (t - z + 2) + k - 1)$, from which we conclude $t < n/2$ and $z \leq k$: a contradiction.

[\Leftarrow part] This part follows directly from the very existence of the Ω^k -based k -set agreement algorithm described in Section 4.1 and proved in Section 4.3. $\square_{Theorem 6}$

5 Additivity of the Failure Detector Classes $\diamond\mathcal{S}_x$ and $\diamond\psi^y$

This section presents an algorithm that, given as input any pair of failure detectors of the classes $\diamond\mathcal{S}_x$ and $\diamond\psi^y$, constructs a failure detector of the class Ω^z , provided that $x + y + z > t + 1$. (It is proved in Section 6.1 that this is a necessary requirement for such a construction, thereby showing that the algorithm is optimal.)

The algorithm is made up of two components that we call *wheels* because each “turns” like a gear-wheel until they become synchronized and stop turning. The wheel that is the first to eventually stop is the one whose progress depends on the underlying $\diamond\mathcal{S}_x$ failure detector (“lower” wheel). When it stops, it provides a property that allows the second wheel in turn to eventually stop (“upper” wheel). As we will see, the wheel metaphor comes from the fact that each component is made up of main tasks that “turn”, each scanning a sequence until some property becomes satisfied.

Let us remind that $1 \leq x \leq n$. Moreover, as the class $\diamond\psi^t$ is equivalent to the class of eventual perfect failure detectors we consider only the cases $0 \leq y \leq t$, from which we conclude $t - y + 1 > 0$. Finally, as $z \geq t + 2 - (x + y)$ is a necessary requirement and Ω^1 is the strongest class in the family $(\Omega^z)_{1 \leq z \leq n}$, the only interesting cases for the pair (x, y) are when $t + 2 - (x + y) \geq 1$. Hence, in the following we consider that $t - y + 1 > 0$, $z = t + 2 - (x + y)$ and $t + 2 - (x + y) > 0$.

5.1 The Lower Wheel Component

5.1.1 Description

The aim of this component is to provide each process p_i with a local variable $repr_i$ intended to contain a process identity such that the following property becomes eventually satisfied: there is a set X of x processes that either have crashed, or the variables $repr_i$ of the processes of X that have not crashed contain the

identity ℓx of one of them that is a correct process. This process is their common representative (leader). The variable $repr_i$ of a process p_i that does not belong to X has to be equal to the identity i of p_i .

To attain this goal the processes use their local sets $suspected_i$ that collectively satisfy the completeness and limited scope eventual accuracy properties defining the class $\diamond\mathcal{S}_x$. Let \mathcal{X} be the finite set of all the sets of x processes that can be built from the set $\Pi = \{p_1, \dots, p_n\}$. Let nb_x denote the number of combinations of x elements in a set of n elements. \mathcal{X} has nb_x elements. Let us organize \mathcal{X} as a sequence, and let $\mathcal{X}[k]$ be its k th element, $1 \leq k \leq nb_x$. Within $\mathcal{X}[k]$, let us arrange the x processes it is made up of in some predefined (arbitrary) order: $\ell_1^k, \dots, \ell_x^k$. This means that the infinite sequence $\mathcal{X}[1], \mathcal{X}[2], \dots, \mathcal{X}[nb_x], \mathcal{X}[1], \mathcal{X}[2], \dots, \mathcal{X}[nb_x], \mathcal{X}[1], \dots$ gives rise to an infinite sequence of process identities, namely, $\ell_1^1, \dots, \ell_x^1, \ell_1^2, \dots, \ell_x^2, \ell_1^3, \dots$ (see Figure 6). This sequence is assumed to be initially known by all the processes in order they can scan it in the same order.

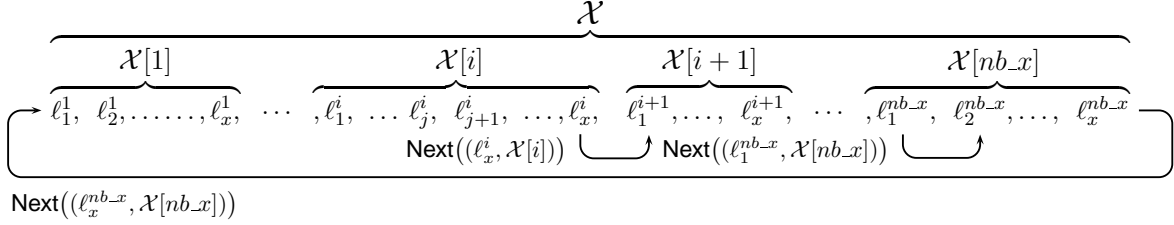


Figure 6: The $\text{Next}()$ function on the logical ring (ℓ, X)

In addition to its output $repr_i$, each process p_i manages a local set X_i and a local variable ℓx_i . It starts with X_i initialized to $\mathcal{X}[1]$, and ℓx_i initialized to ℓ_1^1 (the first process of $\mathcal{X}[1]$). Then, it uses the function $\text{Next}(-, -)$ defined as follows to progress along the infinite sequence of process identities. $\text{Next}(\ell_y^k, \mathcal{X}[k])$ outputs the pair $(\ell_{y+1}^k, \mathcal{X}[k])$ if $y < x$ and the pair $(\ell_1^{k+1}, \mathcal{X}[k+1])$ if $y = x$ (with $k+1$ being replaced by 1 when $k = nb_x$).

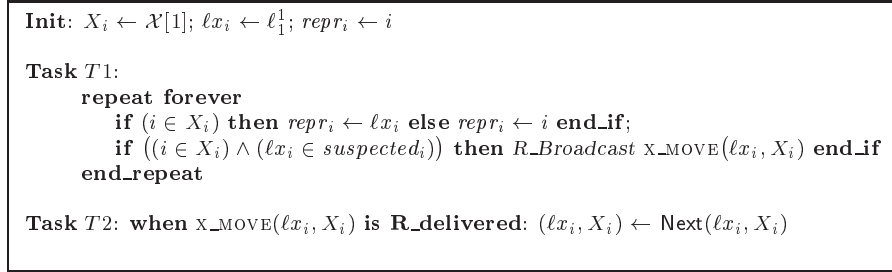


Figure 7: From $\diamond\psi^y + \diamond\mathcal{S}_x$ to Ω^z : lower wheel component (code for p_i)

The behavior of the lower wheel component of a process p_i is described in Figure 7. It is made up of two simple tasks. The processes scan the infinite sequence of sets generated from \mathcal{X} until they stabilize. X_i represents the set of x processes that are currently in charge of extracting a common representative ℓx_i from this set. To do it, each process p_i that belongs to X_i uses its set $suspected_i$ provided by the underlying failure detector of the class $\diamond\mathcal{S}_x$. If the processes of X_i succeed in not suspecting one of them -whose identity is kept by p_i in ℓx_i -, they stop sending $x_MOVE()$ messages. Differently, if a process p_j of the set X_i suspects its current "leader" ℓx_j , it uses the reliable broadcast primitive to send the message $x_MOVE(\ell x_j, X_i)$ indicating that, from its point of view, ℓx_j cannot be their common representative. A process p_j delivers a message $x_MOVE(\ell x, X)$ only when $\ell x = \ell x_i$ and $X_i = X$; it then proceeds to the next process identity (according to the infinite sequence), and possibly to the next candidate set $\mathcal{X}[k+1]$ if $X_i = X = \mathcal{X}[k]$ and $\ell x = \ell x_i$ is the last process of $\mathcal{X}[k]$.

Let us finally consider the case where the processes progress until they consider a set X such that the x processes that constitute X have crashed. It is easy to see that each no-crashed process p_i continues looping inside task $T1$ without sending messages, and is such that $repr_i = i$.

5.1.2 Proof of the lower wheel component

The proof considers an arbitrary run of the algorithm described in Figure 7. C denotes the set of processes that are correct in that run. Moreover, var_i^τ denotes the value of the local variable var_i at time τ .

Lemma 4 $\forall i \in C$, there are a pair (λ_i, σ_i) and a time τ_i such that $\forall \tau \geq \tau_i : (\ell x_i^\tau, X_i^\tau) = (\lambda_i, \sigma_i)$.

Proof We claim (*Claim C1*) that there is a pair (ℓ, X) such that the number of $x_MOVE(\ell, X)$ messages that are sent is finite. Let us assume (by way of contradiction) that there is no pair (λ_i, σ_i) such that after some time $(\ell x_i, X_i) = (\lambda_i, \sigma_i)$ remains true forever. As the pairs $(\ell x, X)$ are arranged in a logical ring (see Figure 6), it follows from the way p_i updates its local pair $(\ell x_i, X_i)$ that the sequence of the successive values of the local variables $(\ell x_i, X_i)$ is $(\ell_1^1, \mathcal{X}[1]), (\ell_1^2, \mathcal{X}[1]), \dots, (\ell_x^{nb-x}, \mathcal{X}[nb-x]), (\ell_1^1, \mathcal{X}[1])$, etc. Consequently, $(\ell x_i, X_i)$ takes each values $(\ell_\alpha^\beta, \mathcal{X}[\beta]), 1 \leq \alpha \leq x, 1 \leq \beta \leq nb-x$ infinitely often. In particular, p_i executes $(\ell x_i, X_i) \leftarrow Next(\ell, X)$ infinitely often. But this contradicts the *Claim C1* that states that the number of $x_MOVE(\ell, X)$ messages that are sent is finite. It follows that there are a pair (λ_i, σ_i) and a time τ_i such that $\forall \tau \geq \tau_i : (\ell x_i^\tau, X_i^\tau) = (\lambda_i, \sigma_i)$.

Claim C1: There is a pair (ℓ, X) such that the number of $x_MOVE(\ell, X)$ messages that are sent is finite.

Proof of Claim C1. We consider two cases according to the number f of actual process crashes.

- Case 1: $f \geq x$. Let X be a set of x processes that are faulty and ℓ be the identity of an arbitrary process in X . As only processes that belongs to X can send $x_MOVE(\ell, X)$ messages, it follows from the fact all these processes eventually stop taking steps that the number of $x_MOVE(\ell, X)$ messages sent is finite.
- Case 2: $f < x$. Due to the limited scope eventual accuracy property of the class $\diamond\mathcal{S}_x$, there are a set $X \subseteq \Pi$ of size x and a correct process $p_\ell \in X$ such that, after some time τ , no process that belongs to the set X suspects p_ℓ . Since (1) only process that belongs to X can send $x_MOVE(\ell, X)$ messages, and, (2) a process $p_i \in X$ broadcasts a $x_MOVE(\ell, X)$ message only if $\ell \in suspected_i$, it follows that after time τ , no message $x_MOVE(\ell, X)$ can be broadcast, which implies that the number of such messages is finite. *End of the Proof of Claim C1*.

□ *Lemma 4*

Corollary 1 *The protocol is quiescent (i.e., eventually all the processes stop sending x_MOVE messages).*

Proof Let us assume (for contradiction) that there is a correct process p_i that never stop sending x_MOVE messages. Due to Lemma 4, there is a time τ after which $(\ell x_i, X_i)$ remains permanently equal to the constant pair (λ_i, σ_i) . Consequently, after time τ , p_i keeps on broadcasting $x_MOVE(\lambda_i, \sigma_i)$. It follows then from the validity and termination properties of the reliable broadcast primitive that there is a time $\tau' > \tau$ at which p_i executes $(\ell x_i, X_i) \leftarrow Next(\lambda_i, \sigma_i)$, contradicting Lemma 4. □ *Corollary 1*

Lemma 5 $\forall i, j \in C : (\lambda_i, \sigma_i) = (\lambda_j, \sigma_j)$. (*In the following, (λ, σ) denotes that pair.*)

Proof Due to the properties of the reliable broadcast primitive, p_i and p_j deliver the same multiset of $x_MOVE(\ell, X)$ messages. Moreover, it follows from Corollary 1 that this multiset is finite. Due to the fact that p_i and p_j consume the messages according to the same ring order, and the fact that the common multiset of delivered messages is finite, it follows that $(\lambda_i, \sigma_i) = (\lambda_j, \sigma_j)$. □ *Lemma 5*

Lemma 6 $(\sigma \cap C \neq \emptyset) \Rightarrow (\lambda \in C)$.

Proof Let us assume (by contradiction) that $\sigma \cap C \neq \emptyset$ and λ is the identity of a faulty process. Let p_i be a process that belongs to $\sigma \cap C$. Due to the strong completeness property of the class $\diamond\mathcal{S}_x$, it exists a time τ_1 after which the local predicate $\lambda \in \text{suspected}_i$ remains permanently satisfied. Moreover, it follows from lemmas 4 and 5 that, from some time τ_i , the predicate $(\ell x_i, X_i) = (\lambda, \sigma)$ remains permanently true. There is consequently a time $\tau \geq \max(\tau_1, \tau_i)$ at which p_i broadcasts a message $\text{x_MOVE}(\lambda, \sigma)$. When p_i delivers this message, it executes $(\ell x_i, X_i) \leftarrow \text{Next}(\lambda, \sigma)$, contradicting Lemma 4. $\square_{\text{Lemma 6}}$

Theorem 7 *The algorithm described in Figure 7 ensures the existence of a set X and a time τ such that $\forall \tau' \geq \tau$, the following holds:*

1. $|X| = x$,
2. $i \in \Pi - X \Rightarrow \text{repr}_i = i$,
3. $\forall i, j \in X \cap C : \text{repr}_i = \text{repr}_j = \rho \in C \cap X$.

Proof Let $\tau = \max\{\tau_i : i \in \Pi\}$ where τ_i is the time introduced in Lemma 4, and σ and λ be the set and the process identity defined in Lemma 5. Let us first observe that due to its definition (σ is a set X_i) we have $|\sigma| = x$ (Item 1). Let p_i be a correct process. If $i \in \Pi - X$, then as the value of repr_i does not change after time τ (Lemma 4 and Task T1), it follows that $\text{repr}_i = i$ is permanently true from time τ (Item 2). Moreover, it directly follows from Lemma 5 and task T1 that all the correct processes p_j belonging to the set σ have permanently the same representative $\text{repr}_j = \lambda$ from time τ . Finally, due to Lemma 6, λ is the identity of a correct (Item 3). Taking $X = \sigma$, $\tau = \max\{\tau_i : i \in \Pi\}$ and $\rho = \lambda$ completes the proof of the theorem. $\square_{\text{Theorem 7}}$

5.2 The Upper Wheel Component

5.2.1 Principles and description

The “upper wheel” component consists of four tasks T3-T6 (Figure 8)³. Similarly to the lower wheel component, it uses a sequence, that we call \mathcal{L} , including all the possible sets of size $z = (t + 2) - (x + y)$ generated from the n processes composing the system. \mathcal{L} is known by all the processes. Let nb_L be the length of this sequence, and $\mathcal{L}[k]$ its k th element. The function $\text{Next}(\mathcal{L}[k])$ returns $\mathcal{L}[k + 1]$ when $k < \text{nb}_L$, and $\mathcal{L}[1]$ when $k = \text{nb}_L$.

```

Init:  $L_i \leftarrow \mathcal{L}[1]$ 

Task T3:
(0 1) repeat forever
(0 2)   Broadcast INQUIRY();
(0 3)   wait until ( corresponding RESPONSE() received from  $\geq n - \text{nb}_{c_i}$  processes )
                                     %  $\text{nb}_{c_i}$  can dynamically change
(0 4)   let  $\text{rec\_from}_i = \{id_j \text{ received previously at line 03}\}$ ;
(0 5)   if ( $\text{rec\_from}_i \cap L_i = \emptyset$ ) then R_Broadcast L_MOVE( $L_i$ ) end_if
(0 6) end_do

Task T4: when L_MOVE( $L_i$ ) is R_delivered: ( $L_i$ )  $\leftarrow$  Next( $L_i$ )

Task T5: when INQUIRY() is received from  $p_j$ : send RESPONSE( $\text{repr}_i$ ) to  $p_j$ 

Task T6: when  $\text{trusted}_i$  is read by the upper layer: return( $L_i$ )

```

Figure 8: From $\diamond\psi^y + \diamond\mathcal{S}_x$ to Ω^z : upper wheel component (code for p_i)

³A version of this component, based on $\diamond\phi^y$, is described in [21]. It is much more involved than the one presented in Figure 8.

The transformation, described in Figure 8, relies on the following principles. (Let us recall that nb_c_i is the read-only local variable that p_i is provided with by the underlying failure detector of the class $\diamond\psi^y$.) The aim is for p_i to compute the value of the set $trusted_i$ provided to the upper layer (Task $T6$), namely, a set of z processes that eventually includes (at least) one correct process. So, starting from the set $L_i = \mathcal{L}[1]$, the processes scan (in the same order) the infinite sequence of sets $\mathcal{L}[1], \mathcal{L}[2], \dots, \mathcal{L}[nb_L], \mathcal{L}[1], \dots$ (tasks $T3$ and $T4$). When p_i is working with a set L_i , it proceeds as follows.

- First, p_i strives to know if L_i contains a correct process. To that end, it repeatedly broadcasts an inquiry message (task $T3$, line 02). When a process p_j receives such a message it sends back to p_i the identity of its representative as defined by the lower wheel component (task $T5$).
- Then, p_i waits for responses from $n - nb_c_i$ processes. Let us observe that, as eventually $nb_c_i = \max(t - y, f)$ (where f is the number of faulty processes in the considered run), p_i eventually receives $n - \max(t - y, f)$ response messages (the value nb_c_i provided by the failure detector of the class $\diamond\psi^y$ is repeatedly read until the waiting condition becomes true).
- Finally, when it has received enough responses, p_i defines rec_from_i as the set of processes from which responses have been received (line 04). If one of these processes belongs to the current set L_i , p_i keeps the current value of L_i . Otherwise, it considers that the processes of L_i are faulty, and broadcasts consequently a message $L_MOVE(L_i)$ to inform all the processes that they has to proceed to the next set for L_i .

To capture the intuition that underlies the fact that the two wheels synchronize and the processes stabilize on the same set L , let us observe that, due to the property eventually ensured on the $repr_j$ local variables by the lower wheel component, there is a time after which all the $RESPONSE(id)$ messages carry identities of correct processes. It follows that if the set L_i currently investigated by the processes does not change, that set includes at least one correct process and we have obtained the property required by $trusted_i$.

5.2.2 Proof of the upper wheel component

The proof is very similar to the proof of the lower wheel algorithm. Its structure is the same, and some of its parts are also the same. This is a direct consequence of the fact that both components are based on the same “wheel” principle. The proof considers an arbitrary run of the algorithm. As before, C denotes the set of processes that are correct in that run, and var_i^τ denotes the value of the local variable var_i of at time τ .

Lemma 7 $\forall i \in C$, there is a set Λ_i and a time τ_i such that $\forall \tau \geq \tau_i : L_i^\tau = \Lambda_i$.

Proof We claim (*Claim C2*) that there is a set L such that the number of $L_MOVE(L)$ messages that are sent is finite. This claim, used to prove the lemma, is proved later.

Let p_i be a correct process and let us assume (by way of contradiction) that there is no set Λ_i such that after some time $L_i = \Lambda_i$ remains true forever. It follows from the way that each p_i updates its local variable L_i , that the sequence of successive values taken by each L_i is $\mathcal{L}[1], \mathcal{L}[2], \dots, \mathcal{L}[nb_L], \mathcal{L}[1], \dots$ ⁴. Consequently, L_i takes each value $\mathcal{L}[\alpha], 1 \leq \alpha \leq nb_L$ infinitely often. In particular, p_i executes $L_i \leftarrow Next(L)$ infinitely often. Since this occurs when p_i delivers a $L_MOVE(L)$ message, this contradicts the Claim *C2* that states that a finite number of such messages are sent. It follows that there is a set Λ_i and a time τ_i such that $\forall \tau \geq \tau_i : L_i^\tau = \Lambda_i$.

Claim C2: There is a set L such that the number of $L_MOVE(L)$ messages that are sent is finite.

Proof of Claim C2.

Let us consider the time τ at which the lower wheel stops turning. More precisely, there is a time τ , a set $X \subseteq \Pi, |X| = x$ and a process identity $\lambda \in X$ (Theorem 7) such that:

1. $\forall i \in \Pi - X, \forall \tau' \geq \tau : repr_i^{\tau'} = i$ and,

⁴This follows from the fact that each process visits the sets of \mathcal{L} according to the same deterministic order defined from a logical ring, as in Figure 6, where $\mathcal{X}[\beta]$ is replaced by $\mathcal{L}[\beta'], 1 \leq \beta' \leq nb_L$.

2. (a) $X \cap C \neq \emptyset$: $\exists \lambda \in C \cap X$ such that $\forall i \in X, \forall \tau' \geq \tau : repr_i^{\tau'} = \lambda$ or,
- (b) $X \cap C = \emptyset$: all processes that belong to X have crashed by time τ .

Let us consider a set L of $z = (t+2) - (x+y)$ processes defined as follows (see Figure 9): (1) $|X \cap L| = 1$, (2) if $X \cap C \neq \emptyset$ then, $X \cap L = \{\lambda\}$ and (3) L contains the identity of a correct process. It is easy to see that such a set L does exist. Moreover, let us observe that there is $\ell, \ell \in L$, such that p_ℓ is a correct process and eventually $repr_\ell = \ell$.

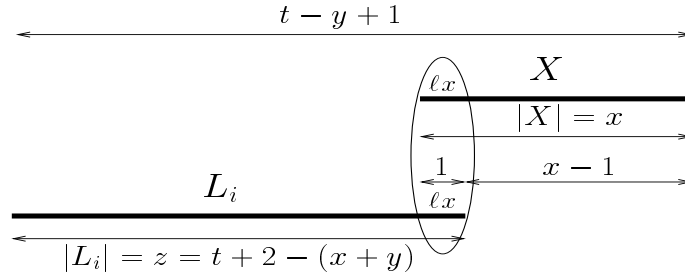


Figure 9: When the upper wheel stops looking for a new L_i set

We examine two cases according to the actual number f of process crashes. In each case, we show that, after some time defined by the case assumption, no $L_MOVE(L)$ message is sent.

- Case 1: $f \geq t - y + 1$. Due to the eventual convergence property of the class $\diamond\psi^y$, there is a time τ' after which $nb_c_i = n - f$ remains forever true at each correct process p_i . Let τ'' be a time at which the f faulty processes have crashed and the messages they sent to the correct processes have been received and processed.

Let $\tau_0 = \max(\tau, \tau', \tau'')$, i.e., after τ_0 , no process crashes and the outputs of both the lower wheel component and the $\diamond\psi^y$ failure detector do no longer change. Let p_i be a correct process. After time τ_0 , each time p_i updates rec_from_i , we have $rec_from_i = C$ (this is because, after τ_0 , p_i waits for $n - f$ response messages and the f processes that are faulty have crashed before τ_0).

As L contains the identity of a correct process p_ℓ such that $repr_\ell = \ell$, it follows that $L \cap rec_from_i \neq \emptyset$. Consequently, no message $L_MOVE(L)$ can be sent after time τ_0 , which implies that the number of these messages is finite.

- Case: $f < t - y + 1$.

In that case, due to the eventual convergence property of the class $\diamond\psi^y$, there is a time τ' after which at each process p_i , $nb_c_i = t - y$ remains forever true. Let τ_0 be a time after which the outputs of both the failure detector of the class $\diamond\psi^y$ and the lower wheel component do not change at each process.

Let us consider an execution of the repeat loop started after τ_0 by a correct process p_i . We first show that after p_i has updated rec_from_i at line 04, there is $j \in L \cup X$ such that $repr_j \in rec_from_i \cap L$. To update rec_from_i , p_i waits for $n - nb_c_i = n - (t - y)$ responses. Moreover, due to the definition of L , we have $|L \cup X| = |L| + |X| - 1 = 1 + (t - y)$. Consequently, among the $n - (t - y)$ responses taken into account by p_i to update rec_from_i , there is a response sent by a process p_j such that $j \in L \cup X$. We show that $repr_j \in L$. If $j \in X$, $repr_j = \lambda \in L$. Otherwise, $j \in L - X$, from which we have $repr_j = j \in L$.

Hence, after time τ_0 , a process that is waiting for responses always receives such a message from a process p_j that belongs to $L \cup X$ and this message carries a process identity $repr_j$ such that $repr_j \in L$. It then follows from lines 04-05 that, after some time, no process can broadcast a message $L_MOVE(L)$.

End of the Proof of Claim C2.

□ Lemma 7

Corollary 2 *Eventually all processes stop sending L_MOVE messages.*

Proof Let us assume (by contradiction) that it exists a correct process p_i that never stops sending `L_MOVE` messages. Due to Lemma 7, there is a time τ_i after which L_i remains permanently equal to the constant set Λ_i . Consequently, after time τ_i , p_i keeps on broadcasting `L_MOVE`(Λ_i). It follows then from the validity and the termination properties of the reliable broadcast primitive that there is a time $\tau' > \tau_i$ at which p_i executes $L_i \leftarrow \text{Next}(\Lambda_i)$, contradicting Lemma 7. $\square_{\text{Corollary 2}}$

Remark. The fact that there is a time after which no `L_MOVE`(L) messages are exchanged, does not imply that the algorithm is quiescent. This is because the correct processes keep on sending forever `INQUIRY`() messages, and answering them by sending back `RESPONSE`() messages. Differently, the lower wheel component uses only `X_MOVE`() messages.

Lemma 8 $\forall i, j \in C : \Lambda_i = \Lambda_j$. (In the following, Λ denotes that set.)

Proof Due to the properties of the reliable broadcast primitive, p_i and p_j deliver the same multiset of `L_MOVE`(L) messages. Moreover, it follows from Corollary 2 that this multiset is finite. Due to the fact that p_i and p_j consume the messages according to the same ring order, and the fact that the common multiset of delivered messages is finite, it follows that $\Lambda_i = \Lambda_j$. $\square_{\text{Lemma 8}}$

Theorem 8 The sets trusted_i implemented by the algorithm described in Figure 8 satisfy the property defining the class Ω^z .

Proof Due to Lemma 8, there is a time after which all the processes have permanently the same set Λ , $|\Lambda| = z = t + 2 - (x + y)$. It remains to show that $\Lambda \cap C \neq \emptyset$.

Let us assume for contradiction that $\Lambda \cap C = \emptyset$. Let p_i be a correct process. Due to the properties ensured by the lower wheel (Theorem 7), there is a time after which any message `RESPONSE`(*repr*) contains the identity of a correct process. From the assumption that Λ contains only faulty processes, it follows that there is a time τ_1 after which p_i cannot receive a `RESPONSE` message that carries the identity of a process belonging to Λ . Moreover, there is a time τ_i after which the predicate $L_i = \Lambda$ is permanently true (Lemma 7). Consequently, there is a time $\tau \geq \max(\tau_1, \tau_i)$ at which the predicate in the **if** statement of line 05 is not satisfied (i.e., at time τ , we have $\text{rec_from}_i \cap \Lambda = \emptyset$). It follows then that p_i broadcasts a message `L_MOVE`(Λ). When p_i delivers such a message, it executes $L_i \leftarrow \text{Next}(\Lambda)$. The fact that this occurs after the time τ_i contradicts Lemma 7. $\square_{\text{Theorem 8}}$

6 Lower Bounds and (Ir)Reducibility Results

This section states first a lower bound related to the addition of failure detector classes (Figure 2). It then proves the (ir)reducibility results stated in the grid depicted in Figure 1. As the classes ψ^y and ϕ^y ($\diamond\psi^y$ and $\diamond\phi^y$) are equivalent (section 3), we sometimes use ϕ^y ($\diamond\phi^y$) instead of ψ^y ($\diamond\psi^y$) in the proofs.

6.1 A Lower bound when Adding $\diamond\mathcal{S}_x$ and $\diamond\psi^y$

This section shows that $(x + y + z > t + 1)$ is a lower bound when one wants to add failure detectors of the class $\diamond\mathcal{S}_x$ and $\diamond\psi^y$ to build a failure detector of the class Ω^z .

Theorem 9 Let us consider any system $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x, \diamond\psi^y]$. $(\diamond\mathcal{S}_x + \diamond\psi^y \rightsquigarrow \Omega^z) \Leftrightarrow (x + y + z > t + 1)$.

Proof [\Leftarrow part] This part follows directly from the two wheels algorithm previously described in Sections 5.1.1 and 5.2.1, and proved in sections 5.1.2 and 5.2.2.

[\Rightarrow part] The proof of this part is by contradiction and considers the stronger system $\mathcal{AS}_{n,t}[\mathcal{S}_x, \psi^y]$. As $\mathcal{S}_x \subseteq \diamond\mathcal{S}_x$ and $\psi^y \subseteq \diamond\psi^y$, any impossibility result established in $\mathcal{AS}_{n,t}[\mathcal{S}_x, \psi^y]$ holds in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x, \diamond\psi^y]$.

Let us assume that there is an algorithm \mathcal{T} that builds a failure detector of the class Ω^z in $\mathcal{AS}_{n,t}[\mathcal{S}_x, \psi^y]$ with $x + y + z \leq t + 1$. The contradiction is based on the following observations:

- **Observation O1:** Let f be the number of actual failures. When $f \leq t - y$, the only information that a failure detector of the class ψ^y can provide is the fact that the number of failures is $\leq t - y$.
Proof of O1. Consider a run where $f \leq t - y$. Let $E \subseteq \Pi$. Due to the safety property of the class ψ^y , at each process p_i , the value of $nb_{\mathcal{C}_i}$ is always $t - y$. Consequently the value of $nb_{\mathcal{C}_i}$ does not depend on which processes has crashed. *End of the Proof of O1.*
- **Observation O2:** There is no algorithm that solves the k -set agreement problem in $\mathcal{AS}_{n,t}[\mathcal{S}_x]$ when $t \geq k + x - 1$.
Proof of O2. This is a lower bound for solving the k -set agreement problem in $\mathcal{AS}_{n,t}[\mathcal{S}_x]$ established in [14]. *End of the Proof of O2.*

Let us now consider the transformation \mathcal{T} . In any run where $f \leq t - y$, it follows from O1 that \mathcal{T} can rely on ψ^y only to know that the number of failures is $\leq t - y$. This implies that \mathcal{T} can be used to build a failure detector of the class Ω^z in $\mathcal{AS}_{n,t-y}[\mathcal{S}_x]$. Moreover, it exists an algorithm \mathcal{A} that solves the z -set agreement problem in $\mathcal{AS}_{n,t-y}[\Omega^z]$ (such an algorithm is described in Section 4). Consequently, by combining transformation \mathcal{T} and algorithm \mathcal{A} , one can solve the z -set agreement problem in $\mathcal{AS}_{n,t-y}[\mathcal{S}_x]$. Hence, it follows from O2 that the constraint $t - y < z + x - 1$ has to be satisfied, from which we obtain $x + y + z > t + 1$: a contradiction. $\square_{\text{Theorem 9}}$

The following corollary is an immediate consequence of the proof of Theorem 9.

Corollary 3 *Let us consider any system $\mathcal{AS}_{n,t}[\mathcal{S}_x, \psi^y]$, $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x, \psi^y]$ or $\mathcal{AS}_{n,t}[\mathcal{S}_x, \diamond\psi^y]$. In any of these systems, it exists an algorithm that builds a failure detector of the class Ω^z if and only if $(x + y + z) > t + 1$.*

The following corollary is a consequence of Theorem 9.

Corollary 4 *The two wheels algorithm described in Figures 7 and 8 is optimal with respect to the possible values of x , y and z .*

As $\diamond\mathcal{S}_1$ (case $x = 1$) provides no information on failures, we directly obtain the following corollary from the two wheel algorithm and Theorem 9.

Corollary 5 *It is possible to build a failure detector of the class Ω^z in $\mathcal{AS}_{n,t}[\psi^y]$ or $\mathcal{AS}_{n,t}[\diamond\psi^y]$ if and only if $y + z > t$.*

Similarly, as $\diamond\psi^0$ (case $y = 0$) provides no information on failures, we also have:

Corollary 6 *It is possible to build a failure detector of the class Ω^z in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x]$ if and only if $x + z > t + 1$.*

6.2 Relations between $\mathcal{S}_x/\diamond\mathcal{S}_x$ and $\psi^y/\diamond\psi^y$

Theorem 10 *Let $1 \leq x \leq t + 1$ and $1 \leq y \leq t$. It is not possible to build a failure of the class ψ^y or $\diamond\psi^y$ in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x]$ or in $\mathcal{AS}_{n,t}[\mathcal{S}_x]$.*

Proof For convenience, the result is proved using the classes $\phi^y/\diamond\phi^y$. The proof considers the “stronger” system $\mathcal{AS}_{n,t}[\mathcal{S}_x]$. As $\mathcal{S}_x \subseteq \diamond\mathcal{S}_x$, the proof remains valid for a system $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x]$. Similarly, as $\phi^y \subseteq \diamond\phi^y$ the proof considers only the “weaker” class $\diamond\phi^y$. The proof is by contradiction. Let us assume that there is a failure detector \mathcal{F} of the class \mathcal{S}_x and an algorithm \mathcal{A} that transforms \mathcal{F} into a failure detector of the class $\diamond\phi^y$. We exhibit a run R in which the eventual safety property of the class $\diamond\phi^y$ is not satisfied.

Let $E \subseteq \Pi$, $|E| = t - y + 1$ and $E \cap C \neq \emptyset$. Let p_c be a correct process that does not belong to set E . Moreover, p_c is never suspected by \mathcal{F} in run R . Let τ_0 be the time at which any $\text{QUERY}(E)$ invoked after time τ_0 returns the value *false*. Such a time exists due to the correctness of algorithm \mathcal{A} and the eventual safety property of the class $\diamond\phi^y$. We consider two runs $R1$ and $R1'$ defined as follows:

- Runs $R1$ and R are indistinguishable by all processes until time τ_0 . A time $\tau_0 + 1$, all processes that belong to E crash. Let $\tau_1 > \tau_0$ be a time at which a process $p_i \in \Pi - E$ invokes $\text{QUERY}(E)$ and obtains the value *true*. Such a time must exist due to liveness property of the class $\diamond\phi^y$.

- Runs $R1'$ and R are indistinguishable by all processes until time τ_0 . In the run $R1'$, all the processes in E are correct, but all the messages they send between times $\tau_0 + 1$ and τ_1 are delayed until time $\tau_1 + 1$.

Moreover, both runs $R1$ and $R1'$ are such that the outputs of the failure detector \mathcal{F} , at each process, are exactly the same between the times 0 and τ_1 . (Let us notice that whatever the output of \mathcal{F} in $R1$, the output of \mathcal{F} can be exactly the same in $R1'$ without violating the properties of the class \mathcal{S}_x . As p_c is correct in $R1$ and $R1'$ and never suspected in $R1$ and $R1'$, limited scope perpetual accuracy is insured. Since strong completeness is an eventual property, it is always satisfied in any finite prefix of any execution.) Clearly, up to time τ_1 , the processes that belong to $\Pi - E$ cannot distinguish the run $R1$ from the run $R1'$. It follows that, in the run $R1'$, an invocation of $\text{QUERY}(E)$ by p_i at time $\tau_1 > \tau_0$ returns the value *true*. But in run $R1'$, $\text{QUERY}(E)$ issued after time τ_0 must return the value *false*: a contradiction. $\square_{\text{Theorem 10}}$

Theorem 11 *Let $0 \leq y < t$ and $1 < x \leq t + 1$. It is not possible to build a failure detector of the class \mathcal{S}_x or $\diamond\mathcal{S}_x$ neither in $\mathcal{AS}_{n,t}[\diamond\psi^y]$ nor in $\mathcal{AS}_{n,t}[\psi^y]$.*

Proof Let us first notice that we need to prove only the impossibility to build a failure detector of the class $\diamond\mathcal{S}_x$ in $\mathcal{AS}_{n,t}[\psi^y]$. The proof is by contradiction and uses the following observations.

- Observation *O1*: Let f be the number of actual failures. When $f \leq t - y$, the only information that a failure detector of the class ψ^y can provide is the fact that the number of failures is $\leq t - y$. (This observation has already been stated and proved in Theorem 9.)
- Observation *O2*: There are algorithms that solve the k -set agreement problem in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x]$. All these algorithms require $t \leq k + x - 2$. (Examples of such algorithms can be found in [14, 24]. The lower bound on t is established in [14].)
- Observation *O3*: The k -set agreement problem can be solved in $\mathcal{AS}_{n,t-y}[\emptyset]$ if and only if $k > t - y$. (The proof of this observation constitutes an important result of fault-tolerant distributed computing. It can be found in [1, 16, 28].)

Let us suppose that there is an algorithm \mathcal{A} that builds a failure detector of the class $\diamond\mathcal{S}_x$ from a failure detector of the class ψ^y . In any run where $f \leq t - y$, it follows from *O1* that \mathcal{A} can rely on ψ^y only to know that the number of failures is $\leq t - y$. Consequently, \mathcal{A} can build a failure detector of the class $\diamond\mathcal{S}_x$ in a system $\mathcal{AS}_{n,t-y}[\emptyset]$. This means that one can use \mathcal{A} to solve the k -set agreement problem with $k = (t - y) - x + 2$, using any algorithm listed in observation *O2* in a system $\mathcal{AS}_{n,t-y}[\emptyset]$. We then conclude from *O3* ($k > t - y$) that $(t - y) - x + 2 > t - y$, i.e., $x \leq 1$, a contradiction with the assumption $1 < x \leq n$ ⁵.

$\square_{\text{Theorem 11}}$

6.3 From Ω^z to $\psi^y/\diamond\psi^y$ or $\mathcal{S}_x/\diamond\mathcal{S}_x$

It has been shown (Corollaries 5 and 6) that it is possible to build a failure detector of the class Ω^z from any failure detector of the classes $\psi^y/\diamond\psi^y$ (resp., $\mathcal{S}_x/\diamond\mathcal{S}_x$) if and only if $x + z > t + 1$ (resp., $y + z > t$). This section shows that it is not possible to build a failure detector of the classes $\psi^y/\diamond\psi^y$ (resp., $\mathcal{S}_x/\diamond\mathcal{S}_x$) from any failure detector of the class Ω^z . The proofs of these impossibilities are based on Theorem 10 and 11.

Theorem 12 *Let $1 \leq y \leq t$ and $1 \leq z \leq t + 1$. It is impossible to build a failure detector of a class $\psi^y/\diamond\psi^y$ in $\mathcal{AS}_{n,t}[\Omega^z]$.*

Proof The proof is by contradiction. Let us assume that there is an algorithm \mathcal{A} that builds a failure detector of a class $\diamond\psi^y$, $1 \leq y \leq t$, from any failure detector of a class Ω^z , $1 \leq z \leq t + 1$. Due to Corollary 6, it is possible to build a failure detector of a class Ω^z in $\mathcal{AS}_{n,t}[\diamond\mathcal{S}_x]$ when $x + z > t + 1$. Combining this construction with the algorithm \mathcal{A} we obtain an algorithm \mathcal{B} that builds a failure detector of the class ψ^y ,

⁵Let us remind that the failure detectors of the classes \mathcal{S}_1 and $\diamond\mathcal{S}_1$ provide no information on failures.

$1 \leq y \leq t$ from a failure detector of the class $\diamond\mathcal{S}_x$. But such an algorithm \mathcal{B} contradicts Theorem 10 that states that there is no such algorithm when $1 \leq x \leq t + 1$ and $1 \leq y \leq t$. $\square_{\text{Theorem 12}}$

Theorem 13 *Let $1 < x, z \leq t$. It is impossible to build a failure detector of the class $\mathcal{S}_x / \diamond\mathcal{S}_x$ in $\mathcal{AS}_{n,t}[\Omega^z]$.*

Proof The proof is similar to the proof of Theorem 12. It is left to the reader. $\square_{\text{Theorem 13}}$

6.4 Optimality in the Grid of Figure 1

It follows from all the previous theorems and lemmas that, when we consider all the failure detector classes depicted in Figure 1, Ω^k is the weakest class that allows solving the k -set agreement problem. This constitutes a first step towards the characterization of the weakest failure detector class for solving that problem. A corresponding Ω^k -based k -set agreement protocol has been described in Section 4.

7 Conclusion

Considering two objects of two types, $O1$ that allows solving the $k1$ -set agreement problem and does not allow solving the $(k1 - 1)$ -set agreement problem, and $O2$ that allows solving the $k2$ -set agreement problem and does not allow solving the $(k2 - 1)$ -set agreement problem, is it possible to combine them so as to solve a stronger version of the k -set agreement problem, i.e., such that $k < \min(k1, k2)$?

Considering the previous question as a guideline, and base objects that are failure detectors, the paper has investigated three families of failure detector classes, namely, $(\diamond\mathcal{S}_x)_{1 \leq x \leq n}$, $(\diamond\psi_{0 \leq x \leq n}^y)$ and $(\Omega^z)_{1 \leq z \leq n}$. Among these failure detector classes, it has shown which ones are equivalent and which ones are not. As an example, the paper has shown that any class in the sub-family $(\diamond\mathcal{S}_x)_{t < x \leq n}$ and the class Ω_1 are equivalent (given any failure detector of one class, it is possible to build a failure detector of the other class). It has also shown that it is impossible to build a failure detector of the class $(\diamond\mathcal{S}_x)_{1 < x \leq n}$ from a failure detector of any class in the sub-family $(\Omega_{1 < z \leq n})$. A main result of the paper is the theorem “ $\diamond\mathcal{S}_x + \diamond\psi^y \rightsquigarrow \Omega^z \Leftrightarrow x + y + z > t + 1$ ” that states that it is possible to combine any failure detector of the class $\diamond\mathcal{S}_x$ and any failure detector of the class $\diamond\psi_y$ to obtain a failure detector belonging to the class Ω^z where $z = (t + 2) - (x + y)$.

The paper has also presented a k -set agreement protocol for message-passing asynchronous systems equipped with Ω^k , and established that the resilience bound $t < n/2$ and the failure detector bound $z \leq k$ are tight for such systems.

The theorem “ $\diamond\mathcal{S}_x + \diamond\psi^y \rightsquigarrow \Omega^z \Leftrightarrow x + y + z > t + 1$ ” shows that, in a system equipped with failure detectors of both classes $\diamond\mathcal{S}_x$ and $\diamond\psi_y$, these failure detector classes are not robust. Their combination allows solving the k -set agreement problem with $z = (t + 2) - (x + y)$, while each of them taken separately cannot. Apparently, this seems to contradict the results on base object composition stated in [2] and [15]. There is no contradiction: both these papers consider base objects that have a sequential specification (and are consequently linearizable), while our base objects are failure detectors that have no sequential specification. This shows an interesting difference according to the fact that the base objects have or not a sequential specification.

References

- [1] Borowsky E. and Gafni E., Generalized FLP Impossibility Results for t -Resilient Asynchronous Computations. *Proc. 25th ACM Symposium on the Theory of Computing (STOC'93)*, ACM Press, pp. 91-100, 1993.
- [2] Borowsky E. and Gafni E., The Implication of the Borowsky-Gafni Simulation on the Set Consensus Hierarchy. *Technical Report 93-0021*, Computer Science Department, University of California at Los Angeles, 1993.
- [3] Chandra T., Hadzilacos V. and Toueg S., The Weakest Failure Detector for Solving Consensus. *Journal of the ACM*, 43(4):685-722, 1996.

- [4] Chandra T.D. and Toueg S., Unreliable Failure Detectors for Reliable Distributed Systems. *Journal of the ACM*, 43(2):225-267, 1996.
- [5] Chaudhuri S., More *Choices* Allow More *Faults*: Set Consensus Problems in Totally Asynchronous Systems. *Information and Computation*, 105:132-158, 1993.
- [6] Chu F., Reducing Ω to $\diamond W$. *Information Processing Letters*, 76(6):293-298, 1998.
- [7] Dutta P. and Guerraoui R., Fast Indulgent Consensus with Zero Degradation. *Proc. 4th European Dependable Computing Conference (EDCC'02)*, Springer-Verlag LNCS #2485, pp. 191-208, 2002.
- [8] Delporte-Gallet C., Fauconnier H. and Guerraoui R., (Almost) All Objects are Universal in Message Passing Systems. *Proc. 19th Symposium on Distributed Computing (DISC'05)*, Springer Verlag LNCS #3724, pp. 184-198, 2005.
- [9] Fischer M.J., Lynch N. and Paterson M.S., Impossibility of Distributed Consensus with One Faulty Process. *Journal of the ACM*, 32(2):374-382, 1985.
- [10] Guerraoui R., Non-Blocking Atomic Commit in Asynchronous Distributed Systems with Failure Detectors. *Distributed Computing*, 15:17-25, 2002.
- [11] Guerraoui R. and Raynal M., The Information Structure of Indulgent Consensus. *IEEE Transactions on Computers*. 53(4), 53(4):453-466, 2004.
- [12] Guerraoui R. and Schiper A., Gamma-accurate Failure Detectors. *Proc. 10th Workshop on Distributed Algorithms (WDAG'96)*, Springer Verlag LNCS #1151, pp. 269-286, 1996.
- [13] Hadzilacos V. and Toueg S., Reliable Broadcast and Related Problems. In *Distributed Systems*, ACM Press, New-York, pp. 97-145, 1993.
- [14] Herlihy M.P. and Penso L. D., Tight Bounds for k -Set Agreement with Limited Scope Accuracy Failure Detectors. *Distributed Computing*, 18(2):157-166, 2005.
- [15] Herlihy M.P. and Rajsbaum S., Set Consensus using Arbitrary Objects. *Proc. 13th ACM Symposium on Principles of Distributed Computing (PODC'94)*, ACM Press, pp. 324-333, 1994.
- [16] Herlihy M.P. and Shavit N., The Topological Structure of Asynchronous Computability. *Journal of the ACM*, 46(6):858-923, 1999.
- [17] Lamport L., The Part-Time Parliament. *ACM Transactions On Computer Systems*, 16(2):133-169, 1998.
- [18] Mostefaoui A., Rajsbaum S. and Raynal M., Conditions on Input Vectors for Consensus Solvability in Asynchronous Distributed Systems. *Journal of the ACM*, 50(6):922-954, 2003.
- [19] Mostefaoui A., Rajsbaum S. and Raynal M., The Combined Power of Conditions and Failure Detectors to Solve Asynchronous Set Agreement. *Proc. 24th ACM Symposium on Principles of Distributed Computing (PODC'05)*, ACM Press, pp. 179-188, 2005.
- [20] Mostefaoui A., Rajsbaum S., Raynal M. and Travers C., From $\diamond W$ to Ω : a Simple Bounded Quiescent Reliable broadcast-based Transformation. *Journal of Parallel and Distributed Computing*. To appear, 2007.
- [21] Mostefaoui A., Rajsbaum S., Raynal M. and Travers C., Irreducibility and Additivity of Set Agreement-oriented Failure Detector Classes (Extended Abstract). *Proc. 25th ACM Symposium on Principles of Distributed Computing (PODC'06)*, ACM Press, pp. 153-162, 2006.
- [22] Mostefaoui A. and Raynal M., Solving Consensus Using Chandra-Toueg's Unreliable Failure Detectors: a General Quorum-Based Approach. *Proc. 13th Symposium on Distributed Computing (DISC'99)*, Springer Verlag LNCS #1693, pp. 49-63, 1999.
- [23] Mostefaoui A. and Raynal M., Unreliable Failure Detector with Limited Scope Accuracy and an Application to Consensus. *Proc. 19th Int'l Conference on Foundations of Software Technology and Theoretical Computer Science (FST&TCS'99)* Springer Verlag LNCS #1738, pp. 329-340, 1999.
- [24] Mostefaoui A. and Raynal M., k -Set Agreement with Limited Accuracy Failure Detectors. *Proc. 19th ACM Symposium on Principles of Distributed Computing (PODC'00)*, ACM Press, pp. 143-152, 2000.

- [25] Mostéfaoui A. and Raynal M., Leader-Based Consensus. *Parallel Processing Letters*, 11(1):95-107, 2001.
- [26] Neiger G., Failure Detectors and the Wait-free Hierarchy. *Proc. 14th ACM Symposium on Principles of Distributed Computing (PODC'95)*, ACM Press, pp. 100-109, 1995.
- [27] Raynal M., A Short Introduction to Failure Detectors for Asynchronous Distributed Systems. *ACM SIGACT News, Distributed Computing Column*, 36(1):53-70, 2005.
- [28] Saks M. and Zaharoglou F., Wait-Free k -Set Agreement is Impossible: The Topology of Public Knowledge. *SIAM Journal on Computing*, 29(5):1449-1483, 2000.
- [29] Schiper A., Early Consensus in an Asynchronous System with a Weak Failure Detector. *Distributed Computing*, 10:149-157, 1997.
- [30] Yang J., Neiger G. and Gafni E., Structured Derivations of Consensus Algorithms for Failure Detectors. *Proc. 17th ACM Symposium on Principles of Distributed Computing (PODC'98)*, pp.297-308, 1998.

A Appendix: A simple addition $\diamond\mathcal{S}_x + \diamond\phi^y \rightsquigarrow \diamond\mathcal{S}_n$ ($x + y > t$)

This appendix presents a simple algorithm that adds the power of ϕ^y and the power of \mathcal{S}_x (resp., $\diamond\phi^y$ and $\diamond\mathcal{S}_x$) to provide a failure detector of the class \mathcal{S}_n (resp., $\diamond\mathcal{S}_n$). (Let us remind that $\mathcal{S}_n = \mathcal{S}$ and $\diamond\mathcal{S}_n = \diamond\mathcal{S}$.) The algorithm is described in Figure 10. As the failure detector classes $\Omega^1 = \Omega$ and $\diamond\mathcal{S}_n = \diamond\mathcal{S}$ are equivalent (they have the same computational power as far as failures are concerned) [3, 6, 20], it follows from Theorem 9 that the algorithm requires $x + y > t$ (which becomes a necessary and sufficient requirement for such a transformation).

To show the versatility of the approach, the algorithm is expressed in the shared memory model. It can be easily translated in the message-passing model without adding any requirement on t . Each process p_i has the following local variables:

- $suspected_i$ is a local variable that p_i can only read. It contains the set of processes provided to p_i by its underlying failure detector module of the class \mathcal{S}_x (resp., $\diamond\mathcal{S}_x$). These sets satisfy the properties defining the class \mathcal{S}_x (resp., $\diamond\mathcal{S}_x$): they eventually includes all crashed processes and x of these sets do not include the same correct process from the very beginning in the case of \mathcal{S}_x (or after some unknown but finite time in the case of $\diamond\mathcal{S}_x$).
- $SUSPECTED_i$ is the local set of processes built by the algorithm. The sets $SUSPECTED_i$ of all the processes have to satisfy the properties defining \mathcal{S} (resp., $\diamond\mathcal{S}$). Initially, $SUSPECTED_i = \emptyset$.
- new_i and $prev_i$ are two auxiliary variables. Each is an array of size n initialized to the zero vector.

The shared memory is made up of two arrays denoted $alive[1 : n]$ and $suspect[1 : n]$. Each of their entries is a single writer/multi reader atomic variable. The $alive[i]$ and $suspect[i]$ variables are repeatedly updated by p_i until it (possibly) crashes (see task $T1$ in Figure 10). Their meaning is the following:

- $alive[i]$ is only increased by p_i to indicate it has not crashed. This means that, after a process p_i crashes, the value of $alive[i]$ does not change⁶.
- $suspect[i]$ is used by p_i to inform the other processes about the value of its local $suspected_i$ set.

The task $T2$ of a process p_i repeats forever a set of statements whose aim is to compute the current value of the local set $SUSPECTED_i$ (line 07) whose value is used by the upper layer protocol. To carry out this computation, p_i first reads the shared array $alive[1 : n]$ to know which processes have progressed (the reading of the whole array is not atomic). It reads this array until it knows that all the processes that have not progressed have crashed (lines 02-05). Then, trusting the processes it considers as not crashed (the set $live$), it updates its local set $SUSPECTED_i$ according to the current suspicions made public by these processes.

⁶It is possible to have a bounded implementation for each shared variable $alive[i]$. We do not elaborate on this for two reasons: on one side it would make the protocol much more involved, on another side this is not necessary to prove our goal.

```

task T1: repeat forever
(0 1)      $alive[i] \leftarrow alive[i] + 1; suspect[i] \leftarrow suspected_i$ 
          end_repeat

task T2: repeat forever
(0 2)     repeat for_each  $j \in \{1, \dots, n\}$  do  $new_i[j] \leftarrow alive[j]$  end_for;
(0 3)     let  $live = \{j \mid new_i[j] > prev_i[j]\}$ ;
(0 4)     let  $X = \{1, \dots, n\} \setminus live$ ;
(0 5)     until  $\phi\text{-QUERY}(X)$  end_repeat;
(0 6)      $prev_i \leftarrow new_i$ ;
(0 7)      $SUSPECTED_i \leftarrow (\bigcap_{j \in live} suspect[j]) \setminus live$ 
          end_repeat

```

Figure 10: From $\phi^y + \mathcal{S}_x$ to \mathcal{S} (resp., $\diamond\phi^y + \diamond\mathcal{S}_x$ to $\diamond\mathcal{S}$), (algorithm for p_i)

Theorem 14 *Let $x + y > t$. If the underlying failure detector belongs to the class \mathcal{S}_x (resp., $\diamond\mathcal{S}_x$), the sets $SUSPECTED_i$ built by the ϕ^y -based (resp., $\diamond\phi^y$ -based) algorithm described in Figure 10 define a failure detector of the class \mathcal{S} (resp., $\diamond\mathcal{S}$).*

Proof Let us first show that the inner loop always terminates. Proving this termination is required to claim that the variable $SUSPECTED_i$ is updated at line 07. We consider three cases according to the size of the set parameter X when p_i invokes $QUERY(X)$ at line 05.

- $|X| > t$. In that case, due to the triviality property, the query returns *false*, and p_i enters again the loop. But, as there are at most t faulty processes, each correct process p_j infinitely often increases $alive[j]$ (task T1), and $prev_i[j]$ remains constant within the inner loop, there is a time after which every query issued by p_i is such that $|X| \leq t$. We are then in one of the cases that follow.
- $|X| \leq t - y$. In that case, due to the triviality property, the query returns *true* and p_i exits the inner loop.
- $t - y < |X| \leq t$. If the query returns *false*, p_i enters again the loop. We show that the query eventually returns *true*. Let us consider a process p_j that belongs to \mathcal{S} (this means that $alive[j] = prev_i[j]$ at line 03 of the task T2 executed by p_i). If p_j is correct, there is eventually an inner loop such that $alive[j] > prev_i[j]$ because p_j increases forever $alive[j]$ and $prev_i[j]$ remains constant within the inner loop. This means that eventually such a p_j will disappear from the set X defining the query parameter. It follows that, eventually, the set X used as a query parameter (1) contains only faulty processes or (2) has a size smaller than or equal to $t - y$. Due to the liveness (case 1) or triviality (case 2) property, there is then a query that eventually returns *true*.

Let us now show that, if the sets $suspected_i$ satisfy the strong completeness property, this property is also satisfied by the sets $SUSPECTED_i$. If a process p_k crashes, due to the strong completeness of the sets $suspected_i$, it eventually belongs to the set $suspected_j$ of each non-crashed process p_j . Due to line 01, after some finite time, p_k is always in $suspect_j$ (until p_j possibly crashes). Moreover, as after some time p_k no longer increases $alive[k]$, there is a finite time after which it never belongs to the *live* set computed by any process. Due to line 07, it eventually belongs to (and remains permanently in) the set $SUSPECTED_i$ of any non-crashed process p_i .

The last part of the proof concerns the weak accuracy property. We formulate the proof for going from the class \mathcal{S}_x to the classes \mathcal{S} . (The proof for going from the class $\diamond\mathcal{S}_x$ and $\diamond\phi^y$ to the class $\diamond\mathcal{S}$ is similar, and is consequently omitted.) So, we have to show that, if $x + y > t$ and the sets $suspected_i$ satisfy the limited scope perpetual weak accuracy property (namely, there is a correct process, say p_ℓ , that is not suspected by at least x -correct or faulty- processes), then the sets $SUSPECTED_i$ satisfy perpetual weak accuracy (there is a correct process -namely, p_ℓ again in our transformation- that is no suspected by any process). We consider two cases, according to the size of the set X when a process p_i exits the inner loop.

- $|X| \leq t - y$.
In that case, the exit of the inner loop was due to the triviality property. As $t - y < x$, we have $|X| < x$,

which (due the limited scope perpetual weak accuracy) means that at least one process p_k of the set *live* of p_i is such that p_ℓ never belongs to *suspected_k*, and consequently p_ℓ never belongs to *suspect[k]*. It then follows that p_ℓ can never belong to the intersection computed at line 07, which proves the case.

- $t - y < |X| \leq t$.

In that case, due to the safety property, all the processes in X have crashed. We examine two subcases.

- $t - y < |X| < x$. The proof of this case ($|X| < x$) is the same as the previous one.
- $t - y < x \leq |X|$. In that case, it is possible that all the processes that do not suspect p_ℓ have crashed, and all the remaining processes p_j do suspect p_ℓ (i.e., $p_\ell \in \textit{suspected}_j$). But in that case (noticing that X and *live* define a partition of the whole set of processes), a process that is not in the *live* set of p_i has necessarily crashed (safety and non-triviality properties). So, p_ℓ necessarily belongs to the set *live* of p_i . It follows from line 07 that p_ℓ cannot belong to *SUSPECTED_i*, which proves the case.

□*Theorem 14*