# On the Weakest Failure Detector Ever

R. Guerraoui[1,2]    M. Herlihy[3]    P. Kouznetsov[4]    N. Lynch[1]    C. Newport[1]

[1] Computer Science and Artificial Intelligence Laboratory, MIT
[2] School of Computer and Communication Sciences, EPFL
[3] Computer Science Department, Brown University
[4] Max Planck Institute for Software Systems

## Abstract

Many problems in distributed computing are impossible when no information about process failures is available. So what is the minimal yet non-trivial failure information? In other words, what is the minimal information about failures *needed* to circumvent *any* impossibility and *sufficient* to circumvent *some* impossibility.

This paper proposes a candidate abstraction, denoted $\Upsilon$, to capture this failure information. In every run of the distributed system, $\Upsilon$ eventually informs the processes that *some* set of processes in the system cannot be the set of correct processes in that run. Although seemingly weak, for it might provide random information for an arbitrarily long period of time, and it only excludes one possibility of correct set among many, $\Upsilon$ still captures non-trivial failure information. We show that $\Upsilon$ is *sufficient* to circumvent the celebrated *wait-free set-agreement* impossibility. While doing so, we (a) disprove previous conjectures about the weakest failure detector to solve set-agreement and we (b) prove that solving set-agreement with registers is strictly weaker than solving $n + 1$-process consensus using $n$-process consensus.

We prove that $\Upsilon$ is, in some sense, *necessary* to circumvent any wait-free impossibility. As a corollary, set-agreement is, from a failure information perspective, a minimal wait-free impossible problem in distributed computing.

Our results are generalized through an abstraction $\Upsilon^f$ that we introduce and prove necessary to solve any problem that cannot be solved in an $f$-resilient manner, and yet sufficient to solve $f$-resilient $f$-set-agreement.

# 1 Introduction

**Motivation.**  Fischer, Lynch, and Paterson's seminal result in the theory of distributed computing [8] says that the seemingly easy *consensus* task (a decision task where a collection of processes start with some input values and need to agree on one of the input values), cannot be deterministically solved in an asynchronous distributed system that is prone to process failures, even if processes simply fail by crashing, i.e., prematurely stop taking steps of their algorithms. Later, three independent groups of researchers [17, 12, 2] extended that result by proving the impossibility of *wait-free set-agreement* [6], a decision task where processes need to agree on up to $n$ input values, in an asynchronous shared memory model of $n+1$-processes among which $n$ can crash. This result was then extended to prove the asynchronous impossibility of *$f$-resilient $f$*-set agreement [2], i.e., $f$-set agreement when $f$ processes can crash.

*Asynchrony* refers to the absence of timing assumptions on process speeds and communication delays. Some timing assumptions can typically be made in most real distributed systems however. In the best case, if we assume precise knowledge of bounds on communication delays and process relative speeds, then it is easy to show that all asynchronous impossibilities can be circumvented. Intuitively, such timing assumptions circumvent asynchronous impossibilities by providing processes with *information about failures*, typically through *time-out* (or *heart-beat*) mechanisms.

In general, whereas certain information about failures can indeed be obtained in distributed systems, it is nevertheless reasonable to assume that this information is only partial and might sometimes be inaccurate. Typically, bounds on process speeds and message delays hold only during certain periods of time, or only in certain parts of the system. Hence, the information provided about the failure of a process might not be perfect. The motivation of this work is to identify information about failures that is *necessary* to circumvent *any* (asynchronous) impossibility and yet *sufficient* to circumvent *some* impossibility. In other words, we seek for the minimal *non-trivial* information about failures or, in the parlance of Chandra et al. [4], the weakest failure detector that cannot be implemented in an asynchronous system. By doing so, and assuming that this minimal information is sufficient to circumvent the impossibility of some problem $T$, then $T$ would, from the failure detection perspective, be the weakest impossible problem in asynchronous distributed computing.

**Contributions.**  We focus in this paper on the shared memory model. For presentation simplicity, we also consider first the wait-free case and assume a system with $n+1$ processes among which $n$ can crash. Then we move to the *$f$-resilient* case where $f \leq n$ processes can crash.

We propose to capture minimal information about failures to circumvent wait-free impossibilities in the form of a *failure detector* oracle we introduce here, denoted by $\Upsilon$. This oracle outputs, whenever queried by a process, a non-empty set of processes in the system. The output might be varying for an arbitrarily long period. Eventually, however, the output set should: (a) be the same at all correct processes and (b) not be the exact set of correct processes. Failure detector $\Upsilon$ provides very little information about failures: it basically only excludes one possibility of failures among many, and it does so only eventually. In particular, $\Upsilon$ does not say which set of processes are correct, and the set it outputs might never contain any correct (resp. faulty) process.

To illustrate $\Upsilon$, consider for instance a system of 3 processes, $p_1, p_2, p_3$, and a run where $p_1$ fails whilst $p_2$ and $p_3$ are correct. Oracle $\Upsilon$ can output any set of processes for an arbitrarily long period, it can keep randomly changing this set and can output different sets at different processes. Eventu-

ally however, $\Upsilon$ should permanently output, at $p_2$ and $p_3$, either $\{p_1\}, \{p_2\}, \{p_3\}, \{p_1, p_3\}, \{p_1, p_2\}$ or $\{p_1, p_2, p_3\}$.

We prove that, although seemingly pretty weak, $\Upsilon$ is sufficient to solve $n$-set-agreement with read/write objects (registers), in a system of $n + 1$ processes among which $n$ might crash. In other words, $\Upsilon$ is sufficient to circumvent the celebrated wait-free set-agreement impossibility. At the heart of our algorithm lies the following idea. We use the information provided by $\Upsilon$ to eventually partition the processes into two non-overlapping subsets: *gladiators* that are eventually permanently output by $\Upsilon$, and *citizens* that are not. Very roughly speaking, our algorithm consists for the *gladiators* to compete until they (a) eliminate one of their initial values, which necessarily happens if one of them fails, or (b) escape from their condition if they see the initial value of a *citizen*. Interestingly, what $\Upsilon$ ensures is that at least one of the *gladiators* fail or at least one of the *citizens* is correct.

We use the very existence of our algorithm to (a) disprove the conjecture of [16] about the weakest failure detector to wait-free implement set-agreement and (b) prove that wait-free implementing set-agreement with read/write objects is strictly weaker than solving $n+1$-process consensus using $n$-process consensus.

We then show how our algorithm can be adapted to solve $f$-set-agreement in a system of $n + 1$ processes where $f \leq n$ processes can fail, using a generalization of $\Upsilon$, which we denote $\Upsilon^f$. This oracle outputs a set of processes of size at least $n + 1 - f$ such that (as for $\Upsilon^f$) eventually: the same set is permanently output at all correct processes, and this set is not the exact set of correct processes.

We finally prove that $\Upsilon^f$ encapsulates, in a precise sense, minimal failure information to circumvent any impossibility in an asynchronous shared memory system where $f$ processes can crash. This minimality holds even if the shared memory contains any atomic object type, i.e., beyond just registers. To establish this minimality, we use here a restricted variant of the reduction notion of Chandra, Hadzilacos, and Toueg [4]. We show that any oracle that (a) might output any information (within its range) for an arbitrary period of time, and eventually outputs a permanent value that depends only on which processes failed (not on the computation), as well as (b) helps circumvent *any* impossibility in an asynchronous system with $f$ failures, can be used to compute a set of processes of size $n + 1 - f$ that is not the set of correct processes, i.e., can be used to emulate $\Upsilon^f$. Our proof is surprisingly simple: basically, it extracts the minimal information about failures to solve a decision problem, from the impossibility of the problem itself.

**Related work.** Chandra, Hadzilacos, and Toueg established in [4] the *weakest* failure detector to solve *consensus*, in the form of a failure detector oracle, denoted by $\Omega$. This oracle outputs, whenever queried by a process, a single *leader* process. Eventually, the outputs stabilize on the same correct leader at all processes. $\Omega$ is the weakest failure detector to solve consensus in the sense that (a) there is an algorithm that solves consensus using $\Omega$, and (b) for every oracle $\mathcal{D}$ that provides (only) information about failures such that some algorithm solves consensus using $\mathcal{D}$, there is an algorithm that emulates $\Omega$ using $\mathcal{D}$. In short, every such $\mathcal{D}$ encapsulates at least as much information about failures as $\Omega$. The motivation of our work is to address the general question of the necessary failure information about failures that is needed to circumvent *any* asynchronous impossibility, i.e., beyond consensus.

Not surprisingly, in the case of 2 processes (i.e., the case where wait-free set-agreement coincides with consensus) $\Omega$ and $\Upsilon$ are equivalent. Our minimality result is more restrictive but the proof is

2

significantly simpler than that of [4]: in short, our approach extracts $\Omega$ from the fact of consensus impossibility [8], without having to go through valence arguments as in [4].

In the same vein, and in the general case where $n > 2$, our approach extracts $\Upsilon$ directly from the fact of set-agreement impossibility [17, 12, 2], without having to go through Sperner arguments. In this case, we prove that $\Upsilon$ is strictly stronger than failure detector $\Omega^n$ introduced in [15]. The latter failure detector, outputs, whenever queried by a process, a subset of $n$ processes such that, eventually, it is the same subset at all correct processes and it contains at least one correct process. Failure detector $\Omega^n$ was shown to be sufficient to solve (1) wait-free set-agreement using registers [15], and (2) $n + 1$-process consensus using $n$-process consensus [18]. In fact, $\Omega^n$ was also shown to be necessary to implement $n + 1$-process consensus using $n$-process consensus [10] and conjectured to be necessary to solve set-agreement [16]. It was our long quest to prove this conjecture that led us identify $\Upsilon$ and devise our set-agreement algorithm based on this oracle. We show through our necessity result that set-agreement is, in a precise sense, minimal.

Our minimality result is clearly less general than if we considered the original reduction notion of [4]. Nevertheless, we prove $\Upsilon$ to be minimal among all failure detectors (we are aware of) that have been proposed to capture minimal information to circumvent asynchronous impossibilities in the crash-stop shared memory model [5, 4, 15, 10].

**Roadmap.** The rest of the paper is organized as follows. Section 2 gives some basic definitions needed to state and prove our results. Section 3 defines and discusses $\Upsilon$. Sections 4 describes our set-agreement algorithm using $\Upsilon$. Section 5 proves the minimality of $\Upsilon$ and $\Upsilon^f$. Section 6 concludes the paper by discussing open questions.

## 2   Model

Our model of processes communicating through shared objects and using failure detectors is based on [13, 14, 4]. We recall below the details necessary for describing our results.

**Processes and objects.** The distributed system we consider is composed of a set $\Pi$ of $n + 1$ processes $\{p_1, p_2, .., p_{n+1}\}$. Processes are subject to *crash* failures. A process that never fails is said to be correct.

Process communicate through applying atomic operations on a collection of *shared objects*. When presenting our algorithms, we assume that the shared objects are registers, i.e., they export only read-write operations. The impossibility and necessity parts of our results do not restrict the types of shared objects.

**Failure patterns and failure detectors.** We denote by $\mathcal{D}_i$ the local module for process $p_i$ of failure detector $\mathcal{D}$.

A *failure pattern* $F$ is a function from a global time range $\mathbb{T}$ to $2^\Pi$, where $F(t)$ denotes the set of processes that have crashed by time $t$. Once a process crashes, it does not recover, i.e., $\forall t : F(t) \subseteq F(t + 1)$. We define $correct(F) = \Pi - \cup_{t \in \mathbb{T}} F(t)$, the set of *correct* processes in $F$. Processes in $\Pi - correct(F)$ are called *faulty* in $F$. A process $p \in F(t)$ is said to be *crashed* at time $t$. An *environment* is a set of failure patterns. Unless stated otherwise, we assume the environment that includes all failure patterns in which at least one process is correct, i.e., we assume that $n$ or less processes can fail.

A *failure detector history $H$ with range $\mathcal{R}$* is a function from $\Pi \times \mathbb{T}$ to $\mathcal{R}$. Informally, $H(p, t)$ is the value output by the failure detector module of process $p$ at time $t$. A *failure detector $\mathcal{D}$ with range $\mathcal{R}_{\mathcal{D}}$* is a function that maps each failure pattern to a *set* of failure detector histories with range $\mathcal{R}_{\mathcal{D}}$ (usually defined by a set of requirements that these histories should satisfy). $\mathcal{D}(F)$ denotes the set of possible failure detector histories permitted by $\mathcal{D}$ for failure pattern $F$. Note that we do not restrict possible ranges of failure detectors.

**Algorithms.** We define an *algorithm $A$ using a failure detector $\mathcal{D}$* as a collection of deterministic automata, one for each process in the system. $A(p_i)$ denotes the automaton on which process $p$ runs the algorithm $A$. Computation proceeds in atomic *steps* of $A$. In each step of $A$, process $p$

(i) invokes an operation on a shared object and receives a response from the object, *or* queries its failure detector module $\mathcal{D}_i$ and receives a value from $\mathcal{D}_i$ (in the latter case, we say that the step of $p$ is a *query* step), and

(ii) applies its current state, the response received from the shared object *or* the value output by $\mathcal{D}_i$ to the automaton $A(p)$ to obtain a new state.

A step of $A$ is thus identified by a pair $(p_i, x)$, where $x$ is either the value returned by the invoked operation on a shared object and the resulting object state or, if the step is a query step, the failure detector value output at $p$ during that step.

**Configurations and runs.** A *configuration* of $A$ defines the state of each process and each object in the system. In an *initial configuration $I$* of $A$, every process $p_i$ is in an initial state of $A(p_i)$ and every object is in an initial state determined by its object type.

A *run of algorithm $A$ using a failure detector $\mathcal{D}$* is a tuple $R = \langle F, H, I, S, T \rangle$ where $F$ is a failure pattern, $H \in \mathcal{D}(F)$ is a failure detector history, $I$ is an initial configuration of $A$, $S$ is an *infinite* sequence of steps of $A$, and $T \subseteq \mathbb{T}$ is an *infinite* list of non-decreasing time values indicating when each step of $S$ has occurred such that:

(1) $S$ respects the specifications of all shared objects and all process automata, given their initial states in $I$;

(2) For all $1 \leq k \leq |S|$, if $S[k] = (p_i, x)$, then $p_i$ has not crashed by time $T[k]$, i.e., $p_i \notin F(T[k])$;

(3) For all $1 \leq k \leq |S|$, if $x \in \mathcal{R}_{\mathcal{D}}$, then $x$ is the value of the failure detector module of $p_i$ at time $T[k]$, i.e., $d = H(p_i, T[k])$;

(4) For all $1 \leq k, l \leq |S|$, if $T[k] = T[l]$, then $S[k]$ and $S[l]$ are steps of different processes.

(5) Every correct (in $F$) process takes infinitely many steps in $S$.

A *partial run of $A$* is a finite prefix of a run of $A$. [1]

A problem is a set of runs. An algorithm $A$ solves a problem $M$ *using a failure detector $\mathcal{D}$*, if every run of $A$ using $\mathcal{D}$ is in $M$.

---

[1] A more formal definition of a run of an algorithm using a failure detector can be found in [4, 10].

4

**Comparing failure detectors.** If, for failure detectors $\mathcal{D}$ and $\mathcal{D}'$, there is an algorithm $T_{\mathcal{D}' \to \mathcal{D}}$ using $\mathcal{D}'$ that *extracts* the output of $\mathcal{D}$, i.e. implements a distributed variable $\mathcal{D}$-*output* such that in every run $R = \langle F, H', I, S, T \rangle$ of $T_{\mathcal{D}' \to \mathcal{D}}$, there exists $H \in \mathcal{D}(F)$ such that for all $p_i \in \Pi$ and $t \in \mathbb{T}$, $H(p_i, t) = \mathcal{D}$-*output*$_i(t)$ (i.e., the value of $\mathcal{D}$-*output* output at $p_i$ at time $t$), then we say that $\mathcal{D}$ *is weaker than* $\mathcal{D}'$. If $\mathcal{D}$ is weaker than $\mathcal{D}'$ but $\mathcal{D}'$ is *not* weaker than $\mathcal{D}$, then we say that $\mathcal{D}$ *strictly weaker than* $\mathcal{D}'$. If $\mathcal{D}$ and $\mathcal{D}'$ are weaker than each other, we say they are equivalent.

If $\mathcal{D}$ is weaker than $\mathcal{D}'$, then $\mathcal{D}'$ provides at least as much information about failures as $\mathcal{D}$: every problem that can be solved using $\mathcal{D}$ can also be solved using $\mathcal{D}'$. $\mathcal{D}$ is the weakest failure detector to solve a problem $M$ if there is an algorithm that solves $M$ using $\mathcal{D}$ and $\mathcal{D}$ is weaker than any failure detector that can be used to solve $M$. If the weakest failure detector to solve a problem $A$ is strictly weaker than the weakest failure detector to solve a problem $B$, then we say that $A$ is strictly weaker than $B$, i.e., $A$ requires strictly less failure information than $B$.

# 3 The Candidate Failure Detector

We introduce failure detector $\Upsilon$, which outputs a non-empty set of processes ($\mathcal{R}_\Upsilon = 2^\Pi - \{\emptyset\}$), such that for every failure pattern $F$ and every failure detector history $H \in \Upsilon(F)$, eventually:

(1) the same set $S \in 2^\Pi - \{\emptyset\})$ is permanently output at all correct processes.

(2) this set $S$ is not the set of correct processes in $F$, i.e., $S \notin \text{correct}(F)$.

In a system of 2 processes, $\Upsilon$ and $\Omega$ [4] are equivalent. (Recall that $\Omega$ outputs a *leader* process so that eventually the same correct leader is output at all correct processes). Basically, to get $\Upsilon$ from $\Omega$, output the complement of $\Omega$ in $\Pi$. On the other hand, to get $\Omega$ from $\Upsilon$, output the complement of $\Upsilon$ if this is a singleton, and output the process itself otherwise.

$\Omega$ was generalized to a failure detector $\Omega^n$ [15], which outputs a subset of processes of size $n$ so that, eventually, the same subset containing at least one correct process is permanently output at all correct processes. (Clearly, $\Omega^1$ is $\Omega$.) The complement of $\Omega^n$ in $\Pi$ is a legal output for $\Upsilon$. Hence, $\Upsilon$ is weaker than $\Omega^n$. The converse is however not true in our default environment where $n$ processes can fail, as we show below.

**Theorem 1** $\Upsilon$ *is strictly weaker than* $\Omega^n$ *if* $n \geq 2$.

**Proof.** Emulating $\Omega^n$ boils down to eventually identifying, in every run, a process $p_c$ that is *not the only* correct process in that run. In the following, we show that this identification is impossible with $\Upsilon$.

We proceed by contradiction and we assume that we can extract the output of $\Omega^n$ from $\Upsilon$, i.e., there is an algorithm $T$ that, using $\Upsilon$, eventually outputs the same $p_c$ at every correct process such that $\Pi - \{p_c\}$ contains at least one correct process. To establish a contradiction, we construct a run of $T$ in which the extracted output never stabilizes.

We consider the set of runs of $T$ in which $\Upsilon$ permanently outputs $\{p_1, p_2, .., p_n\}$ at all processes. This is a legitimate output of $\Upsilon$ if either $p_{n+1}$ is correct or there is at least one faulty process in $\{p_1, p_2, .., p_n\}$.

Consider a run $R_1$ in which $p_{n+1}$ is the only correct process. Clearly, $T$ eventually permanently outputs a process $p_{i_1} \in \{p_1, p_2, .., p_n\}$ in run $R_1$. Let $R'_1$ be a prefix of $R_1$ in which $T$ outputs $p_{i_1}$ at some process.

Now let $R_2$ be an extension of $R_1'$ in which (1) $p_{n+1}$ takes at least one step after the last step of $R_1'$, and (2) $p_{i_1}$ is the only correct process. Note that, since $n \geq 2$, at least one process in $\{p_1, p_2, .., p_n\}$ is faulty in $R_2$, and it is still legitimate for $\Upsilon$ to output $\{p_1, p_2, .., p_n\}$ in $R_2$. Thus, in $R_2$, $T$ should eventually permanently output a process $p_{i_2} \in \{p_1, \ldots, p_{i-1}, p_{i+1}, \ldots, p_{n+1}\}$.

Let $R_2'$ be a prefix of $R_2$ and an extension of $R_1'$ in which some process outputs $p_{i_2}$ after $R_1'$. Now let $R_3$ be an extension of $R_2'$ in which (1) $p_{n+1}$ takes at least one step after $R_2'$, and (2) $p_{i_2}$ is the only correct process. Again, in $R_3$, $T$ should eventually permanently output a process $p_{i_3} \in \{p_1, \ldots, p_{i-2}, p_{i+2}, \ldots, p_{n+1}\}$.

Following this procedure, we obtain an infinite run $R$ in which $p_{n+1}$ takes infinitely many steps. Thus, $\{p_1, p_2, \ldots, p_n\}$ cannot be the set of correct processes in $R$, i.e., it is legitimate for $\Upsilon$ to output $\{p_1, p_2, \ldots, p_n\}$. Thus, $R$ is a run of $T$ using $\Upsilon$. But the extracted output never stabilizes in $R$ — a contradiction. $\qquad\square$

# 4  Set-Agreement

## 4.1  The problem

In the $k$-set-agreement problem, processes need to agree on at most $k$ values out of a possibly larger set of values. Let $V$ be the value domain such that $\bot \notin V$. Every process $p_i$ starts with an initial value $v$ in $V$ (we say $p_i$ *proposes* $v$), and aims at reaching a state in which $p_i$ irrevocably commits on a decision value $v'$ in $V$ (we say $p_i$ *decides* $v'$). Every run of a $k$-set-agreement algorithm satisfies the following properties: (1) *Termination:* Every correct process eventually decides a value; (2) *Agreement:* At most $k$ values are decided; (3) *Validity:* Any value decided is a value proposed.

In the following, we first focus on solving $n$-set-agreement in a system of $n + 1$ processes. We sometimes also talk about wait-free implementing set-agreement. This problem is impossible if processes can only communicate using registers, $n$ processes can crash, and no information about failures is available [17, 12, 2].

We show how to circumvent this impossibility using $\Upsilon$: we describe a protocol that solves $n$-set-agreement using registers and $\Upsilon$, while tolerating the failure of $n$ processes. Basically, wait-free implementing set-agreement consists for the processes to *exclude* at least one proposed value among the $n + 1$ possible ones. Our protocol achieves this by using the output of $\Upsilon$ to eventually split the processes into two non-overlapping subsets: those in the subset output by $\Upsilon$, and which we call *gladiators*, and those outside that subset, and which we call *citizens*. Intuitively, *gladiators* do not decide any value until either they make sure one of them gives up its value, which is guaranteed to happen if one of them crashes, or they see a value of a *citizen*, in which case they simply decide that value. The property eventually ensured by $\Upsilon$ is that either at least one of the *gladiators* crash or at least one of the *citizens* is correct.

Besides putting this intuition to work, technical difficulties handled by our protocol include coping with the facts that (1) $\Upsilon$ might output random sets for an arbitrarily long periods of time, providing divergent and temporary information about who is *gladiator* and who is *citizen*, and (2) *citizens* might be faulty. A key procedure we use to handle these difficulties is the *k-converge* routine, introduced in [18]. A process calls *k-converge* with an input value in $V$ and gets back an output value in $V$ and a boolean $c$. We say that the process *picks* $v$ and, if $c = true$, we say that the process *commits* $v$. The *k-converge* routine ensures the following properties: (1) *C-Termination:*

every correct process picks some value; (2) *C-Validity*: if a process picks $v$ then some process invoked *k-converge* with $v$; (3) *C-Agreement*: If some process commis to a value, then then at most $k$ values are picked; (4) *Convergence*: If there are at most $k$ different input values, then every process that picks a value commits. The *k-converge* routine can be implemented, for any $k$, using registers [18]. By definition, 0-*converge*$(v)$ always returns $(v, false)$.

## 4.2 The protocol

The abstract pseudo-code of the protocol that solves $n$-set agreement using $\Upsilon$ and registers is described in Figure 1.

---

Shared abstractions:
  Registers $D$, $D[\,]$, initially $\bot$; *Binary registers Stable*$[\,]$, initially *true*
  *Convergence instances: n-converge*$[\,]$, *j-converge*$[\,][\,]$, for all $j = 0, \ldots, n$

Code for every process $p_i$:

```
1      v_i := the input value of p_i; r := 0
2      repeat
3        r := r + 1
4        (v_i, c) := n-converge[r](v_i)
5        if c = true then
6            D := v_i; return (v_i)
7        S := query(Υ_i)
8        if p_i ∉ S then
9            D[r] := v_i
10       else
11           k := 0
12           repeat
13              k := k + 1
14              (v_i, c) := (|S| − 1)-converge[r][k](v_i)
15              if c = true then D[r] := v_i
16              if S ≠ query(Υ_i) then Stable[r] := false
17           until D ≠ ⊥ or D[r] ≠ ⊥ or ¬Stable[r]
18           if D[r] ≠ ⊥ then
19              v_i := D[r]
20    until D ≠ ⊥
21    return (D)
```

---

Figure 1: $\Upsilon$-based set agreement protocol.

The protocol proceeds in rounds. In every round $r$, the processes first try to reach agreement using $n$-convergence (line 4). If a process $p_i$ commits to a value $v$, then $p_i$ writes $v$ in register $D$ and returns $v$. If $p_i$ fails to commit (which can only happen if all $n + 1$ processes take part in the $n$-convergence instance), then $p_i$ queries its module of $\Upsilon$. Let $S$ be the returned value.

Now $p_i$ cyclically executes the following procedure (lines 12–17). If $p_i$ does not belong to $S$ ($p_i$ believes it is a *citizen*), then $p_i$ writes its value in a shared register $D[r]$ and proceeds to the next round. Otherwise ($p_i$ believes it is a *gladiator*), $p_i$ takes part in the $(|S| - 1)$-convergence protocol trying to eliminate one of the values concurrently proposed by processes in $S$. (Recall that, by definition, 0-*converge*$(v)$ always returns $(v, false)$.) The procedure is repeated as long as none of the conditions in line 17 is satisfied, i.e., (a) no process participating in the current round $r$ reports that its module of $\Upsilon$ has not yet stabilized, (b) $(|S| - 1)$-convergence does not commit to a value,

7

and (c) no non-$\perp$ value is found in $D[r]$ or $D$ (line 17). If $p_i$ finds $D[r] \neq \perp$, then $p_i$ adopts the value in $D[r]$ and proceeds to round $r + 1$. If $p_i$ finds $D \neq \perp$ then $p_i$ returns $D$.

Remember that there is a time after which $\Upsilon$ permanently outputs, at all correct processes, the same set $S$ that is not the set of correct processes: $S$ either contains a faulty process or there is a correct process outside $S$. Thus, no process can be blocked in round $r$ by repeating forever the procedure described above: eventually, either some process outside $S$ writes its value in $D[r]$, or some process is faulty in $S$ and $(|S| - 1)$-convergence returns a committed value.

As a result, eventually, there is a round in which at least one input value is eliminated: either some process in $S$ adopts a value from outside $S$, or processes in $S$ commit to at most $|S| - 1$ input values. In both cases, every process that participates in $n$-convergence in round $r + 1$ (line 4) commits one of at most $n$ "survived" values.

**Theorem 2** *The algorithm in Figure 1 solves $n$-set agreement using $\Upsilon$ and registers.*

**Proof.** Consider an arbitrary run $R$ of the algorithm in Figure 1.

Validity immediately follows from the protocol and the C-Validity property of $k$-converge.

Agreement is implied by the fact that every decided value is first committed by $n$-convergence (line 4). Indeed, let $r$ be the first round in which some process $p_i$ commits to a value after invoking $n$-*converge*$[r]$. By the C-Agreement property of $n$-convergence, every process that invoked $n$-*converge*$[r]$ adopted at most $n$ different values. Thus, no more than $n$ different values can ever be written in register $D$. Since a process is allowed to decide a value only if the value was previously written in $D$ (lines 6 and 21), at most $n$ different values can be decided.

Now consider Termination. We observe first that no process can decide unless $D$ contains a non-$\perp$ value, and if $D \neq \perp$, then every correct process eventually decides. This is because the converge instances are non-blocking and every correct process periodically checks whether $D$ contains a non-$\perp$ value and, if there is one, returns the value (lines 20 and 17). Assume now, by contradiction, that $D = \perp$ forever and, thus, no process ever decides in $R$.

Let $S$ be the stable output of $\Upsilon$ in $R$, i.e., at every correct process, $\Upsilon$ eventually permanently outputs $S$. Whenever a process observes that the output of $\Upsilon$ is not stable in round $r$, it sets register $Stable[r]$ to *true* (line 16) and proceeds to the next round. Further, if a process finds $D[r] \neq \perp$, then eventually every correct process finds $D[r] \neq \perp$ and proceeds to the next round. Moreover, by our assumption, no process ever writes in $D$ and returns in line 6. Thus, there exists a round $r$ such that every correct process reaches $r$, and the observed output of $\Upsilon$ at every process that reached round $r$ has stabilized on $S$.

Recall that $S$ is a non-empty set of processes that is *not* the set of correct processes in $R$, i.e., $S \neq \emptyset$ and $S \neq C$, where $C$ is the set of correct processes in $R$. Thus, two cases are possible: (1) $C \subsetneq S$, and (2) $C - S \neq \emptyset$.

In case (1), there is at least one faulty process in $S$. Since every faulty process eventually crashes, there exists $k \in \mathbb{N}$, such that at most $|S| - 1$ values are proposed to $(|S| - 1)$-*converge*$[r][k]$. By the Convergence property of the $(|S| - 1)$-*converge* procedure, every correct process eventually commits to a value, writes it in $D[r]$ and proceeds to round $r + 1$.

In case (2), there is at least one correct process $p_j$ outside $S$. Thus, $p_j$ eventually reaches round $r$ and writes its current value in $D[r]$. Thus, every correct process eventually reads the value, adopts it and proceeds to round $r + 1$.

In both cases, every correct process reaches round $r + 1$. By the algorithm, every process that reaches round $r + 1$ adopted a value previously written in $D[r]$.

8

A process is allowed to write a value in $D[r]$ only if (a) the process is in $\Pi - S$, or (b) a process is in $S$ and the value is committed in $(|S| - 1)$-*converge*$[r][k]$ for some $k$. By the C-Agreement and C-Validity properties of $(|S| - 1)$-convergence and because every value returned by an instance of $(|S| - 1)$-*converge*$[r][k]$ is adopted (line 14), at most $|S| - 1$ distinct values can be written in $D[r]$ by processes in $S$. Thus, at most $n + 1 - |S| + |S| - 1 = n$ distinct values can ever be found in $D[r]$. Hence, at most $n$ distinct values can be proposed to $n$-convergence (line 4) in round $r + 1$. By the Convergence property of $n$-convergence, every correct process commits and decides — a contradiction.

Thus, eventually, every correct process decides. $\qquad\square$

**Remark.** Our algorithm actually solves a stronger version of set-agreement that terminates even if not every correct process *participates*, i.e., proposes a value and executes the protocol. Indeed, assume (by slightly changing the model) that some (possibly correct) process does not participate in a given run of the algorithm in Figure 1. Thus, in round 1, at most $n$ different values are proposed to $n$-*converge* (line 4) and, by the Convergence property of $n$-*converge*, every correct participant commits to a value. Thus, every correct *participant* returns in line 6 of round 1.

As a corollary to Theorems 1 and 2, we disprove the conjecture of [16] by showing that:

**Corollary 3** $\Omega^n$ *is not the weakest failure detector to wait-free implement set-agreement using registers.*

As a corollary to Theorems 1 and 2, and the fact that $\Omega^n$ is the weakest failure detector to implement $n + 1$-process consensus using $n$-process consensus [10], we get the following:

**Corollary 4** *Wait-free implementing set-agreement using registers is strictly weaker than implementing $n + 1$-process consensus using $n$-process consensus.*

## 4.3  $f$-Resilient Set-Agreement

For pedagogical purposes, we focused so far on the environment where $n$ out of $n + 1$ processes can crash. In this section, we consider the more general environment where $f$ processes can crash, and $0 < f < n + 1$. More specifically, we consider the environment $\mathcal{E}^f$ that consists of all failure patterns $F$ such that $faulty(F) \leq f$.

By reduction to the impossibility of wait-free set agreement, Borowsky and Gafni showed that $f$-set agreement is impossible in $\mathcal{E}^f$ [2]. We present a failure detector, which generalizes $\Upsilon$, and which circumvents this impossibility. This failure detector, which we denote by $\Upsilon^f$, outputs a set of processes of size at least $n + 1 - f$ ($\mathcal{R}_{\Upsilon^f} = \{S \subseteq \Pi : |S| \geq n + 1 - f\}$), such that, for every failure pattern $F \in \mathcal{E}^f$ and every failure detector history $H \in \Upsilon^f(F)$, eventually (as for $\Upsilon$): (1) the same set is permanently output at all correct processes, and (2) this set is not the set of correct processes in $F$. Clearly, $\Upsilon^n$ is $\Upsilon$.

Failure detector $\Omega^f$ can also be used to solve $f$-resilient $f$-set agreement (a simple variation of the consensus algorithm in [15] will work). It is easy to see that $\Upsilon^f$ is weaker than $\Omega^f$ in $\mathcal{E}^f$: to emulate $\Upsilon^f$, every process simply outputs the complement of $\Omega^f$ in $\Pi$. Eventually the correct processes obtain the same set of $n + 1 - f$ processes that is not the set of correct processes (the output of $\Omega^f$ eventually includes at least one correct process).

It is also straightforward to extract $\Omega^1 = \Omega$ from $\Upsilon^1$ in $\mathcal{E}^1$. In the reduction algorithm, every process $p_i$ periodically writes ever-growing timestamps in the shared memory. If $\Upsilon_i^1$ outputs a

proper subset of $\Pi$ (of size $n$), then $p_i$ elects the process $p_\ell = \Pi - \Upsilon_i$, otherwise, if $\Upsilon^1$ outputs $\Pi$ (i.e., exactly one process is faulty), then $p_i$ elects the process with the smallest id among $n$ processes with the highest timestamps. Eventually, the same correct process is elected by the correct processes — the output of $\Omega$ is extracted. However, in general, $\Upsilon^f$ is strictly weaker than $\Omega^f$:

**Theorem 5** $\Upsilon^f$ *is strictly weaker than* $\Omega^f$ *in* $\mathcal{E}^f$ *if* $2 \leq f \leq n$. *(Proof in the appendix)*

We present in the appendix a generalized $f$-resilient $f$-set-agreement algorithm using $\Upsilon^f$. The algorithm essentially follows the lines of our "wait-free" algorithm described in Figure 1, except that now the set $S$ of $n + 1 - f$ or more *gladiators* (processes that are eventually permanently output by $\Upsilon^f$) have to be able to eventually commit on at most $|S| + f - n - 1$ distinct values, so that, together with at most $n + 1 - |S|$ values chosen by the *citizens*, there would eventually be at most $f$ distinct values in the system. To achieve this, we add a simple mechanism based on the use of *atomic-snapshots* [1] which ensures that whenever $S$ contains at least one faulty processes, then at most $|S| + f - n - 1$ values are eventually chosen by the members of $S$.

**Theorem 6** *There is an algorithm that implements $f$-set agreement using $\Upsilon^f$ and registers in $\mathcal{E}^f$.*

# 5 The Necessity of $\Upsilon^f$

We establish here that $\Upsilon^f$ is, in a certain sense, minimal in systems where up to $f$ processes can crash, implying that $\Upsilon$ is also minimal when up to $n$ processes can crash.

We introduce the notion of a *dummy* failure detector, which always outputs the same value (i.e., its range is a singleton $\{d\}$). Clearly, a dummy failure detector $\mathcal{D}$ can be *emulated* in an asynchronous system. If a problem can be solved in $\mathcal{E}^f$ using a *dummy* failure detector, then we say that the problem is *$f$-resilient*. Otherwise, we say that the problem is *$f$-resilient impossible*. We say that a failure detector is *$f$-non-trivial* if it can be used to solve an $f$-resilient impossible problem in $\mathcal{E}^f$.

Establishing our minimality result goes through delimiting the scope of failure detectors within which $\Upsilon^f$ is minimal. We say that a failure detector, $\mathcal{D}$ with range $\mathcal{R}_\mathcal{D}$, is *eventually stable* if it satisfies the following properties:

- The same value is eventually permanently output by $\mathcal{D}$ at all correct processes. Formally, for every failure pattern $F$ and every $H \in \mathcal{D}(F)$, there exists a value $d \in \mathcal{R}_\mathcal{D}$ and $t \in \mathbb{N}$ such that for all $t' \geq t$ and $p_i \in correct(H)$, $H(p_i, t') = d$ (we say that $d$ is stable in $H$).[2]

- The stable output of $\mathcal{D}$ depends only on the set of correct processes. Formally, for all failure patterns $F$ and $H \in \mathcal{D}(F)$, if $d$ is stable in $H$, then for all $F'$ such that $correct(F') = correct(F)$, there exists $H' \in \mathcal{D}(F')$ such that $d$ is also stable in $H'$.

- $\mathcal{D}$ is allowed to output any value in its range in every finite prefix of every history of $\mathcal{D}$. Formally, for all failure patterns $F$, histories $H \in \mathcal{D}(F)$, values $d \in \mathcal{R}_\mathcal{D}$ and times $t \in \mathbb{T}$, there exists $H' \in \mathcal{D}(F)$, such that, for all $p_i \in \Pi$ and $t' \in \mathbb{T}$, $H'(p_i, t') = d$ if $t' \leq t$ and $H'(p_i, t') = H(p_i, t)$.

---

[2] Our lower bound proofs actually work also for "locally stable" failure detectors that eventually permanently output a "stable" value at every correct process (the stable values output at different correct processes can be different though).

**Theorem 7** $\Upsilon^f$ *is weaker than any $f$-non-trivial eventually stable failure detector.*

**Proof.** Let $\mathcal{D}$ be any eventually stable failure detector that can be used to solve an $f$-resilient impossible problem. Let $\mathcal{R}_\mathcal{D}$ be the range of $\mathcal{D}$, and let $d$ be any value in $\mathcal{R}_\mathcal{D}$.

First we show that there exists a set of processes $S \in \mathcal{R}_{\Upsilon^f}$ (i.e., $|S| \geq n + 1 - f$) such that for all failure patterns $F$ where $correct(F) = S$ and all histories $H \in \mathcal{D}(F)$, $d$ is not the stable value in $H$.

Suppose not, i.e., there exists a value $d \in \mathcal{R}_\mathcal{D}$ such that, for all $S \in \mathcal{R}_{\Upsilon^f}$, there exists a failure pattern $F$ and a history $H \in \mathcal{D}(F)$ such that $correct(F) = S$ and $d$ is the stable value of $H$. Since $\mathcal{D}$ is stable, for every $F$, $\mathcal{D}(F)$ contains a history in which $d$ is always output at every process.

But then every history of a dummy failure detector that always outputs $d$ is a history of $\mathcal{D}$, i.e., the dummy failure detector can implement $\mathcal{D}$. This contradicts the assumption that $\mathcal{D}$ is $f$-non-trivial.

Thus, for any $d \in \mathcal{R}_\mathcal{D}$, there exists $S \in \mathcal{R}_{\Upsilon^f}$ such that $d$ cannot be the stable value whenever $S$ is the set of correct processes. Since the elements of $\mathcal{R}_{\Upsilon^f}$ can be totally ordered, we can define a function $\sigma$ that maps every $d \in \mathcal{R}_\mathcal{D}$ to the smallest $S \in \mathcal{R}_{\Upsilon^f}$ such that whenever $d$ is stable, $S$ that cannot be the current set of correct processes.

The reduction algorithm $T_{\mathcal{D} \to \Upsilon^f}$ works as follows. Every process periodically queries its module of $\mathcal{D}$ and for every returned value $d$ outputs $S = \sigma(d)$. In every run of $T_{\mathcal{D} \to \Upsilon^f}$, the produced output eventually stabilizes at a set $S \in \mathcal{R}_{\Upsilon^f}$, that is not the set of correct processes. That is, the output of $\Upsilon^f$ is extracted. $\qquad\square$

# 6   Concluding Remarks

We stated in this paper that $\Upsilon$ (resp. $\Upsilon^f$) is weaker than any eventually stable failure detector that circumvents a wait-free (resp. $f$-resilient) impossibility.

Generalizing this result by establishing the minimality of $\Upsilon$ within a wider scope of failure detectors is left for future research. Nevertheless, and interestingly, most failure detectors (we are aware of) that have been proposed to capture minimal information to circumvent asynchronous impossibilities in the shared memory model are eventually stable or have eventually stable equivalents [5, 4, 15, 10]. Thus, our minimality result suggests that set agreement is, in a strict sense, a minimal impossible decision task. This conjecture contrasts with the recent observation of Gafni et al. [9] that the renaming task is strictly weaker than set agreement. We believe this discrepancy between our results is due to the very nature of failure detector-based protocols. For example, Delporte et al. showed in [7] that it is sufficient to provide enough information to reach agreement between any two processes to be able to solve consensus among *all* processes. Recall that this statement is not true if we consider an asynchronous system equipped with powerful shared objects: 2-consensus objects are not able to solve $n$-consensus where $n > 2$ [11, 3].

Another interesting aspect of our minimality result is that it holds regardless of which shared objects are used to circumvent an impossibility. Indeed, the only fact we use to extract the output of $\Upsilon^f$ is the very impossibility to solve a given problem in a given model. On the other hand, our $\Upsilon$-($\Upsilon^f$-)based algorithms work in the "weakest" shared memory model where processes communicate through registers.

# References

[1] Yehuda Afek, Hagit Attiya, Danny Dolev, Eli Gafni, Michael Merritt, and Nir Shavit. Atomic snapshots of shared memory. *Journal of the ACM*, 40(4):873–890, 1993.

[2] Elizabeth Borowsky and Eli Gafni. Generalized FLP impossibility result for $t$-resilient asynchronous computations. In *Proceedings of the 25th ACM Symposium on Theory of Computing*, pages 91–100, May 1993.

[3] Tushar D. Chandra, Vassos Hadzilacos, Prasad Jayanti, and Sam Toueg. Wait-freedom vs. $t$-resiliency and the robustness of wait-free hierarchies. In *ACM Symposium on Principles of Distributed Computing*, pages 334–343, August 1994.

[4] Tushar Deepak Chandra, Vassos Hadzilacos, and Sam Toueg. The weakest failure detector for solving consensus. *Journal of the ACM*, 43(4):685–722, July 1996.

[5] Tushar Deepak Chandra and Sam Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225–267, March 1996.

[6] Soma Chaudhuri. More choices allow more faults: Set consensus problems in totally asynchronous systems. *Information and Computation*, 105(1):132–158, 1993.

[7] Carole Delporte-Gallet, Hugues Fauconnier, and Rachid Guerraoui. (almost) all objects are universal in message passing systems. In *International Symposium on Distributed Computing*, pages 184–198, 2005.

[8] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, April 1985.

[9] Eli Gafni, Sergio Rajsbaum, and Maurice Herlihy. Subconsensus tasks: Renaming is weaker than set agreement. In *International Symposium on Distributed Computing*, pages 329–338, 2006.

[10] Rachid Guerraoui and Petr Kouznetsov. On failure detectors and type boosters. In *Proceedings of the 17th International Symposium on Distributed Computing*, pages 292–305. Springer-Verlag, 2003.

[11] Maurice Herlihy. Wait-free synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):123–149, January 1991.

[12] Maurice Herlihy and Nir Shavit. The asynchronous computability theorem for $t$-resilient tasks. In *Proceedings of the 25th ACM Symposium on Theory of Computing*, pages 111–120, May 1993.

[13] Maurice Herlihy and Jeannette M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.*, 12(3):463–492, 1990.

[14] Prasad Jayanti. Robust wait-free hierarchies. *Journal of the ACM*, 44(4):592–614, 1997.

[15] Gil Neiger. Failure detectors and the wait-free hierarchy. In *14th ACM Symposium on Principles of Distributed Computing*, 1995.

[16] Michel Raynal and Corentin Travers. In search of the holy grail: Looking for the weakest failure detector for wait-free set agreement. Technical Report TR 06-1811, INRIA, August 2006.

[17] Michael Saks and Fotios Zaharoglou. Wait-free $k$-set agreement is impossible: The topology of public knowledge. In *Proceedings of the Twenty fifth ACM Symposium on Theory of Computing*, pages 101–110, May 1993.

[18] Jiong Yang, Gil Neiger, and Eli Gafni. Structured derivations of consensus algorithms for failure detectors. In *Proceedings of the 17th ACM Symposium on Principles of Distributed Computing*, pages 297–306, 1998.

# Appendix: $f$-Resilient Set-Agreement

## $\Upsilon^f$ is strictly weaker than $\Omega^f$

**Theorem 5** $\Upsilon^f$ *is strictly weaker than* $\Upsilon^f$ *in* $\mathcal{E}^f$ *if* $2 \leq f \leq n$.

**Proof.** By contradiction, assume that there exists an algorithm $T$ using $\Upsilon^f$ that in every run with at least $n + 1 - f$ correct processes, eventually outputs at every correct process the same set of processes $L$ such that $|L| = f$ and $L$ contains at least one correct process. To establish a contradiction, we construct a run of $T$ in which the extracted output of never stabilizes.

We consider the set of runs of $T$ in which $\Upsilon^f$ permanently outputs $S = \{p_1, p_2, .., p_n\}$ at all processes. Recall that this is a legitimate output if either $p_{n+1}$ is correct or there is at least one faulty process in $\{p_1, p_2, .., p_n\}$.

Consider a run $R_1$ in which $\{p_1, p_2, .., p_{n+1-f}\}$ is the set of correct processes. Clearly, $T$ eventually permanently outputs a set $L_1$ in run $R_1$. Let $R_1'$ be a prefix of $R_1$ in which $T$ outputs $L_1$ at some process.

Now let $R_2$ be an extension of $R_1'$ in which (1) every process takes at least one step after the last step of $R_1'$, and (2) $\Pi - L_1$ is the set of correct processes. Note that, since $2 \leq f \leq n$, $S \neq \Pi - L_1$, and it is still legitimate for $\Upsilon^f$ to output $S$ in $R_2$. Thus, in $R_2$, $T$ should eventually permanently output a set $L_2 \neq L_1$.

Let $R_2'$ be a prefix of $R_2$ and an extension of $R_1'$ in which some process outputs $L_2$ after $R_1'$. Now let $R_3$ be an extension of $R_2'$ in which (1) every takes at least one step after the last step of $R_2'$, and (2) $\Pi - L_2$ is the set of correct processes. Again, $S \neq \Pi - L_2$ and, in $R_3$, $T$ should eventually permanently output a set $L_3 \neq L_2$.

Following this procedure, we obtain an infinite run $R$ in which every process takes infinitely many steps. Thus, $S = \{p_1, p_2, \ldots, p_n\}$ cannot be the set of correct processes in $R$, i.e., it is legitimate for $\Upsilon^f$ to always output $S$. Thus, $R$ is a run of $T$ using $\Upsilon^f$. But the extracted output of $\Omega^f$ never stabilizes in $R$ — a contradiction. $\qquad\square$

## The protocol

**Theorem 6** *The algorithm in Figure 2 implements $f$-resilient $f$-set agreement using $\Upsilon^f$ and registers.*

**Proof.** The Agreement and Validity properties are immediate from the algorithm.

Termination is proved along the lines of the proof of our "wait-free" algorithm described in Figure 1, except that now we have a new potentially blocking loop (in lines 17–17). Suppose, by contradiction, that there is a run $R$ of our algorithm in which some correct process never decides. Using the arguments presented in the proof of Theorem 2, we can show that $D$ always contains $\perp$. Further, there exists a round $r$ such that every correct process reached round $r$, and the observed output of $\Upsilon^f$ at every process that reached round $r$ has stabilized on some set $S$ in round $r$. By the properties of $\Upsilon^f$, $S$ is of size at least $n + 1 - f$ and $S$ is not the set of correct processes in $R$.

Suppose, by contradiction that some correct process is blocked in the loop of lines 17–19, while executing sub-round $k$ of round $r$. It is easy to see that, if a correct process exits the loop, then eventually every correct process is freed too. Thus, our assumption implies that every correct process is blocked in the loop of lines 17–19, while executing sub-round $k$ of round $r$.

Shared abstractions:
    Registers $D$, $D[\,]$, initially $\perp$
    Vectors of registers $A[\,][\,]$, initially $\perp$
    Binary registers $Stable[\,]$, initially $true$
    Convergence instances: $n$-$converge[\,]$, $j$-$converge[\,][\,]$, for all $j = 0, \ldots, n$

Code for every process $p_i$:

```
1      v_i := the input value of p_i
2      r := 0
3      repeat
4         r := r + 1
5         (v_i, c) := f-converge[r](v_i)
6         if c = true then
7             D := v_i
8             return (v_i)
9         S := query(Υ_i)
10        if p_i ∉ S then
11            D[r] := v_i
12        else
13            k := 0
14            repeat
15               k := k + 1
16               A[r][k][i] := v_i
17               repeat
18                  V := atomic-snapshot(A[r][k])
19               until D ≠ ⊥ or D[r] ≠ ⊥ or ¬Stable[r] or V contains ≥ n + 1 − f non-⊥ entries
20               if D ≠ ⊥ then
21                  return(D)
22               else if D[r] ≠ ⊥ then
23                  v_i := D[r]
24               else if Stable[r] then
25                  v_i := min non-⊥ value in V
26                  (v_i, c) := (|S| + f − n − 1)-converge[r][k](v_i)
27                  if c = true then
28                     D[r] := v_i
29                  if S ≠ query(Υ_i) then
30                     Stable[r] := false
31            until D ≠ ⊥ or D[r] ≠ ⊥ or ¬Stable[r]
32            if D[r] ≠ ⊥ then
33               v_i := D[r]
34      until D ≠ ⊥
35      return (D)
```

Figure 2: $\Upsilon^f$-based $f$-resilient $f$-set agreement protocol.

Hence, $\Pi - S$ contains no correct process: otherwise, some correct process in $\Pi - S$ would eventually write a non-$\perp$ value in $D[r]$ in line 11, and every correct process would eventually exit the loop. Thus, every correct process $p_i$ belongs to $S$ and eventually writes a non-$\perp$ value in $A[r][k][i]$ (line 16). But since there are at least $n + 1 - f$ correct processes in $R$, $A[r][k]$ eventually contains at least $n + 1 - f$ non-$\perp$ entries and, thus, the condition in line 19 is eventually satisfied — a contradiction.

Thus, in every sub-round $k$ of round $r$, each correct process eventually exits the loop in lines 17–19 and, since $D$ is never $\bot$, reaches line 23 (if $D[r] \neq \bot$) or line 26 (otherwise).

Note that at most $n + 1 - |S|$ different non-$\bot$ values can be written in $D[r]$ by processes not in $S$. On the other hand, a process in $S$ is allowed to write $v$ in $D[r]$ only if it has committed on $v$ in some instance of $(f + |S| - n - 1)$-$converge[r][k]$. By the C-Agreement and Validity properties of $(f+|S|-n-1)$-convergence and the fact that every value returned by $(f+|S|-n-1)$-$converge[r][k]$ is adopted, at most $f + |S| - n - 1$ distinct values can ever written by processes in $S$. Thus, at most $n + 1 - |S| + f + |S| - n - 1 = f$ distinct values can ever be written in $D[r]$.

Suppose, $D[r] \neq \bot$ at some point in $R$. Thus, eventually every process either fails or adopts one of at most $f$ values written in $D[r]$ (line 23 or 33), and then proceeds to round $r + 1$. Hence, by the Convergence property of $f$-convergence, every correct process commits a value after invoking $f$-$converge[r + 1]$ and decides — a contradiction.

Now suppose that $D[r] = \bot$ forever. By the algorithm, there are no correct processes outside $S$ and, thus, there is at least one faulty process in $S$ (otherwise, $S$ would be the set of correct processes, violating the properties of $\Upsilon^f$). Let $k$ be a sub-round of round $r$ in which no faulty process participates (every faulty process fails before starting the sub-round). Since there is at least one faulty process in $S$, at most $|S| - 1$ values can be written in $A[r][k]$.

Now consider all sets that can be returned by $atomic\text{-}snapshot(A[r][k])$ in line 18. Every such set contains at least $n + 1 - f$ and at most $|S| - 1$ non-$\bot$ values. Moreover, by the properties of atomic snapshot [1], all these sets are related by containment. Thus, there can be at most $|S| - 1 - (n + 1 - f) + 1 = |S| + f - n - 1$ distinct sets, and, thus, at most $|S| + f - n - 1$ different values can be computed by the processes in line 25. Hence, by the Convergence property of $(|S| + f - n - 1)$-$converge[r][k]$, every correct process that invokes the operation, commits on a value and writes it in $D[r]$ — a contradiction.

Thus, eventually, every correct process decides. $\qquad\square$