



HAL
open science

Revisiting Immediate Snapshot

Carole Delporte, Hugues Fauconnier, Sergio Rajsbaum, Michel Raynal

► **To cite this version:**

Carole Delporte, Hugues Fauconnier, Sergio Rajsbaum, Michel Raynal. Revisiting Immediate Snapshot. 2016. hal-01315342v2

HAL Id: hal-01315342

<https://inria.hal.science/hal-01315342v2>

Preprint submitted on 31 May 2016 (v2), last revised 27 Jun 2016 (v4)

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.

Revisiting Immediate Snapshot

Carole Delporte[†], Hugues Fauconnier[†], Sergio Rajsbaum[°], Michel Raynal^{*,‡}

[†] IRIF, Université Paris Diderot, Paris, France

[°] Instituto de Matemáticas, UNAM, México D.F, 04510, México

* Institut Universitaire de France

[‡] IRISA, Université de Rennes, 35042 Rennes, France

Tech Report #2036, 14 pages, May 2016
IRISA, University of Rennes 1, France

May 30, 2016

Abstract

Immediate snapshot is the basic communication object on which relies the read/write distributed computing model made up of n crash-prone asynchronous processes, called *iterated* distributed model. Each iteration step (usually called a *round*) uses a new immediate snapshot object, which allows the processes to communicate and cooperate. More precisely, the x -th immediate snapshot object can be used by a process only when it executes the x -th round. An immediate snapshot object can be implemented by an $(n - 1)$ -resilient algorithm, i.e. an algorithm that tolerates up to $(n - 1)$ process crashes (also called wait-free algorithm).

Considering a t -crash system model (i.e. a model in which up to t processes are allowed to crash), this paper is on the construction of an extension of immediate snapshot objects to t -resiliency. In the t -crash system model, at each round each process may be ensured to get values from at least $n - t$ processes, and t -immediate snapshot has the properties of classical immediate snapshot (1-immediate snapshot) but ensures that each process will get values from at least $n - t$ processes. Its main result is the following. While there is a (deterministic) t -resilient read/write-based algorithm implementing t -immediate snapshot in a t -crash system when $t = n - 1$, there is no t -resilient algorithm in a t -crash model when $t \in [1..(n-2)]$. This means that the notion of t -resilience is inoperative when one has to implement immediate snapshot for these values of t : the model assumption “at most $t < n - 1$ processes may crash” does not provide us with additional computational power allowing for the design of genuine t -resilient algorithms (genuine meaning that such a t -resilient algorithm would work in the t -crash model, but not in the $(t + 1)$ -crash model). To show these results, the paper relies on well-known distributed computing agreement problems such as consensus and k -set agreement.

Keywords: Asynchronous system, Atomic read/write register, Consensus, Distributed computability, Immediate snapshot, Impossibility, Iterated model, k -Set Agreement, Linearizability, Process crash failure, Snapshot object, t -Resilience, Wait-freedom.

1 Introduction

The iterated immediate snapshot system model The *iterated immediate snapshot* (IIS) model is a distributed computing model introduced by Borowsky and Gafni in [5, 7]. It consists of n asynchronous processes, among which any subset of up to $(n - 1)$ processes may crash¹, which execute a sequence of asynchronous rounds. One and only one immediate snapshot (IS) object is associated with each round, which allows the processes to communicate with during this round. More precisely, for any $x > 0$, a process accesses the x -th immediate snapshot only when it executes the x -th round, and it accesses it only once.

From an abstract point of view, an IS object $IMSP$, can be seen as an initially empty set, which can then contain at most n pairs (one per process), each made up of a process index and a value. This object provides the processes with a single operation denoted `write_snapshot()`, that each process may invoke only once. The invocation $IMSP.write_snapshot(v)$ by a process p_i adds the pair $\langle i, v \rangle$ to $IMSP$ and returns a set of pairs belonging to $IMSP$ such that the sets returned to the processes that invoke `write_snapshot()` satisfy specific inclusion properties. It is important to notice that, in the IIS model, the processes access the sequence of IS objects one after the other, in the same order, and asynchronously.

The noteworthy feature of the IIS model is the following. It has been shown by Borowsky and Gafni in [7], that this model is equivalent to the usual read/write wait-free model ($(n - 1)$ -crash model) for task solvability with the wait-freedom progress condition (any non-faulty process obtains a result). Its advantage lies in the fact that its runs are more structured and easier to analyze than the runs in the basic read/write shared memory model [23]. It is also the basis of combinatorial topology approach of [15] for distributed systems. Hence, IS objects constitute the algorithmic foundation of distributed iterated computing models.

It has been shown in [26] that trying to enrich the IIS model with (non trivial) failure detectors is inoperative. This means that, for example, enriching IIS with the failure detector Ω (which is the weakest failure detector that allows consensus to be solved in the basic read/write communication model [10, 20]) does not allow to solve consensus in such an enriched IIS model. However, it has been shown in [25] that it is possible to capture the power of a failure detector (and other partially synchronous systems) in the IIS model by appropriately restricting its set of runs, giving rise to the *Iterated Restricted Immediate Snapshot* (IRIS) model. This approach has been further investigated in [28].

The IIS model has many interesting features among which the following two are noteworthy. The first is on the foundation side of distributed computing, namely IIS established a strong connection linking distributed computing and topology [6, 15, 16, 17, 29]. The second one lies on the algorithmic and programming side, namely IIS allows for a recursive formulation of algorithms solving distributed computing problems. This direction, initiated in [5, 13], has also been investigated in [24, 27].

t -Crash model and t -resilient algorithms The previous basic read/write model and IIS model consider that all but one process may crash. Differently, a t -crash model assumes that at most t processes may crash, i.e., by assumption, at least $(n - t)$ of them never crash. As already said, an algorithm designed for such a model is said to be t -resilient.

One of the most fundamental results of distributed computing is the impossibility to design a 1-resilient consensus algorithm in the 1-crash n -process model, be the communication medium an asynchronous message-passing system [12] or a read/write shared memory [21]. Differently, other problems, such as renaming (introduced in the context of t -resilient message-passing systems where $t < n/2$ [3]), can be solved by $(n - 1)$ -resilient algorithms in the $(n - 1)$ -crash read/write shared memory model (such

¹From a terminology point of view, we say *t-failure model* (in the present case *t-crash model*) if the model allows up to t processes to fail. We keep the term *t-resilience* for algorithms. The $(n - 1)$ -crash model is also called *wait-free* model [14]. Several progress conditions have been associated with $(n - 1)$ -resilient algorithms: wait-freedom, non-blocking, or obstruction-freedom.

renaming algorithms are described in several textbooks, e.g. [4, 27, 30]).

Contribution of the paper When considering the t -crash n -process model where $t < n - 1$, and assuming that each correct process writes a value, a process may wait for values written by $(n - t)$ processes without risking being blocked forever. This naturally leads to the notion of a t -crash n -process iterated model, generalizing the IIS model to any value of t .

To this end the paper introduces the notion of a k -immediate snapshot object, which generalizes the immediate snapshot object. More precisely, when considering a t -immediate snapshot object in a t -crash n -process model, an invocation of `write_snapshot()` by a process returns a set including at least $(n - t)$ pairs (while it would return a set of x pairs with $1 \leq x \leq n$ if the object was an IS object). Hence, a t -immediate snapshot object allows processes to obtain as much information as possible from the other processes while guaranteeing progress.

The obvious question is then the implementability of a t -immediate snapshot object in the t -crash n -process model. This question is answered in this paper, which shows that it is impossible to implement a t -IS object in a t -crash n -process model when $0 < t < n - 1$. More precisely we prove that implementing t -IS object is equivalent² to implement consensus when $t < n/2$ and enables to implement $(2t - n + 2)$ set agreement when $n/2 \leq t < n - 1$.

At first glance, this impossibility result may seem surprising. An IS object is a snapshot object (a) whose operations `write()` and `snapshot()` are glued together in a single operation `write_snapshot()`, and (b) satisfying an additional property linking the sets of pairs returned by concurrent invocations (called *Immediacy* property, Section 2.2). Then, as already indicated, a t -IS object is an IS object such that the sets returned by `write_snapshot()` contain at least $(n - t)$ pairs (*Output size* property, Section 2.4). The same Output size property on the sets returned by a snapshot object can be trivially implemented in a t -crash n -process model. Let us call t -snapshot such a constrained snapshot object. Hence, while a t -snapshot object can be implemented in the t -crash n -process model, a t -IS object cannot when $0 < t < n - 1$.

Roadmap As previously indicated, the paper is on the computability power of t -IS objects in the t -crash computing model, for $t < n - 1$. Made up of 6 sections, it has the following content.

- Section 2 introduces the basic crash-prone read/write system model, immediate snapshot, a k -set agreement, and k -immediate snapshot (k -IS). It also proves a theorem which captures the additional computational power of k -immediate snapshot with respect to immediate snapshot.
- Assuming a majority of processes never crash, i.e. a t -crash read/write model in which $t < n/2$, Section 3 shows that it is impossible to implement t -immediate snapshot in such a model. The proof is a reduction of the consensus problem to t -immediate snapshot.
- Assuming a t -crash read/write model in which $n/2 \leq t < n - 1$, Section 4 shows that it is impossible to implement t -immediate snapshot in such a model. The proof is a reduction of the $(2t - n + 2)$ -set agreement problem to t -immediate snapshot.
- By a simulation argument, Section 5 shows that consensus is not solvable with t -immediate snapshot when $n/2 \leq t < n$ proving that the computational power of t -immediate snapshot when $0 < t < n/2$ is strictly stronger than the computational power of t -immediate snapshot when $n/2 \leq t < n$.

Finally, Section 6 concludes the paper.

²A is equivalent to B if A can be (computationally) reduced to B and reciprocally.

2 Immediate Snapshot, k -Set Agreement, and k -Immediate Snapshot

2.1 Basic read/write system model

Processes The computing model is composed of a set of $n \geq 3$ sequential processes denoted p_1, \dots, p_n . Each process is asynchronous which means that it proceeds at its own speed, which can be arbitrary and remains always unknown to the other processes.

A process may halt prematurely (crash failure), but executes correctly its local algorithm until it possibly crashes. The model parameter t denotes the maximal number of processes that may crash in a run. A process that crashes in a run is said to be *faulty*. Otherwise, it is *correct* or *non-faulty*. Let us notice that, as a faulty process behaves correctly until it crashes, no process knows if it is correct or faulty. Moreover, due to process asynchrony, no process can know if an other process crashed or is only very slow.

It is assumed that (a) $0 < t < n$ (at least one process may crash and at least one process does not crash), and (b) any process, until it possibly crashes, executes the algorithm assigned to it.

Communication layer The processes cooperate by reading and writing Single-Writer Multi-Reader (SWMR) atomic read/write registers [19]. This means that the shared memory can be seen as a set of arrays $A[1..n]$ where, while $A[i]$ can be read by all processes, it can be written only by p_i .

Notation The previous computation model is denoted $\mathcal{CARW}_{n,t}[\emptyset]$ (Crash Asynchronous Read/Write). A model constrained by a predicate on t (e.g. $t < x$) is denoted $\mathcal{CARW}_{n,t}[t < x]$. Hence, as we assume at least one process does not crash, $\mathcal{CARW}_{n,t}[t < n]$ is a synonym of $\mathcal{CARW}_{n,t}[\emptyset]$, which (as always indicated) is called *wait-free* model. When considering t -crash models, $\mathcal{CARW}_{n,t}[t \leq \alpha]$ is less constrained than $\mathcal{CARW}_{n,t}[t < \alpha - 1]$.

Shared objects are denoted with capital letters. The local variables of a process p_i are denoted with lower case letters, sometimes suffixed by the process index i .

2.2 One-shot immediate snapshot object

The immediate snapshot (IS) object was informally presented in the introduction. It can be seen as a variant of the snapshot object introduced in [1, 2]. While a snapshot object provides the processes with two operations (`write()` and `snapshot()`) which can be invoked separately by a process (usually `write()` before `snapshot()`), a immediate snapshot provides the processes with a single operation `write_snapshot()`. One-shot means that a process may invoke `write_snapshot()` at most once.

Definition Let $IMSP$ be an IS object. It is a set, initially empty, that will contain pairs made up of a process index and a value. Let us consider a process p_i that invokes $IMSP.write_snapshot(v)$. This invocation adds the pair $\langle i, v \rangle$ to $IMSP$ (contribution of p_i to $IMSP$), and returns to p_i a set, called view and denoted $view_i$, such that the sets returned to the processes collectively satisfy the following properties.

- Termination. A correct process that invokes `write_snapshot()` returns from its invocation.
- Self-inclusion. $\forall i : \langle i, v \rangle \in view_i$.
- Validity. $\forall i : (\langle j, v \rangle \in view_i) \Rightarrow p_j$ invoked `write_snapshot(v)`.
- Containment. $\forall i, j : (view_i \subseteq view_j) \vee (view_j \subseteq view_i)$.
- Immediacy. $\forall i, j : (\langle i, v \rangle \in view_j) \Rightarrow (view_i \subseteq view_j)$.

Immediacy can be re-stated as: $\forall i, j : ((\langle i, - \rangle \in view_j) \wedge (\langle j, - \rangle \in view_i)) \Rightarrow (view_i = view_j)$.

Implementations of an IS object in the wait-free model $\mathcal{CARW}_{n,t}[t < n]$ are described in [5, 13, 24, 27]. While both a one-shot snapshot object and an IS object satisfy the Self-inclusion, Validity and Containment properties, only an IS object satisfies the Immediacy property. This additional property creates an important difference, from which follows that, while a snapshot object is atomic (operations on an IS object can be linearized [18]), an IS object is not atomic (its operations cannot always be linearized). However, an IS object is set-linearizable (set-linearizability allows several operations to be linearized at the same point of the time line [9, 22]).

The iterated immediate snapshot (IIS) model In this model (introduced in [7]), the shared memory is composed of a (possibly infinite) sequence of IS objects: $IMSP[1], IMSP[2], \dots$. These objects are accessed sequentially and asynchronously by the processes according to the following round-based pattern executed by each process p_i . The variable r_i is local to p_i ; it denotes its current round number.

```

 $r_i \leftarrow 0; \ell_{s_i} \leftarrow$  initial local state of  $p_i$  (including its input, if any);
repeat forever % asynchronous IS-based rounds
   $r_i \leftarrow r_i + 1;$ 
   $view_i \leftarrow IMSP[r_i].write\_snapshot(\ell_{s_i});$ 
  computation of a new local state  $\ell_{s_i}$  (which contains  $view_i$ )
end repeat.

```

As indicated in the Introduction, when considering distributed tasks (as formally defined in [8, 17]), the IIS model and $\mathcal{CARW}_{n,t}[t < n]$ have the same computational power [7].

2.3 k -Set agreement

k -Set agreement was introduced by S. Chaudhuri [11] to investigate the relation linking the number of different values that can be decided in an agreement problem, and the maximal number of faulty processes. It generalizes consensus which corresponds to the case $k = 1$.

A k -set agreement object is a one-shot object that provides the processes with a single operation denoted $propose_k()$. This operation allows the invoking process p_i to propose a value it passes as an input parameter (called *proposed* value), and obtain a value (called *decided* value). The object is defined by the following set of properties.

- Termination. If a process invokes $propose_k()$ and does not crash, it returns from its invocation.
- Validity. A decided value is a proposed value.
- Agreement. No more than k different values are decided.

It is shown in [6, 17, 29] that the problem is impossible to solve in $\mathcal{CARW}_{n,t}[k \leq t]$.

2.4 k -Immediate Snapshot

A k -immediate snapshot object (denoted k -IS) is an immediate snapshot object with the following additional property.

- Output size. The set $view$ obtained by a process is such that $|view| \geq n - k$.

Theorem 1 *A k -IS object cannot be implemented in $\mathcal{CARW}_{n,t}[k < t]$.*

Proof To satisfy the output size property, the view obtained by a process p_i must contain pairs from $(n - k)$ different processes. If t processes crash (e.g. initially) a process can obtain at most $(n - t)$ pairs.

If $t > k$, we have $n - t < n - k$. It follows that, after it has obtained pairs from $(n - t)$ processes, a process can remain blocked forever waiting for the $(t - k)$ missing pairs. $\square_{\text{Theorem 1}}$

Considering the system model $\mathcal{CARW}_{n,t}[t < n - 1]$, the next theorem characterizes the power of a t -IS object in term of the Containment property.

Theorem 2 *Considering the system model $\mathcal{CARW}_{n,t}[t < n - 1]$, and a t -IS object, let us assume that all correct processes invoke `write_snapshot()`. No process obtains a view with less than $(n - t)$ pairs. Moreover, if the size of the smallest view obtained by a process is ℓ ($\ell \geq n - t$), there is a set S of processes such that $|S| = \ell \geq n - t$ and each process of S obtains the smallest view or crashes during its invocation of `write_snapshot()`.*

Proof It follows from the Output size property of the t -IS object that no view contains less than $(n - t)$ pairs. Let $view$ be the smallest view returned by a process, and let $\ell = |view|$. We have $\ell \geq n - t$. Moreover, due to (a) the Immediacy property (namely $(\langle i, - \rangle \in view) \Rightarrow (view_i \subseteq view)$) and (b) the minimality of $view$, it follows that $view_i = view$. As this is true for each process whose pair participates in $view$, and $\ell = |view|$, it follows that there is a set S of processes such that $|S| = \ell \geq n - t$ and each of its processes obtains the view $view$, or crashed during its invocation of `write_snapshot()`. Due to the Containment property, the others processes crash or obtain views which strictly include $view$. $\square_{\text{Theorem 2}}$

3 t -Immediate Snapshot is Impossible in $\mathcal{CARW}_{n,t}[0 < t < n/2]$

This section shows that it is impossible to implement a t -IS object when $0 < t < n/2$.

From t -IS to consensus in $\mathcal{CARW}_{n,t}[t < n/2]$ Algorithm 1 reduces consensus to t -IS in the system model $\mathcal{CARW}_{n,t}[t \leq n/2]$. As at most $t < n/2$ process may crash, at least $n - t > n/2$ processes invoke the consensus operation `propose1`(v).

operation `propose1`(v) is

- (1) $view_i \leftarrow IMSP.write_snapshot(v); VIEW[i] \leftarrow view_i;$
- (2) `wait`($|\{j \text{ such that } VIEW[j] \neq \perp\}| = t + 1$);
- (3) **let $view$ be** the smallest of the previous $(t + 1)$ views;
- (4) `return`(smallest proposed value in $view$)

end operation.

Algorithm 1: Solving consensus in $\mathcal{CARW}_{n,t}[t < n/2, t\text{-IS}]$ (code for p_i)

In addition to a t -IS object denoted $IMSP$, the processes access an array $VIEW[1..n]$ of SWMR atomic registers, initialized to $[\perp, \dots, \perp]$. The aim of $VIEW[i]$ is to store the view obtained by p_i from the a t -IS object $IMSP$,

When it calls `propose1`(v), a process p_i invokes first the t -IS object, in which it deposits the pair $\langle i, v \rangle$, and obtains a view from it, that it writes in $VIEW[i]$ to make it publicly known. (line 1). Then, it waits (line 2) until it sees the views of at least $(t + 1)$ processes (as $n - t \geq t + 1$, p_i cannot block forever and at least one of these views is from a correct process). Process p_i extracts then of these views the one with the smallest cardinality (line 3), and finally returns proposed value contained in this smallest view (line 4).

Theorem 3 *Algorithm 1 reduces consensus to t -IS in $\mathcal{CARW}_{n,t}[t \leq n/2]$.*

Proof Let us first prove the consensus Termination property. As $n-t \geq t+1$, and there are at least $(n-t)$ correct processes, it follows that at least $(n-t)$ entries of $VIEW[1..n]$ are eventually different from \perp . Hence, no correct process can remain blocked forever at line 2, which proves consensus Termination.

Let us now consider the consensus Agreement property. It follows from Theorem 2 that there is a set of at least $\ell \geq n-t$ processes, that obtained the same view min_view (or crashed before returning from $write_snapshot()$), and this view is the smallest view obtained by a process and its size is $|min_view| = \ell$. As $\ell \geq n-t$ and $(n-t) + (t+1) > n$, it follows from the waiting predicate of line 2, that, any process that executes line 3, obtains a copy of min_view , and consequently we have $view = min_view$ at line 3. It follows that no two processes can decide different values.

Finally, the consensus Validity property follows from the fact that any pair contained in a view is composed of a process index and the value proposed by the corresponding process. $\square_{Theorem\ 3}$

Corollary 1 *Implementing a t -IS object in $\mathcal{CARW}_{n,t}[t < n/2]$ is impossible.*

Proof The proof is an immediate consequence of Lemma 3, and the fact that consensus cannot be solved in $\mathcal{CARW}_{n,t}[t < n/2]$ [21]. $\square_{Corollary\ 1}$

With consensus it is easy to implement k -immediate snapshot: each process writes its value proposed to k -immediate snapshot in its own register and read these registers until it gets at least $n-k$ values, then by successive consensus processes agree on increasing set of values. The first time a consensus decides a set containing (p, v) then p adopts that set as the output of the k -immediate snapshot. Then we get:

Theorem 4 *consensus is equivalent to k -immediate snapshot in $\mathcal{CARW}_{n,t}[t < n/2]$*

4 t -Immediate Snapshot is Impossible in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1]$

This section shows that it is impossible to implement a t -IS object in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1]$. To this end, it presents a reduction of k -set agreement (in short k -SA) to t -IS for $k = 2t - n + 2$ (e.g., a reduction of $(n-2)$ -SA agreement to $(n-2)$ -IS in $\mathcal{CARW}_{n,t}[t = n-2]$).

From t -IS to $(2t-k+2)$ -set agreement in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1, t$ -IS] Algorithm 2 reduces $(2t-n+2)$ -set agreement to t -IS in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1]$. As at most t process may crash, at least $(n-t)$ processes invoke the k -SA operation $propose_k()$. This algorithm is very close to Algorithm 1. Its main difference lies in the replacement of $(t+1)$ by $(n-t)$ at line 2.

operation $propose_{2t-n+2}(v)$ **is**
(1) $view_i \leftarrow IMSP.write_snapshot(v); VIEW[i] \leftarrow view_i;$
(2) wait($|\{j \text{ such that } VIEW[j] \neq \perp\}| = n-t$);
(3) **let** $view$ **be** the smallest of the previous $(n-t)$ views;
(4) return(smallest proposed value in $view$)
end operation.

Algorithm 2: Solving $(2t-n+2)$ -set agreement in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1, t$ -IS] (code for p_i)

Theorem 5 *Algorithm 2 reduces $(2t-n+2)$ -set agreement to t -IS in $\mathcal{CARW}_{n,t}[n/2 \leq t < n-1]$.*

Proof Let $k = 2t - n + 2$.

Let us first consider the k -SA Termination property. There are at least $(n-t)$ correct processes, and each of them first invokes $IMSP.write_snapshot()$ and then writes the view it obtained in the shared

array $VIEW$ (line 1). Hence, at least $(n - t)$ entries of $VIEW$ are eventually different from \perp , from which follows that no process can block forever at line 2.

Let us now consider the k -SA Validity property. It follows from the Containment property of the t -IS object that any set of views deposited in $VIEW$ is not empty. Therefore, the view selected by a process at line 3 is not empty. As a view can only contain pairs, each including a proposed value (line 1), the k -SA Validity property follows.

Let us finally consider the k -SA Agreement property. Let us first observe that, due to the t -IS Containment property and Theorem 2, at most $n - (n - t) + 1 = t + 1$ different views can be written in the array $VIEW[1..n]$. Let $V(1)$ the smallest of these views (which contains $\ell \geq n - t$ pairs), $V(2)$ the second smallest, etc., until $V(t + 1)$ the greatest one. There are two cases according to the $(n - t)$ non- \perp views obtained by a process p_i at line 2. Let us remind that, as $n \leq 2t$, we have $n - t \leq t$.

- Case 1. The view $V(1)$ belongs to the $(n - t)$ views obtained by p_i . In this case, p_i selects $V(1)$ at line 3 and decides at line 4 the smallest proposed value contained in $V(1)$.
- Case 2. The view $V(1)$ does not belong to the $(n - t)$ views obtained by p_i . Hence, the $(n - t)$ views obtained by any process of Case 2 belong to $\{V(2), \dots, V(t + 1)\}$.

It follows that the $m = (n - t) - 1$ biggest views in $\{V(2), \dots, V(t + 1)\}$ will never be selected by the processes that are in Case 2, and consequently the set of these processes obtain at most $t - m = t - ((n - t) - 1) = 2t - n + 1$ different smallest views. Hence, these processes may decide at most $2t - n + 1$ different values at line 4.

When combining the two cases, at most $k = 2t - n + 2$ different values can be decided, which concludes the proof of the theorem. $\square_{\text{Theorem 5}}$

Corollary 2 *Implementing a t -IS object in $\mathcal{CARW}_{n,t}[n/2 \leq t < n - 1]$ is impossible.*

Proof As $t \leq n - 2$, we have $2t - n + 2 \leq t$. The proof is an immediate consequence of Theorem 5, and the fact that $(2t - n + 2)$ -set agreement cannot be solved in $\mathcal{CARW}_{n,t}[n/2 \leq t < n - 1]$ [5, 17, 29].

$\square_{\text{Corollary 2}}$

5 t -Immediate Snapshot and Consensus in $\mathcal{CARW}_{n,t}[n/2 \leq t < n - 1]$

Let us first remark that (as immediate snapshot objects) k -immediate snapshot objects are not linearizable. As a t -immediate snapshot o contains values from at least $(n - t)$ processes, at least $(n - t)$ processes must have invoked the operation `write_snapshot()` on o for any invocation of `write_snapshot()` to be able to terminate. It follows that there is a time τ at which $(n - t)$ processes have invoked the operation `write_snapshot()` on the k -immediate snapshot o and have not yet returned. We then say that these $(n - t)$ processes are *inside* their k -immediate snapshot o . Hence the following lemma:

Lemma 1 *If an invocation of `write_snapshot()` on a k -immediate snapshot object o terminates, there is a time τ at which at least $(n - t)$ processes are inside this k -immediate snapshot object o .*

Theorem 6 *There is no t -resilient consensus algorithm using t -immediate snapshot in $\mathcal{CARW}_{n,t}[n/2 \leq t < n - 1]$.*

Proof To prove the theorem, let us consider first the case $n = 2t$. The proof is by contradiction. Let us assume that \mathcal{A} is a t -resilient consensus algorithm for a set of processes $\{p_1, \dots, p_n\}$ which use a

t -immediate snapshot object in a system where $n = 2t$. The contradiction is obtained by simulating \mathcal{A} with two processes Q_0 and Q_1 , such that Q_0 and Q_1 solve consensus despite the possible crash of one of them. As there is no wait-free consensus algorithm for 2 processes, it follows that such a consensus algorithm \mathcal{A} based on t -immediate snapshot objects does not exist. The simulation is described in Algorithm 3.

Let A_0 and A_1 be a partition of $\{p_1, \dots, p_n\}$: $|A_0| = |A_1| = t$, $\{p_1, \dots, p_n\} = A_0 \cup A_1$, and $A_0 \cap A_1 = \emptyset$.

Code for Q_i ($i \in \{0, 1\}$):

- (1) **for all** p_j in A_i : initialize v_{p_j} with the initial value of Q_i ;
- (2) **repeat forever**
- (3) **for each** p in A_i in a round robin way **do**
- (4) **if** the next step of p is $is(o, v)$ (i.e. `write_snapshot(v)` on the immediate snapshot object o)
- (5) **then** $prop_i[o] \leftarrow prop_i[o] \cup \{(p, v)\}$;
- (6) **if** $REG[i][o] = \perp$
- (7) **then if** $REG[1-i][o] \neq \perp$
- (8) **then** $REG[i][o] \leftarrow REG[1-i][o] \cup \{(p, v)\}$;
- (9) simulation step $is(o, v)$ for p which returns $REG[i][o]$
- (10) **end if**
- (11) **else** $REG[i][o] \leftarrow REG[i][o] \cup \{(p, v)\}$;
- (12) simulation step $is(o, v)$ for p which returns $REG[i][o]$
- (13) **end if**
- (14) **else** simulate the next step of p ;
- (15) **if** p decides v in this step **then** Q_i decides v **end if**
- (16) **end if**;
- (17) **if** $(|prop_i(o)| = t) \wedge (REG[i][o] = \perp)$ **then** $REG[i][o] \leftarrow IMSP[o].write_snapshot(prop_i(o))$
- (18) **end if**
- (19) **end for**
- (20) **end repeat.**

Algorithm 3: Simulation of \mathcal{A} by Q_i ($i \in \{0, 1\}$) for $n = 2t$

Let A_0 and A_1 be a partition of $\{p_1, \dots, p_n\}$ such that each of A_0 and A_1 has t elements. Q_0 simulates the processes in A_0 , while Q_1 simulates the processes in A_1 . In the simulation, if Q_i is correct and makes an infinite number of steps, then each process in A_i makes an infinite number of (simulated) steps, and consequently the processes of A_i are correct in the simulated run. If Q_i crashes, its crash entails (in the simulated run) the crashes of all the processes in A_i . Note that, as at most t simulated processes may crash in a simulated run, if all processes of A_i crash, no process of A_{1-i} crashes.

In the following, given a simulated process p , $is(o, v)$ denotes the invocation of `write_snapshot(v)` by p on the t -immediate snapshot o . We assume the t -immediate snapshot objects are one-shot objects (each process invokes an object o at most once). The underlying idea of the simulation is that a 1-immediate snapshot object accessed by Q_0 and Q_1 allows them to simulate a t -immediate snapshot object shared by the simulated processes p_1, \dots, p_n .

The 1-immediate snapshot object associated with the simulated t -immediate snapshot object o , is denoted $IMSP[o]$. In addition to these 1-immediate snapshot objects, the simulator processes Q_0 and Q_1 of the simulation Algorithm 3 manage the following variables.

- $REG[0, 1][o]$ is an array made up of two atomic read/write registers associated with each simulated t -immediate snapshot object o . $REG[i][o]$ is written by Q_i and read by both Q_i and Q_{1-i} . It contains (at least) the values written in o by the processes simulated by Q_i (lines 8 and 11). If Q_i has not already simulated an immediate snapshot operation on o while Q_{1-i} has, $REG[i][o]$ is initialized to the result of the immediate snapshot on o made by the processes of A_{1-i} simulated by Q_{1-i} (lines 6-8).
- $prop_i[o]$ is a local variable of Q_i containing the values written in the t -immediate snapshot o by the simulated processes in A_i (line 5). When the next step of all the simulated processes is a t -

immediate snapshot on o , Q_i gives the initial value of $REG[i][o]$ (line 17). In the next t executions of the loop, when Q_i considers the simulated process p , this value will be returned to p (line 12) by the simulation of immediate snapshot invocation on o issued by p .

The central point of the simulation lies in the way the t -immediate snapshot objects are simulated. For this, only when the next step of *all* the simulated processes in A_i are $o.write_snapshot()$ (t -immediate snapshot operation on the *same* object o) the simulator Q_i performs an immediate snapshot on the corresponding 1-immediate snapshot object $IMSP[o]$ shared by Q_0 and Q_1 , with the values written by the processes in A_i in this t -immediate snapshot on o . The result of this immediate snapshot contains either all the values from all simulated processes, or only the values of the processes in A_i . Moreover, all processes of Q_i obtain the same result, and Q_i also writes this result value into $REG[i, o]$ (line 17).

Let us now consider the case in which the next step of the processes in A_i is not a t -immediate snapshot operation on the same object. If the next step of some process $p \in A_i$ is a t -immediate snapshot on object o and no t -immediate snapshot on o by processes in A_i have already returned from their invocations, we prove that there is a time τ at which all processes in A_0 , or all processes A_1 , are *inside* the t -immediate snapshot object o . To this end, let us assume that there is no time at which all processes in A_i are inside a t -immediate snapshot object o . By Lemma 1 there is a time τ at which a set of at least t processes, say C , are inside a t -immediate snapshot o . At this time, as –by assumption– at least one process in A_i is not inside a t -immediate snapshot, it follows that at least one process of A_{1-i} is inside a t -immediate snapshot. But let us then consider the run in which all processes in A_i crash (in particular all processes in A_i may be considered as crashed before they invoked the t -immediate snapshot). Hence for this run, C contains no process in A_i and, as $|C| \geq t$, C is equal to A_{1-i} .

From this observation we deduce that either there is a time for which the next step of all $p \in A_i$ is a t -immediate snapshot on o , or there is a time at which the next step of all $p \in A_{1-i}$ is a t -immediate snapshot on o . Hence, Q_i or Q_{1-i} performs an immediate snapshot on o . If Q_{1-i} performs an immediate snapshot on o , then the result of the t -immediate snapshot on o for each processes in A_{1-i} is the set V made up of the values written by the processes in A_{1-i} . After that, Q_i can read V from a shared variable, and is able to compute the result of a t -immediate snapshot on o (the result is V union the set of values of processes in A_i for which Q_i has simulated the t -immediate snapshot on o). Hence, if $p \in A_i$ is stuck in the simulation on an object o , either Q_{1-i} eventually makes an immediate snapshot on o and Q_i eventually simulates the t -immediate snapshot on o for p , or eventually the next step of all processes in A_i is a t -immediate snapshot on o and Q_i can compute the result of this t -immediate snapshot on o .

To extend the result to $2t > n$, we partition $\{p_1, \dots, p_n\}$ in 3 sets A_0, A_1, D such that $|A_0| = n - t$, $|A_1| = n - t$, $|D| = 2t - n$. Then, we run the previous simulation algorithm \mathcal{A} where all processes in D are initially dead, Q_0 simulates the set of processes of A_0 , and Q_1 simulates the processes of A_1 . With this simulation, Q_0 and Q_1 realizes a wait-free consensus, which is known to be impossible. $\square_{Theorem 6}$

6 Conclusion

This paper addressed the design of t -tolerant algorithms building a t -immediate snapshot (t -IS) object. Such an object in an immediate snapshot object (defined by Termination, Self-inclusion, Containment, and Immediacy properties), in a t -crash asynchronous system. Hence, it is required that each set returned to a process contains at least $(n - t)$ pairs. Immediate snapshot corresponds to $(n - 1)$ -immediate snapshot.

The paper has shown that, while it is possible to build an $(n - 1)$ -IS object in the asynchronous read/write $(n - 1)$ -crash model, it is impossible to build a t -IS object in an asynchronous read/write t -

crash model when $0 < t < n - 1$. It follows that the notion of an IIS distributed model seems inoperative for these values of t .

Hence, we have two contrasting impossibility results in asynchronous read/write t -crash n -process systems: consensus is impossible as soon as $t > 0$, while t -immediate snapshot is impossible as soon as $t < n - 1$.

Acknowledgments

This work has been partially supported by the French ANR project DISPLEXITY devoted to the study of Computability and Complexity in distributed computing.

References

- [1] Afek Y., Attiya H., Dolev D., Gafni E., Merritt M. and Shavit N., Atomic snapshots of shared memory. *Journal of the ACM*, 40(4):873-890 (1993)
- [2] Anderson J., Multi-writer composite registers. *Distributed Computing*, 7(4):175-195 (1994)
- [3] Attiya H., Bar-Noy A., Dolev D., Peleg D., and Reischuk R., Renaming in an asynchronous environment. *Journal of the ACM*, 37(3):524-548 (1990)
- [4] Attiya H. and Welch J., *Distributed computing: fundamentals, simulations and advanced topics*, (2d Edition), Wiley-Interscience, 414 pages (2004)
- [5] Borowsky E. and Gafni E., Immediate atomic snapshots and fast renaming. *Proc. 12th ACM Symposium on Principles of Distributed Computing (PODC'93)*, pp. 41-50 (1993)
- [6] Borowsky E. and Gafni E., Generalized FLP impossibility results for t -resilient asynchronous computations. *Proc. 25th ACM Symposium on Theory of Computation (STOC'93)*, California (USA), pp. 91-100 (1993)
- [7] Borowsky E. and Gafni E., A simple algorithmically reasoned characterization of wait-free computations. *Proc. 16th ACM Symposium on Principles of Distributed Computing (PODC'97)*, ACM Press, pp. 189-198 (1997)
- [8] Borowsky E., Gafni E., Lynch N. and Rajsbaum S., The BG distributed simulation algorithm. *Distributed Computing*, 14:127-146 (2001)
- [9] Castañeda A., Rajsbaum S., and Raynal M., Specifying concurrent problems: beyond linearizability and up to tasks. *Proc. 29th Symposium on Distributed Computing (DISC'15)*, Springer LNCS 9363, pp. 420-435 (2015)
- [10] Chandra T., Hadzilacos V., and Toueg S., The weakest failure detector for solving consensus. *Journal of the ACM*, 43(4):685-722 (1996)
- [11] Chaudhuri S., More choices allow more faults: set consensus problems in totally asynchronous systems. *Information and Computation*, 105(1):132-158 (1993)
- [12] Fischer M.J., Lynch N.A., and Paterson M.S., Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374-382 (1985)
- [13] Gafni E. and Rajsbaum S., Recursion in distributed computing. *Proc. 12th Int'l Conference on Stabilization, Safety, and Security of Distributed Systems (SSS'10)*, Springer LNCS 6366, pp. 362-376 (2010)
- [14] Herlihy M. P., Wait-free synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):124-149 (1991)
- [15] Herlihy M.P., Kozlov D., and Rajsbaum S., *Distributed computing through combinatorial topology*, Morgan Kaufmann/Elsevier, 336 pages, ISBN 9780124045781 (2014)
- [16] Herlihy M., Rajsbaum S., and Raynal M., Power and limits of distributed computing shared memory models. *Theoretical Computer Science*, 509:3-24 (2013)

- [17] Herlihy M. P. and Shavit, N., The topological structure of asynchronous computability. *Journal of the ACM*, 46(6):858-923 (1999)
- [18] Herlihy M. P. and Wing J. M., Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463-492 (1990)
- [19] Lamport L., On interprocess communication, Part I: basic formalism. *Distributed Computing*, 1(2):77-85 (1986)
- [20] Lo W.-K. and Hadzilacos V., Using failure detectors to solve consensus in asynchronous shared memory systems. *Proc. 8th Int'l Workshop on Distributed Algorithms (WDAG'94)*, Springer LNCS 857, pp. 280-295 (1994)
- [21] Loui M. and Abu-Amara H., Memory requirements for agreement among unreliable asynchronous processes. *Advances in Computing Research*, 4:163-183, JAI Press (1987)
- [22] Neiger G., Set-linearizability. Brief announcement in *Proc. 13th ACM Symposium on Principles of Distributed Computing (PODC'94)*, ACM Press, page 396 (1994)
- [23] Rajsbaum S., Iterated shared memory models. *Proc. 9th Latin American Symposium on Theoretical Informatics (LATIN'10)*, Springer LNCS 6034, pp. 407-416 (2010)
- [24] Rajsbaum, S. and Raynal, M., An introductory tutorial to concurrency-related distributed recursion. *Bulletin of the European Association of TCS*, 111:57-75 (2013)
- [25] Rajsbaum S., Raynal M., and Travers C., The iterated restricted immediate snapshot model. *Proc. 14th Annual Int'l Conference on Computing and Combinatorics (COCOON'08)*, Springer LNCS 5092, pp. 487-497 (2008)
- [26] Rajsbaum, S., Raynal, M., and Travers, C., An impossibility about failure detectors in the iterated immediate snapshot model. *Information Processing Letters*, 108(3):160-164 (2008)
- [27] Raynal M., *Concurrent programming: algorithms, principles and foundations*. Springer, 515 pages, ISBN 978-3-642-32026-2 (2013)
- [28] Raynal M. and Stainer J., Increasing the power of the iterated immediate snapshot model with failure detectors. *Proc. 19th Int'l Colloquium on Structural Information and Communication Complexity (SIROCCO'12)*, Springer LNCS 7355, pp. 231-242 (2012)
- [29] 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)
- [30] Taubenfeld G., Synchronization algorithms and concurrent programming. *Pearson Prentice-Hall*, 423 pages, ISBN 0-131-97259-6 (2006)

A Building an $(n - 1)$ -IS Object in the $(n - 1)$ -Crash Model

For a completeness purpose, this appendix presents Algorithm 4, which implements an $(n - 1)$ -IS object in the $(n - 1)$ -crash model (wait-free read/write model). This algorithm is due to Borowsky and Gafni [5]. Its explanation that follows is from [27].

Algorithm 4 uses two arrays of SWMR atomic registers denoted $REG[1..n]$ and $LEVEL[1..n]$ (only p_i can write $REG[i]$ and $LEVEL[i]$). A process p_i first writes its value in $REG[i]$. Then the core of the implementation of `write_snapshot()` is based on the array $LEVEL[1..n]$. This array, initialized to $[n + 1, \dots, n + 1]$, can be thought of as a ladder, where initially a process is at the top of the ladder, namely at level $(n + 1)$. Then it descends the ladder, one step after the other, according to predefined rules until it stops at some level (or crashes). While descending the ladder, a process p_i registers its current position in the ladder in the atomic register $LEVEL[i]$.

After it has stepped down from one ladder level to the next one, a process p_i computes a local view (denoted $view_i$) of the progress of the other processes in their descent of the ladder. This view contains the processes p_j seen by p_i at the same or a lower ladder level (i.e. such that $level_i[j] \leq LEVEL[i]$).

Then, if the current level ℓ of p_i is such that p_i sees at least ℓ processes in its view (i.e. processes that are at its level or a lower level) it stops at the level ℓ of the ladder. Finally, p_i returns a set of pairs determined from the values of $view_i$. Each pair is a process index and the value written by the corresponding process.

```

operation write_snapshot( $v_i$ ) is
   $REG[i] \leftarrow v_i$ ;
  repeat  $LEVEL[i] \leftarrow LEVEL[i] - 1$ ;
    for  $j \in \{1, \dots, n\}$  do  $level_i[j] \leftarrow LEVEL[j]$  end for;
     $view_i \leftarrow \{j : level_i[j] \leq LEVEL[i]\}$ ;
  until ( $|view_i| \geq LEVEL[i]$ ) end repeat;
  return ( $\{j, REG[j]\}$  such that  $j \in view_i$ )
end operation.

```

Algorithm 4: Borowsky-Gafni's write_snapshot() algorithm (code for p_i) [5]

The set $view_i$ of a process that terminates the algorithm, satisfy the following main property: if $|view_i| = \ell$, then p_i stopped at the level ℓ , and there are ℓ processes whose current level is $\leq \ell$. From this property, follow the Self-inclusion, Containment and Immediacy properties (stated in Section 2.2).

B An Ad hoc Proof of 1-IS Impossibility in $\mathcal{CARW}_{n,t}[t = 1]$

This section provides a customized proof for the impossibility of 1-IS in $\mathcal{CARW}_{n,t}[t = 1]$ (1-resilient read/write model). The next lemma is a simple re-statement of Theorem 2 for $t = 1$.

Lemma 2 *Considering the system model $\mathcal{CARW}_{n,t}[t = 1]$, let $view_i$ be the set returned by process p_i when it invokes the 1-IS object. The sets obtained by the processes are such that:*

- (a): $\forall i : |view_i| = n$ (and consequently all sets are equal), or
- (b): $(n - 1)$ sets are equal and such that $|view_j| = n - 1$, and the other set $view_i$ is such that $|view_i| = n$ or p_i crashed before returning it.

From 1-IS to consensus in $\mathcal{CARW}_{n,t}[t = 1]$ Let $\mathcal{CARW}_{n,t}[t = 1, 1\text{-IS}]$ denote the system model $\mathcal{CARW}_{n,t}[t = 1]$ enriched with an algorithm implementing 1-IS objects. Algorithm 5 is a reduction of consensus to 1-IS in such a system model. Let us remember that, as at most one process may crash, at least $(n - 1)$ processes invokes the consensus operation $propose_1()$.

As in previous reductions, there is an array of SWMR atomic registers $VIEW[1..n]$, whose aim is to store the view obtained by the processes.

```

operation propose1( $v$ ) is
(1)  $view_i \leftarrow IMSP.write\_snapshot(v)$ ;  $VIEW[i] \leftarrow view_i$ ;
(2) if ( $|VIEW[i]| = n - 1$ ) then return( $\min(VIEW[i])$ )
(3) else wait( $\exists j \neq i : VIEW[j] \neq \perp$ );
(4) if ( $|VIEW[j]| = n$ ) then return( $\min(VIEW[i])$ )
(5) else return( $\min(VIEW[j])$ )
(6) end if
(7) end if
end operation.

```

Algorithm 5: Solving consensus in $\mathcal{CARW}_{n,t}[t = 1, 1\text{-IS}]$ (code for p_i)

The algorithm works as follows. When p_i invokes the consensus operation $propose_1(v)$, it first invokes $IMSP.write_snapshot(v)$ and deposits the view it obtains in its SWMR register $VIEW[i]$ (line 1).

If $VIEW[i]$ contains $(n - 1)$ pairs (each made up of a process index and a proposed value), p_i selects the smallest of the proposed values present in these pairs and decides it (statement `return()` at line 2). Otherwise, due to Lemma 2, $VIEW[i]$ contains n pairs. In this case, p_i waits until another process p_j obtained a view and deposited it in $VIEW[j]$ (line 3). If $VIEW[j]$ contains n pairs, it follows from Lemma 2, that no view contains less than n pairs. Hence, p_i decides the smallest proposed value contained in these n pairs (line 4). Otherwise, $VIEW[j]$ contains $(n - 1)$ pairs, and p_i decides the smallest proposed value contained in these $(n - 1)$ pairs (line 5).

Lemma 3 *Algorithm 5 reduces consensus to 1-IS in $\mathcal{CARW}_{n,t}[t = 1]$.*

Proof Due to Lemma 2 on the The size of the views obtained by the processes ($(n - 1)$ or n) There are two cases.

- The size of all the views is n (Item (a) of Lemma 2). In this case, the predicate of line 2 is false at any process, which consequently executes the “else” part of the “if” statement”. As all processes have deposited a value in the 1-IS object $IMSP$ (otherwise the view size would be less than n), the `wait()` statement of line 3 eventually terminates, and $|VIEW[j]| = n$. Hence, the predicate of line 3 is satisfied, and as all views are equal (Lemma 2), all processes decide the same value.
- The size of the views is such that a process p_k obtains a view $VIEW[k]$ with $(n - 1)$ pairs. Due to Lemma 2, $(n - 1)$ processes obtains the very same view. The predicate of line 2 is then true at any of these processes, which, as they have the same view, decide the same value when they execute the `return()` statement of line 2. The other process, say p_ℓ , is such that $|VIEW[\ell]| = n$. Hence, it executes the “else” part of the “if” statement, and (for the same reason as above) cannot block forever at line 3. As it is the only process whose view has size n , it proceeds to line 5, and decides the smallest proposed value contained in $VIEW[j]$. Due to Item (b) of Lemma 2, this is the value decided by the $(n - 1)$ other processes, which obtained a view of size $(n - 1)$.

It follows that, in both cases, each correct process decides (Termination), no two different values are decided (Agreement), and the decided value is a proposed value (Validity). $\square_{\text{Lemma 3}}$

Theorem 7 *Implementing a 1-IS object in $\mathcal{CARW}_{n,t}[t = 1]$ is impossible.*

Proof The proof is an immediate consequence of Theorem 3, and the fact that consensus cannot be solved in $\mathcal{CARW}_{n,t}[t = 1]$ [21]. $\square_{\text{Theorem 7}}$