

# Communication Predicates: A High-Level Abstraction for Coping with Transient and Dynamic Faults<sup>\*†</sup>

Martin Hutle  
*martin.hutle@epfl.ch*

André Schiper  
*andre.schiper@epfl.ch*

*École Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland*

## Abstract

*Consensus is one of the key problems in fault tolerant distributed computing. A very popular model for solving consensus is the failure detector model defined by Chandra and Toueg. However, the failure detector model has limitations. The paper points out these limitations, and suggests instead a model based on communication predicates, called HO model. The advantage of the HO model over failure detectors is shown, and the implementation of the HO model is discussed in the context of a system that alternates between good periods and bad periods. Two definitions of a good period are considered. For both definitions, the HO model allows us to compute the duration of a good period for solving consensus. Specifically, the model allows us to quantify the difference between the required length of an initial good period and the length of a non initial good period.*

## 1. Introduction

Consensus is one of the key problems in fault tolerant distributed computing. Consensus is related to replication and appears when implementing atomic broadcast, group membership, etc. The problem is defined over a set of processes  $\Pi$ , where each process  $p_i \in \Pi$  has an initial value  $v_i$ : All processes must agree on a common value that is the initial value of one of the processes.

Consensus can be impossible to solve, as established by the FLP impossibility result [13]. Later it has been shown that consensus can be solved in a partially synchronous system with a majority of correct processes [12]. Roughly speaking, a partially synchronous system is a system that is initially asynchronous, but eventually becomes synchronous.<sup>1</sup> Moreover, in a partially synchronous system

<sup>\*</sup>Research funded by the Swiss National Science Foundation under grant number 200021-111701.

<sup>†</sup>This is an updated version of the paper that appeared at DSN 2007, including the full proofs in the appendix.

<sup>1</sup>This is not the only definition of a partially synchronous system.

links are initially lossy, but eventually become reliable.

The notion of failure detectors has been suggested a few years later [5]. The failure detector model is defined as an asynchronous system “augmented” with failure detectors, which are defined by some completeness and some accuracy property (see [5] for details). Over the years failure detectors have become very popular. The model is today widely accepted and has become the model mostly used for expressing consensus algorithms. However, the failure detector model has limitations.

First, failure detectors are not an abstraction of the partially synchronous model (even though this claim has sometimes been made). The reason is that in the partially synchronous model links are initially lossy, while the use of failure detector to solve a problem requires perpetual reliable links.<sup>2</sup> When using failure detectors, either the system must provide reliable links, or reliable links need to be implemented on top of the unreliable system links. As a consequence, the capability of algorithms of tolerating message loss — as it is the case for the Paxos algorithm [19] — cannot be expressed naturally in the failure detector model. Only a variant of Paxos that assumes reliable links can be expressed using failure detectors, as done, e.g., in [4].

Second, failure detectors are not well suited to solve consensus in the crash-recovery model, with or without stable storage [1]. In the crash-recovery model, a process can crash and later recover. This is in contrast to the crash-stop model, in which process crashes are permanent. Intuitively, one would think that solving consensus in the crash-stop model or in a crash-recovery model should not lead to major algorithmic differences. However, the comparison of (i) the  $\diamond S$  consensus algorithm in the crash-stop model [5] with (ii) the corresponding algorithm in the crash-recovery model with stable storage [1] shows that the crash-recovery algorithm is a much more complicated protocol than the corresponding crash-stop algorithm. Moreover, the complexity of the crash-recovery consensus algorithm makes it

<sup>2</sup>Failure detectors lead to the following programming pattern: Process  $p$  (i) waits for a message from process  $q$  or (ii) suspects  $q$ . If  $q$  is not suspected while the message is lost,  $p$  is blocked.

hard to see that the crash-recovery algorithm is based on the same basic ideas as the crash-stop algorithm. This leads to the following question: Is there an inherent gap between the crash-stop and the crash-recovery model that would explain the higher complexity of the crash-recovery consensus algorithm?

Third, failure detectors cannot handle Byzantine failures. The reason is that the definition of a Byzantine behavior is related to an algorithm: It is impossible to achieve a complete separation of failure detectors from the algorithm using them. To overcome this problem, the notion of muteness detectors has been suggested [9, 10, 18]. However, it is not clear what system model could allow the implementation of muteness detectors, which is an inherent limitation of the approach.

These arguments suggest that failure detectors might not be the ultimate answer to the consensus problem. As an alternative to failure detectors, one could program directly at the level of the partially synchronous system model. However, this model provides too low level abstractions. It is indeed useful to provide higher level abstractions for expressing consensus algorithms. The goal of this paper is to show that another abstraction, namely *communication predicates*, provides a better abstraction than failure detectors for solving consensus. Specifically, the paper brings an answer to the question raised in [17], about quantifying the time it takes the environment to reach round synchronization after the system has stabilized.

The paper is structured as follows. Section 2 serves as a motivation to the introduction of communication predicates. Communication predicates are defined in Section 3. The implementation of communication predicates is presented in Section 4. Related work is discussed in Section 5, and Section 6 concludes the paper. Note that the paper is restricted to benign faults; Byzantine faults will be addressed in another paper.

## 2. Fault taxonomy

In this section we discuss the taxonomy of faults, with the goal to understand the limitation of failure detectors. The discussion will serve as the basis for the introduction of the notion of *communication predicates*.

### 2.1. Failure detectors and the paradox of the classical fault taxonomy

Let us come back to the second limitation of failure detectors (see Section 1), namely the gap between solving consensus with failure detectors in the crash-stop model and in the crash-recovery model. Our goal is to explain this gap, and so to understand the limited context in which failure detectors provide a good abstraction.

When looking at process failures, the classical fault taxonomy distinguishes, from the most benign to the most severe, (i) crash faults, (ii) send-omission faults, (iii) general-omission faults (which includes receive-omission faults), and (iv) malicious faults [22]. It can be observed that this taxonomy does not distinguish crash faults without recovery (the crash-stop model) and crash faults with recovery (the crash-recovery model). So, one would expect little difference when solving consensus in either of these two models. However, as already mentioned, this is not the case with failure detectors:

- In the crash-stop model, a standard solution to consensus is the rotating coordinator algorithm that requires the failure detector  $\diamond S$  and a majority of processes [5].
- Extending this solution to the crash-recovery model is not easy. It requires the definition of new failure detectors, and the algorithm becomes more complex [1]. This can be observed by comparing the two algorithms that are given in the appendix of [16].

This observation leads to the following question: What is the key issue, not captured by the classical fault taxonomy, that explains the gap between the crash-stop and crash-recovery consensus algorithm? The key issue is in the distinction between *permanent* faults and *transient* faults. Crash-stop is a model with permanent (crash) faults, while crash-recovery is a model with transient (crash) faults. A fault taxonomy that does not distinguish between permanent and transient fault is not able to explain the limitation of the failure detector model. In the next section we suggest another new fault taxonomy that makes the distinction between permanent and transient fault explicit.

### 2.2. Alternative fault taxonomy (for benign faults)

An alternative process fault taxonomy can be organized along two dimensions. The first dimension distinguishes between the already discussed *permanent* (P) and *transient* faults (T). The second dimension distinguishes faults that can hit any process in the system from faults that hit only a subset of the processes. We use the term *static* (S) for faults that can hit only a fixed subset of processes and *dynamic*<sup>3</sup> (D) for all other cases, i.e., faults that can hit all processes.

Combining this two dimensions leads to four classes of process faults:

- *SP*: at most  $f$  processes out of  $n$  are faulty ( $f < n$ ); a faulty process is permanently faulty.
- *ST*: at most  $f$  processes out of  $n$  are faulty ( $f < n$ ); faults are transient.

<sup>3</sup>This notion of static/dynamic faults was also used by [21].

- *DP*: all processes can be faulty; faults are permanent.
- *DT*: all processes can be faulty; faults are transient.

Among this classes, SP is clearly the most restrictive, whereas DT is the most general one. The crash-stop fault in the classical taxonomy corresponds the SP class. The send-omission and general-omission faults are transient faults. If we assume that only a subset of processes suffer from send-omission or general-omission faults, then send-omission and general-omission faults are classified as ST. Otherwise, these faults are classified as DT.

This alternative taxonomy is able to capture the distinction between the crash-stop model and the crash-recovery model: The crash-stop model corresponds to the class SP, whereas the crash-recovery model can be classified either as ST (if some processes never crash) or as DT. Failure detectors are well-suited to handle the SP fault class, but not to handle dynamic faults. Communication predicates will allow us to handle SP and DT faults in the same way.

### 2.3. Transmission faults

It is usual to distinguish between process faults and link faults. However, the distinction becomes irrelevant with DT faults. To see this, consider process  $p$  sending message  $m$  to process  $q$ . Process  $q$  might not receive  $m$  if (i)  $p$  suffers from a send-omission fault, (ii) the link loses  $m$ , or (iii)  $q$  suffers from a receive-omission fault. In case (i)  $p$  is the faulty component, in case (ii) the link  $l_{pq}$  is the faulty component, in case (iii)  $q$  is faulty. However, if the fault is transient, it may not occur later, for another message  $m'$  sent by  $p$  to  $q$ . For this reason, it makes no sense to put the responsibility of the fault on one of the components (process  $p$ , process  $q$ , or link  $l_{pq}$ ). This observation leads to consider only *transmission faults*:<sup>4</sup> a transmission fault is a fault that results in the non reception of some message  $m$ .

As we will see in Section 3, communication predicates are based on the notion of transmission faults. As such, communication predicates—contrary to failure detectors—are able to handle SP and DT fault classes uniformly.

## 3. Communication predicates and algorithms

### 3.1. Communication predicates

Communication predicates are defined in the context of a communication-closed round model. An algorithm for this model comprises, for each round  $r$  and process  $p \in \Pi$ , a sending function  $S_p^r$  and a transition function  $T_p^r$ . At beginning of a round  $r$ , every process sends a message to all

<sup>4</sup>The term is taken from [21], in which transmission faults are considered in the context of synchronous systems.

---

### Algorithm 1 The *OneThirdRule* algorithm [6].

---

```

1: Initialization:
2:    $x_p \leftarrow v_p$ 
3: Round  $r$ :
4:    $S_p^r$  :
5:     send  $\langle x_p \rangle$  to all processes
6:    $T_p^r$  :
7:     if  $|HO(p, r)| > 2n/3$  then
8:       if the values received, except at most  $\lfloor \frac{n}{3} \rfloor$ , are equal
          to  $\bar{x}$  then
9:          $x_p \leftarrow \bar{x}$ 
10:      else
11:         $x_p \leftarrow$  smallest  $x_q$  received
12:      if more than  $2n/3$  values received are equal to  $\bar{x}$  then
13:        DECIDE( $\bar{x}$ )

```

---

according to  $S_p^r(s_p)$ , where  $s_p$  is  $p$ 's state at the beginning of the round. At the end of a round  $r$ ,  $p$  makes a state transition according to  $T_p^r(\vec{\mu}, s_p)$ , where  $\vec{\mu}$  is the partial vector of all messages that have been received by  $p$  in round  $r$ .

We denote by  $HO(p, r)$  the support of  $\vec{\mu}$ , i.e., the set of processes (including itself) from which  $p$  receives a message at round  $r$ :  $HO(p, r)$  is the *heard of* set of  $p$  in round  $r$ . If  $q \notin HO(p, r)$ , then the message sent by  $q$  to  $p$  in round  $r$  was subject to a transmission failure. Communication predicates are expressed over the sets  $(HO(p, r))_{p \in \Pi, r > 0}$ . For example,

$$\exists r_0 > 0, \forall p, q \in \Pi : HO(p, r_0) = HO(q, r_0)$$

ensures the existence of some round  $r_0$  in which all processes hear of the same set of processes. Another example is a communication predicate that ensures that in every round  $r$  all processes hear of a majority of processes ( $n$  is the number of processes):

$$\forall r > 0, \forall p \in \Pi : |HO(p, r)| > n/2.$$

Let  $\mathcal{A} = \langle S_p^r, T_p^r \rangle$  be an HO algorithm. A problem is solved by a pair  $\langle \mathcal{A}, \mathcal{P} \rangle$ , where  $\mathcal{P}$  is a communication predicate. The consensus problem is specified by the following conditions:

- *Integrity*: Any decision value is the initial value of some process.
- *Agreement*: No two processes decide differently.
- *Termination*: All processes eventually decide.

The termination condition requires all processes to decide; a weaker condition is considered later. An example of a consensus algorithm is given by Algorithm 1.<sup>5</sup> The sending function is specified in lines 4–5. When the transition

<sup>5</sup>We have chosen this algorithm, rather than Paxos or another algorithm, for its simplicity. It allows us to keep the algorithmic part as simple as possible.

function (lines 6–13) is called, messages are available such that the predicate on the HO sets is guaranteed to hold. The consensus problem is solved by Algorithm 1 and the communication predicate  $\mathcal{P}_{otr}$ , given in Table 1 (next page).

**Theorem 1.** *The pair  $\langle \text{Algorithm 1}, \mathcal{P}_{otr} \rangle$  solves consensus.*

*Proof.* Algorithm 1 never violates the safety properties of consensus, namely integrity and agreement. For agreement, if some process decides  $v$  at line 13 of round  $r$ , then in any round  $r' \geq r$ , only  $v$  can be assigned to any  $x_p$ , and hence only  $v$  can be decided. Predicate  $\mathcal{P}_{otr}$  ensures the liveness property of consensus (termination). The first part of  $\mathcal{P}_{otr}$ , namely the existence of some round  $r_0$  in which all processes in  $\Pi$  have the set HO equal to some (large enough) set  $\Pi_0$ , ensures that at the end of round  $r_0$  all processes in  $\Pi$  adopt the same value for  $x_p$ . The second part of  $\mathcal{P}_{otr}$  forces every process  $p \in \Pi$  to make a decision at the end of round  $r_p$ .  $\square$

Note that  $\mathcal{P}_{otr}$  allows rounds in which no messages are received.

### 3.2. Restricted scope communication predicates

Section 2.3 has introduced the “transmission fault” abstraction, which covers various types of faults. One instantiation is to assume that transmission faults abstract link faults, send-omission faults and receive-omission faults, but not process crashes (i.e., processes do not crash). In this case the predicate  $\mathcal{P}_{otr}$ , which expresses a condition that must hold for all processes  $p \in \Pi$ , is perfectly adapted. This interpretation of transmission faults is also consistent with the termination condition for consensus that requires all processes to decide.

Let us now assume that transmission faults include in addition process crashes (without recovery). As already mentioned in [6], from the viewpoint of an HO algorithm this is still not a problem, since a crashed process does not send any messages and is thus indistinguishable from one that receives all messages but sends no messages. This holds no more if we implement the HO machine in a system where processes may exhibit any sort of benign faults. The problem can be addressed by restricting the scope of  $\mathcal{P}_{otr}$  to the subset  $\Pi_0$ , as defined by  $\mathcal{P}_{otr}^{restr}$ , see Table 1 (next page).

Predicate  $\mathcal{P}_{otr}^{restr}$  sets a requirement only for processes in  $\Pi_0$ , and so ensures termination only for processes in  $\Pi_0$ . If processes in  $\Pi_0$  do not crash, while processes in  $\Pi \setminus \Pi_0$  crash, then  $\langle \text{Algorithm 1}, \mathcal{P}_{otr}^{restr} \rangle$  allow all processes that do not crash to decide. So we have:

**Theorem 2.** *The pair  $\langle \text{Algorithm 1}, \mathcal{P}_{otr}^{restr} \rangle$  ensures the validity and agreement property of consensus. Moreover, all processes in  $\Pi_0$  eventually decide.*

*Proof.* Proof of Theorem 1, by replacing  $\Pi$  with  $\Pi_0$ .  $\square$

### 3.3. Crash-recovery model

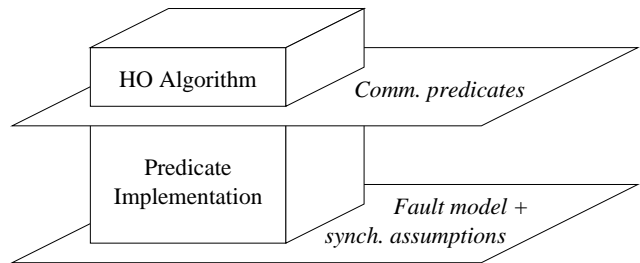
Algorithm 1 with predicate  $\mathcal{P}_{otr}^{restr}$  solves consensus with process crashes (crash-stop), link faults, send-omission, and receive-omission faults. In Section 2.1 we pointed out the gap between solving consensus with failure detectors in the crash-stop vs. the crash-recovery model. The gap disappears with the transmission fault abstraction and communication predicates.

Without any changes, Algorithm 1 can be used in the crash-recovery model. Handling of recoveries is done at a lower layer (cf. Section 4).

### 4. Achieving predicate $\mathcal{P}_{otr}^{restr}$ in good periods

We discuss now the implementation of the communication predicate  $\mathcal{P}_{otr}^{restr}$  introduced in Section 3. Figure 1 shows the algorithmic HO layer, the predicate implementation layer that we discuss now, and the interface between these two layers defined by communication predicates. This illustration shows also that the implementation of the predicates relies on assumptions about the underlying system (these assumptions define the fault and synchrony hypothesis). Note that “transmission faults” is an abstraction relevant to the upper layer: This abstraction does not appear at the lower layer.

In our implementation model, the system alternates between *good* and *bad* periods. In a good period the synchrony and fault assumptions hold; in a bad period the behavior of the system is arbitrary (but malicious behavior is excluded). The idea is here to compute the minimal duration of a good period that allows us to implement the communication predicates, i.e., the minimal duration of a good period that allow Algorithm 1 to solve consensus.



**Figure 1.** The two layers.

$$\mathcal{P}_{otr} :: \exists r_0 > 0, \exists \Pi_0, |\Pi_0| > 2n/3 : (\forall p \in \Pi : HO(p, r_0) = \Pi_0) \wedge (\forall p \in \Pi, \exists r_p > r_0 : |HO(p, r_p)| > 2n/3) \quad (1)$$

$$\mathcal{P}_{otr}^{restr} :: \exists r_0 > 0, \exists \Pi_0, |\Pi_0| > 2n/3 : (\forall p \in \Pi_0 : HO(p, r_0) = \Pi_0) \wedge (\forall p \in \Pi_0, \exists r_p > r_0 : HO(p, r_p) \supseteq \Pi_0) \quad (2)$$

**Table 1. Communication predicates**

#### 4.1. System model

Our system model is inspired by [12]; the differences are pointed out at the end of the section. We consider a message-passing system, and assume the existence of a fictitious global real-time clock that measures time with values from  $\mathbf{R}$  (see the remark on the next page for the reason for considering values from  $\mathbf{R}$  rather than integers). The clock is used only for analysis and is not accessible to the processes. Processes execute a sequence of atomic steps, which are either *send* steps or *receive* steps. As in [12], steps take no time (atomic steps), but time elapses between steps.<sup>6</sup> The network can take a *make-ready* step that is introduced to distinguish a message ready for reception from a message in transit: (i) Every process has two sets of messages called *network<sub>p</sub>* and *buffer<sub>p</sub>*; (ii) a *make-ready* step transfers a message from the first to the second set. Send steps, receive steps, and *make-ready* steps are defined to adequately model a real system:

- In a *send* step, a process  $p$  sends a message to either a single process or to all other processes and makes some local computation. More precisely, if  $p$  executes  $send_p(m)$  to all, then  $m$  is put into *network<sub>s</sub>*, for all  $s \in \Pi$ .
- In a *make-ready* step, the network transfers some messages from *network<sub>p</sub>* into *buffer<sub>p</sub>*. More precisely, if the network executes  $make\_ready_p(M)$  for some subset  $M \subseteq network_p$ , all messages  $m \in M$  are removed from *network<sub>p</sub>* and put into *buffer<sub>p</sub>*. Messages in *buffer<sub>p</sub>* are *ready for reception* by process  $p$ .
- In a *receive* step executed at time  $t$ , a process  $p$  may receive a *single* message that was in *buffer<sub>p</sub>* at time  $t$  and makes some local computation. *So  $n$  receive steps are needed to receive  $n$  messages.* If  $buffer_p = \emptyset$  at the time of a receive step, the empty message  $\lambda$  is received. A process  $p$  may specify any policy, according to which the message *buffer<sub>p</sub>* is selected for reception (e.g., “*message with the largest round number first*”).

We consider that the system alternates between good and bad periods. In a bad period, processes can crash and recover and suffer from send and receive omission; furthermore links can lose messages. We distinguish three types

<sup>6</sup>We model a step that “terminates” at time  $t$  as an atomic step that “occurs” at time  $t$ .

of good periods, from the strongest to the weakest. All these definitions refer to a subset  $\pi_0$  of  $\Pi$ . In all the three definitions, the following property  $\pi_0$ -*sync* holds in a good period for processes in  $\pi_0$ :

**$\pi_0$ -sync:** The subsystem  $\pi_0$  is synchronous, i.e., there is a known upper and lower bound on the process speed and a known upper bound on the communication delays among processes in  $\pi_0$ . Formally:

Let  $I$  be an open contiguous time interval and  $R$  a run. Processes and links are synchronous during  $I$  if there exist  $\Phi^+, \Phi^-, \Delta \in \mathbf{R}$  such that:

- In any contiguous sub-interval of  $I$  of length  $\Phi^+$ , every process in  $\pi_0$  takes at least one step.
- In any open contiguous sub-interval of  $I$  of length  $\Phi^-$ , every process in  $\pi_0$  takes at most one step.
- Consider two processes  $p, q \in \pi_0$ . If process  $p$  executes  $send_p(m)$  at time  $t \in I$ , then  $m \in buffer_q$  at time  $t + \Delta$ , provided that  $t + \Delta \in I$ .

The length of the good period is  $|I|$ . If  $I$  starts at time 0, we say  $I$  is an *initial good period*. We denote  $\Pi \setminus \pi_0$  by  $\bar{\pi}_0$ . We can now define the three types of good periods:

1.  **$\Pi$ -good period:** The property  $\pi_0$ -sync holds for  $\pi_0 = \Pi$ . All processes are up, none of these processes crashes (during the good period).
2. **“ $\bar{\pi}_0$ -down” good period:** The property  $\pi_0$ -sync holds for  $\pi_0 \subseteq \Pi$ . Processes in  $\pi_0$  do not crash. Processes in  $\bar{\pi}_0$  are down and do not recover (during the good period). Moreover, no messages from processes in  $\bar{\pi}_0$  are in transit during the good period.
3. **“ $\bar{\pi}_0$ -arbitrary” good period:** The property  $\pi_0$ -sync holds for  $\pi_0 \subseteq \Pi$ . There are no restrictions on the processes in  $\bar{\pi}_0$  and on the links to and from processes in  $\bar{\pi}_0$  (during the good period processes in  $\bar{\pi}_0$  can crash, recover, be asynchronous; links to and from processes in  $\bar{\pi}_0$  can lose messages, be asynchronous).

Case 2 includes case 1, and case 1 leads to the same implementation as case 2. Thus we distinguish below only between case 2 and case 3. For simplicity, we will use the following notation: We scale all values  $\Phi^+, \Phi^-, \Delta$ , and  $t$

with  $1/\Phi^-$  and use  $\phi = \Phi^+/\Phi^-$  as the normalized upper bound of the process speed,  $\delta = \Delta/\Phi^-$  as the normalized transmission delay, and  $\tau = t/\Phi^-$  as normalized time. Remember that  $\phi$  and  $\delta$  are “known” values, and note that *these values are unit-less*.<sup>7</sup>

**Remark:** For our modeling, we have chosen real-values clocks to represent time. Consider case 3 above, assuming integer clock values instead. By the definition of  $\Phi^+$ , the slowest process in  $\pi_0$  takes at least one step in any interval  $\Phi^+$ . However, with integer clock values, any process can take at most  $\Phi^+$  steps in an interval  $\Phi^+$ , independent how small  $\Phi^-$  is chosen. So, in case 3, processes in  $\overline{\pi_0}$  cannot be arbitrarily fast with respect to processes in  $\pi_0$ . In other words, with integer clock values, processes in  $\overline{\pi_0}$  have some synchrony relation with respect to processes in  $\pi_0$ , which we wanted to exclude under case 3.

#### Differences between our system model and DLS [12]:

In [12] the clocks take integer values. We have explained the reason to consider clocks with real-time values. In [12] a send step allows a process to send a message only to a single destination. Our send primitive allows messages to be broadcast, a facility provided, e.g., by UDP-multicast. In [12], a receive step allows a process to receive several messages. Our receive primitive allows reception of a message from one single process only, which reflects the feature, e.g., of UDP. The reception of messages one by one led us to introduce the make-ready step. Two different synchrony assumptions are considered in [12]: (i) The synchrony bounds are known but hold only eventually; (ii) the synchrony bounds are not known, but hold from the beginning. We considered option (i), which is needed to compute the minimal length of a good period (in the context of the implementation of the communication predicates). In the context of option (i), [12] assumes that the good period holds eventually forever and that the synchrony assumption holds on the whole system. We consider the system alternating between good and bad periods, and synchrony assumptions that hold only on a subset  $\pi_0$ . We also assume the more general crash-recovery model, while [12] considers the crash-stop model. On the other hand, contrary to our fault model, [12] considers also Byzantine faults.

## 4.2. Implementation of $\mathcal{P}_{otr}^{restr}$

We give now algorithms for implementing the predicate  $\mathcal{P}_{otr}^{restr}$  in  $\overline{\pi_0}$ -down and  $\overline{\pi_0}$ -arbitrary good periods. It turns out that both definitions of a good period lead naturally to the implementation of a predicate that is stronger than

<sup>7</sup>For obtaining real-time values, the results in this section have thus to be multiplied by  $\Phi^-$ .

$\mathcal{P}_{otr}^{restr}$ . We define:

$$\begin{aligned} \mathcal{P}_{su}(\Pi_0, r_1, r_2) &:: \forall p \in \Pi_0, \forall r \in [r_1, r_2] : HO(p, r) = \Pi_0 \\ \mathcal{P}_k(\Pi_0, r_1, r_2) &:: \forall p \in \Pi_0, \forall r \in [r_1, r_2] : HO(p, r) \supseteq \Pi_0 \\ \mathcal{P}_{otr}^2(\Pi_0) &:: \exists r_0 > 0 : \mathcal{P}_{su}(\Pi_0, r_0, r_0) \\ &\quad \wedge \mathcal{P}_k(\Pi_0, r_0+1, r_0+1) \\ \mathcal{P}_{otr}^{1/1}(\Pi_0) &:: \exists r_0 > 0, \exists r_1 > r_0 : \mathcal{P}_{su}(\Pi_0, r_0, r_0) \\ &\quad \wedge \mathcal{P}_k(\Pi_0, r_1, r_1) \end{aligned}$$

Predicate  $\mathcal{P}_{su}(\Pi_0, r_1, r_2)$  ensures that rounds from  $r_1$  to  $r_2$  are so called “space uniform” for the processes in  $\Pi_0$ . Predicate  $\mathcal{P}_k(\Pi_0, r_1, r_2)$  ensures a weaker property ( $k$  stands for *kernel*). Predicate  $\mathcal{P}_{otr}^2(\Pi_0)$  ensures two consecutive rounds such that the first satisfies  $\mathcal{P}_{su}(\Pi_0, -, -)$  and the second  $\mathcal{P}_k(\Pi_0, -, -)$ . Predicate  $\mathcal{P}_{otr}^{1/1}(\Pi_0)$  ensures the same property for two rounds that do not need to be consecutive. We clearly have:

$$\begin{aligned} (\exists \Pi_0, \text{ s.t. } |\Pi_0| > 2n/3 : \mathcal{P}_{otr}^2(\Pi_0)) &\Rightarrow \mathcal{P}_{otr}^{restr} \\ (\exists \Pi_0, \text{ s.t. } |\Pi_0| > 2n/3 : \mathcal{P}_{otr}^{1/1}(\Pi_0)) &\Rightarrow \mathcal{P}_{otr}^{restr}. \end{aligned}$$

We give below algorithms for  $\mathcal{P}_{su}(-, -, -)$  and  $\mathcal{P}_k(-, -, -)$ , for both definitions of good periods. We also analyze the timing property of the algorithms under the following two scenarios:

1. Assume that a good period starts at an arbitrary time  $t_G$  resp.  $\tau_G = t_G/\Phi^-$ . We compute, in the worst case, the minimal length of the good period needed to implement the communication predicates. We call this value *minimal length of a good period*.
2. We do the same, assuming that a good period starts from the beginning, i.e.,  $\tau_G = 0$ . We call this value *minimal length of an initial good period*.

Intuitively, scenario 2 allows us to compute the time to solve consensus in the fault-free case, which is often called a “nice” run. Scenario 1 allows us a timing analysis of consensus in “not nice” runs.

### 4.2.1. Ensuring $\mathcal{P}_{otr}^{restr}$ in a “ $\overline{\pi_0}$ -down” good period

Let us consider a “ $\overline{\pi_0}$ -down” good period that is “long enough”, with  $\pi_0$  arbitrary. Algorithm 2 implements  $\mathcal{P}_{su}(\pi_0, -, -)$ . The function  $S_p^{r,p}$  at line 7 returns the message to be sent; the send occurs at line 8. Variable  $i_p$  (line 9, 11) counts the number of receive steps. If  $p$  executes  $x$  steps, at least  $x$  and at most  $x\phi$  (normalized) time has elapsed (see Section 4.1). Process  $p$  executes at most  $\lceil 2\delta + (n+2)\phi \rceil$  receive steps, see line 12 (message reception takes place at line 14; non-empty messages are added to the set  $msgsRcv_p$ , see line 16). Process  $p$  executes receive steps (1) until  $\lceil 2\delta + (n+2)\phi \rceil$  receive steps have been executed, or (2) if  $p$  receives a message from a round  $r'$  larger

---

**Algorithm 2** Ensuring  $\mathcal{P}_{su}(\pi_0, -, -)$  with a “ $\bar{\pi}_0$ -down” good period

---

```

1: Reception policy: Highest round number first
2:  $msgsRcv_p \leftarrow \emptyset$  {set of messages received}
3:  $r_p \leftarrow 1$  {round number}
4:  $next\_r_p \leftarrow 1$  {next round number}
5:  $s_p \leftarrow init_p$  {state of the consensus algorithm}
6: while true do
7:    $msg \leftarrow S_p^{r_p}(s_p)$ 
8:   send  $\langle msg, r_p \rangle$  to all
9:    $i_p \leftarrow 0$ 
10:  while next_r_p = r_p do
11:     $i_p \leftarrow i_p + 1$ 
12:    if  $i_p \geq 2\delta + (n + 2)\phi$  then
13:       $next\_r_p \leftarrow r_p + 1$ ;
14:      receive a message
15:      if message is  $\langle msg, r' \rangle$  from  $q$  then
16:         $msgsRcv_p \leftarrow msgsRcv_p \cup \{\langle msg, r', q \rangle\}$ 
17:        if  $r' > r_p$  then
18:           $next\_r_p \leftarrow r'$ 
19:         $R \leftarrow \{\langle msg', q' \rangle \mid \langle msg', r_p, q' \rangle \in msgsRcv_p\}$ 
20:         $s_p \leftarrow T_p^{r_p}(R, s_p)$ 
21:        forall  $r'$  in  $[r_p + 1, next\_r_p - 1]$  do  $s_p \leftarrow T_p^{r'}(\emptyset, s_p)$ 
22:         $r_p \leftarrow next\_r_p$ 

```

---

than  $r_p$ . In both cases the state transition function  $T_p^{r_p}$  is executed with the set  $R$  of messages received in round  $r_p$  (line 20). Then the state transition function  $T_p^{r_p}$  is executed for all rounds  $r_p + 1$  to  $next\_r_p - 1$  with an empty set of messages.<sup>8</sup>

In order to cope with recoveries after crashes, variables  $r_p$  and  $s_p$  are stored on stable storage. In case of a recovery, the algorithm starts on line 6, with  $msgsRcv_p$  and  $next\_r_p$  reinitialized. Reading variables on stable storage is inefficient. The implementation can be made more efficient by keeping a copy of the variables in main memory: a read operation reads the *in memory* copy, a write operation updates the *in memory* and the *stable storage* copies. Upon recovery, the *in memory* copy is reset with the value of the stable copy.<sup>9</sup>

Algorithm 2 is not optimized regarding space, i.e., the set  $msgsRcv_p$  grows forever. Obviously, messages for round smaller than  $r_p$  can safely be discarded. To keep the presentation short, we did not include this simple optimization.

It should be noted that Algorithm 2 relies exclusively on messages sent by the upper algorithmic layer: Algorithm 2 does not send any additional message.

We prove Algorithm 2 in two steps. First we prove that there exists  $r > 0$  such that, for any  $x > 0$ , Algorithm 2 ensures  $\mathcal{P}_{su}(\pi_0, r, r + x - 1)$ , assuming a “long enough” good

<sup>8</sup>This is required only if  $T_p^{r_p}(\emptyset, s_p) \neq s_p$ . Calling the sending function  $S_p^{r_p}$  is not needed, since the function does not change the state  $s_p$ .

<sup>9</sup>We could express this formally as a variant of Algorithm 2, but the space constraints prevent us from doing this.

period. Then we compute the minimal duration of a good period to ensure  $\mathcal{P}_{otr}^2(\pi_0)$ , and the minimal duration of two good periods to ensure  $\mathcal{P}_{otr}^{1/1}(\pi_0)$ . Note that by the definition of a  $\bar{\pi}_0$ -down good period, all processes in  $\bar{\pi}_0$  are down in a good period, and no messages from these processes are in transit in the good period. In other words, processes in  $\bar{\pi}_0$  can simply be ignored.

**Theorem 3.** *With Algorithm 2, the minimal length of a good period to achieve  $\mathcal{P}_{su}(\pi_0, \rho_0, \rho_0 + x - 1)$  is:*

$$(x + 1)(2\delta + (n + 2)\phi + 1)\phi + \delta + \phi.$$

The proof, also for all other theorems of this paper, can be found in [16]. The following Corollary follows directly from Theorem 3 with  $x=1$  and  $x=2$ , and the fact that  $\mathcal{P}_{su}(-, -, -) \Rightarrow \mathcal{P}_k(-, -, -)$ :

**Corollary 4.** *For implementing  $\mathcal{P}_{otr}^2(\pi_0)$  with Algorithm 2, we need one “ $\bar{\pi}_0$ -down” good period of length*

$$(6\delta + 3n\phi + 6\phi + 3)\phi + \delta + \phi.$$

*For implementing  $\mathcal{P}_{otr}^{1/1}(\pi_0)$  with Algorithm 2, we need two “ $\bar{\pi}_0$ -down” good periods of length*

$$(4\delta + 2n\phi + 4\phi + 2)\phi + \delta + \phi.$$

Corollary B.1 shows an interesting trade-off in terms of the length of a good period. The next theorem gives us the minimal length of an *initial* good period:

**Theorem 5.** *With Algorithm 2, the minimal length of an initial good period to achieve  $\mathcal{P}_{su}(\pi_0, 1, x)$  is:*

$$x(2\delta + (n + 2)\phi + 1)\phi.$$

As already pointed out, Theorem 5 is related to so-called “nice” runs, while Theorem 3 is related to “not nice” runs. This second case has not been addressed in the literature with a time analysis as done here (see Section 5). The results show a factor of approximately 3/2 between the two cases for the relevant value  $x = 2$ .

#### 4.2.2. Ensuring $\mathcal{P}_{otr}^{restr}$ in a “ $\bar{\pi}_0$ -arbitrary” good period

In this section we consider a  $\bar{\pi}_0$ -arbitrary good period. Compared with the previous section, the problem is more complex. We proceed in two steps. First we show how to implement the predicate  $\mathcal{P}_k(\pi_0, -, -)$ . Second, we show how to obtain the predicate  $\mathcal{P}_{su}(\pi_0, -, -)$  from  $\mathcal{P}_k(\pi_0, -, -)$ . Note that we introduce here a parameter  $f$  defined such that  $|\pi_0| = n - f$ . The implementation of  $\mathcal{P}_k(\pi_0, -, -)$  requires  $f < n/2$ .

### a) Implementing $\mathcal{P}_k(\pi_0, -, -)$

The algorithm for implementing  $\mathcal{P}_k(\pi_0, -, -)$  is given as Algorithm 3. It uses two different types of messages, INIT messages and ROUND messages. Processes express the intention to enter a new round  $\rho$  with an  $\langle \text{INIT}, \rho, - \rangle$  message. If a process receives at least  $f + 1$  INIT messages for some round  $\rho$ , it starts round  $\rho$  and sends a  $\langle \text{ROUND}, \rho, - \rangle$  message. A process in round  $\rho$  that receives a ROUND message for a higher round  $\rho'$  enters immediately round  $\rho'$ . This ensures fast synchronization at the beginning of a good period, and is one of the major differences of this algorithm compared to Byzantine clock synchronization algorithms.

The reception policy for Algorithm 3 (line 1) is a little bit more complicated than for Algorithm 2. Algorithm 3 has to ensure that a fast process with a large round number  $r'$  is not able to prevent messages from other processes with lower round numbers  $r < r'$  from being received. The reception policy is as follows: At the  $i$ th receive step, the message with the highest round number from process  $p_{i \bmod n}$  is selected for reception. If no such message exists, an arbitrary message is selected.

As for the previous algorithm, variables  $r_p$  and  $s_p$  are assumed to be on stable storage (possibly with a copy in volatile memory) and the algorithm starts after a recovery in line 6, with  $msgsRcv_p$  and  $next_r_p$  reinitialized.

We prove the following results:

**Theorem 6.** *With Algorithm 3 and  $f < n/2$ , the minimal length of a good period to achieve  $\mathcal{P}_k(\pi_0, \rho_0, \rho_0 + x - 1)$  is*

$$(x+2)[(2\delta+2n\phi+\phi)\phi+\delta+n\phi+2\phi]+(2\delta+2n\phi+\phi)\phi$$

**Theorem 7.** *With Algorithm 3 and  $f < n/2$ , the minimal length of an initial good period to implement  $\mathcal{P}_k(\pi_0, 1, x)$  is:*

$$(x-1)[(2\delta+2n\phi+\phi)\phi+\delta+n\phi+2\phi]+(2\delta+2n\phi+\phi)\phi+\phi.$$

### b) Implementing $\mathcal{P}_{su}(\pi_0, -, -)$ from $\mathcal{P}_k(\pi_0, -, -)$

We show now that  $f+1$  rounds that satisfy  $\mathcal{P}_k(\pi_0, -, -)$ , with  $|\pi_0| = n - f$ , allow us to construct one *macro-round* that satisfies  $\mathcal{P}_{su}(\pi_0, -, -)$ . The “translation” is given by Algorithm 4, which is derived from a similar translation in [6]. Let  $r_1, \dots, r_{f+1}$  denote the sequence of the  $f + 1$  rounds that form a macro-round  $\mathcal{R}$ . In round  $r_1$ , every process  $p$  sends its message for macro-round  $\mathcal{R}$  (line 7). In all subsequent rounds  $r_2, \dots, r_{f+1}$  messages previously received are relayed (line 7). In round  $r_{f+1}$  (i.e.,  $r \equiv 0 \pmod{f+1}$ , see line 9), the set of messages of macro-round  $\mathcal{R}$  to be received by  $p$  are computed (lines 13 and 14).

---

### Algorithm 3 Ensuring $\mathcal{P}_k(\pi_0, , )$ with a “ $\overline{\pi_0}$ -arbitrary” good period

---

```

1: Reception policy: The highest round message from each
   process in a round robin fashion
2:  $msgsRcv_p \leftarrow \emptyset$ 
3:  $r_p \leftarrow 1$ 
4:  $next_r_p \leftarrow 1$ 
5:  $s_p \leftarrow init_p$ 
6: while true do
7:    $msg \leftarrow S_p^{r_p}(s_p)$ 
8:   send  $\langle \text{ROUND}, r_p, msg \rangle$  to all
9:    $i \leftarrow 0$ 
10:  while  $next_r_p = r_p$  do
11:    receive a message
12:    if message is  $\langle \text{ROUND}, msg, r' \rangle$  or  $\langle \text{INIT}, msg, r' + 1 \rangle$ 
      from  $q$  then
13:       $msgsRcv_p \leftarrow msgsRcv_p \cup \{ \langle msg, r', q \rangle \}$ 
14:      if  $r' > r_p$  then
15:         $next_r_p \leftarrow r'$ 
16:      if received  $f+1$  messages  $\langle \text{INIT}, r_p + 1, - \rangle$  from distinct
      processes then
17:         $next_r_p \leftarrow \max\{r_p + 1, next_r_p\}$ 
18:         $i \leftarrow i + 1$ 
19:        if  $i \geq 2\delta + (2n + 1)\phi$  then
20:          send  $\langle \text{INIT}, r_p + 1, msg \rangle$  to all
21:         $R \leftarrow \{ \langle msg', q' \rangle \mid \langle msg', q' \rangle \in msgsRcv_p \}$ 
22:         $s_p \leftarrow T_p^{r_p}(R, s_p)$ 
23:        forall  $r'$  in  $[r_p+1, next_r_p-1]$  do  $s_p \leftarrow T_p^{r'}(\emptyset, s_p)$ 
24:         $r_p \leftarrow next_r_p$ 

```

---

### c) Putting it all together

When combining Algorithm 3 and Algorithm 4, the function  $S_p^{r_p}(\cdot)$  called in Algorithm 3 refers to line 7 of Algorithm 4. Similarly, the function  $T_p^{r_p}(\cdot)$  called in Algorithm 3 refers to the lines 9 to 17 of Algorithm 4. The functions  $S_p^r$  and  $T_p^r$  in Algorithm 4 refer to the sending phase and state transition phase of Algorithm 1.

We compute now the minimal duration of a good period to ensure  $\mathcal{P}_{otr}^2(\pi_0)$  (considering instead  $\mathcal{P}_{otr}^{1/1}(\pi_0)$  is not a valuable alternative here):

1. We need first  $f + 1$  rounds that satisfy  $\mathcal{P}_k(-, -, -)$  to implement one macro-round that satisfies  $\mathcal{P}_{su}(-, -, -)$  (Algorithm 4).
2. Then we need one round that satisfies  $\mathcal{P}_k(-, -, -)$ .

For 1, the worst case happens when the good period starts immediately after the beginning of a macro-round. In this worst case,  $\mathcal{P}_{su}(-, -, -)$  requires two macro-rounds. Since one macro-round consists of  $f+1$  rounds, in the worst case we need  $2(f+1)$  rounds. Item 2 adds one round. So we end up with a minimal duration of  $2f + 3$  rounds. Applying Theorem 6, we get the minimal length of a good period:

$$(2f+5)[(2\delta+2n\phi+\phi)\phi+\delta+n\phi+2\phi]+(2\delta+2n\phi+\phi)\phi.$$



---

**Algorithm 4** Ensuring  $\mathcal{P}_{su}(\pi_0, -, -)$  with  $\mathcal{P}_k(\pi_0, -, -)$  (adapted from [6]).

---

```

1: Variables:
2:  $Listen_p$ , initially  $\Pi$  {set of processes}
3:  $NewHO_p$  {set of processes}
4:  $Known_p$ , initially  $\{\langle S_p^{\mathcal{R}p}(s_p), p \rangle\}$ 
{set of (message, process)}

5: Round  $r$ :
6:  $S_p^r$  :
7: send  $\langle Known_p \rangle$  to all processes

8:  $T_p^r$  :
9:  $Listen_p \leftarrow Listen_p \cap \{q \mid \langle Known_q \rangle \text{ received}\}$ 
10: if  $r \not\equiv 0 \pmod{f+1}$  then
11:    $Known_p \leftarrow Known_p \cup \bigcup_{q \in Listen_p} Known_q$ 
12: else
13:    $NewHO_p \leftarrow \{s \mid \langle -, s \rangle \in Known_q \text{ for } n - f$ 
processes  $q \in Listen_p\}$ 
14:    $R \leftarrow \{\langle msg, s \rangle \mid s \in NewHO_p\}$ 
15:    $s_p \leftarrow T_p^{\mathcal{R}p}(R, s_p)$ 
16:    $Listen_p \leftarrow \Pi$ 
17:    $Known_p \leftarrow \{\langle S_p^{\mathcal{R}p}(s_p), p \rangle\}$ 

```

---

## 5 Related work

The paper addresses several issues that appear in the literature. We now point out the key differences.

The HO model was proposed in [6]. The paper establishes relationship among several communication predicates and identifies the weakest predicate, among the class of predicates with non-empty kernel rounds, for solving consensus. The paper also expresses well-known consensus algorithms (or variants) and new ones in the HO model, with the goal of showing the expressiveness of the model. The implementation of communication predicates is not addressed in [6], nor is the ability of the model to handle uniformly crash-stop and crash-recovery models, and the reason for that. In [7] the HO model is used to express a new consensus algorithm.

The HO model generalizes the round model of [12], but does not reintroduce failure detectors as done in [14] and in [17]. The implementation in [12], contrary to ours, explicitly refers to some “common notion of time” and relies on a distributed clock synchronization algorithm.

It has been sometimes claimed that the partial synchrony model has been superseded by the failure detector model [5]. In our opinion this claim is only partially correct. The models that extend the failure detector model, e.g., [14, 17], all inherit from the limitations of failure detectors pointed out in Section 1.

The issue of performance of consensus following asynchronous periods is considered in [11, 17]. In [17] the focus is on number of rounds rather than time; [11] considers

time. Moreover, in [17] the authors write that being able to quantify the time it takes the environment to reach round synchronization after the system has stabilized is an interesting subject for further studies. This question is answered here. In [11] and [17] the synchronous period is defined only by properties of links: Processes are always considered to be synchronous. This is in contrast to our definition of  $\bar{\pi}_0$ -arbitrary good period, where only a subset of processes are assumed to be synchronous. This definition opens the door to the analysis of the duration of good periods with Byzantine processes. Our algorithm shares some similarities with the Byzantine clock synchronization of [23]. However, the algorithm in [23] assumes reliable links; adapting the algorithm to message loss, we end up with the algorithm of [12].

The notion of good and bad period appears in [8], but the issue of the length of a good period for solving consensus is not addressed. Restricting the scope of synchrony, as we do in good periods, has been considered in other settings, e.g., [15] and [2, 3]. However, in all these papers the issue of synchrony is implicitly restricted to links (i.e., process synchrony is not addressed). This is not the case in our definition of  $\bar{\pi}_0$ -arbitrary good period.

The Paxos algorithm [19] does not assume reliable links and, because of this, works under the crash-recovery model with stable storage. However, the condition for liveness is not expressed by a clean abstraction as done by communication predicates in the HO model (a consensus algorithm *à la Paxos* in the HO model can be found in [6]). The same comment applies to [11], where the system must stabilize before consensus is reached. System stabilization is not required with  $\bar{\pi}_0$ -arbitrary good periods: the HO model provides a clean separation of concerns between the HO algorithmic layer and the predicate implementation layer, which allows a finer definition of good periods, and so a finer timing analysis. As pointed out in Section 3, we have chosen here an algorithm that is simpler than Paxos to illustrate as simply as possible the approach based on communication predicates.

The notion of transmission faults was suggested in [21], however only in the context of synchronous systems. Varying the quorums for “init” and round messages—in the context of  $\bar{\pi}_0$ -arbitrary good periods—was to our knowledge done first in [20, 24], but for other fault scenarios.

## 6 Conclusion

Abstractions are essential when solving difficult problems. Failure detectors provide a nice abstraction for solving the difficult consensus problem; this explains why they have been widely adopted. However, transient and dynamic faults show the limitations of the failure detector approach: For example, solving consensus in the crash-stop model and

in the crash-recovery model leads to significantly different solutions. The HO model provides a different abstraction, namely the “communication predicates”, which allow us to handle uniformly static, dynamic, transient, and permanent faults and so overcome the limitations of failure detectors. Moreover, the HO model allows a nice and concise expression of consensus algorithms.

Similarly to failure detectors, solving consensus in the HO model leads to distinguish two layers: The “algorithmic” layer and the “abstraction” layer (the layer at which the abstraction is implemented). In the case of failure detectors, the abstraction layer must ensure the properties of the failure detectors, based on assumptions of the underlying system. The same holds for communication predicates. However, while communication predicates are based on the very general notion of transmission faults, failure detector assume the limited notion of process crash faults. The communication predicate layer defines a larger “playground” than the failure detector playground, in which more issues can be addressed. Specifically, the communication predicate approach has allowed us to bring an answer the question raised in [17], about quantifying the time it takes to reach round synchronization after the system has stabilized.

**Acknowledgments** We like to thank Bernadette Charron-Bost, Nicolas Schiper, Martin Biely, Josef Widder, Nuno Santos, Sergio Mena, and the anonymous reviewers for their valuable comments that helped us to improve the paper.

## References

- [1] M. Aguilera, W. Chen, and S. Toueg. Failure detection and consensus in the crash-recovery model. *Distributed Computing*, 13(2):99–125, 2000.
- [2] M. K. Aguilera, C. Delporte-Gallet, H. Fauconnier, and S. Toueg. On implementing Omega with weak reliability and synchrony assumptions. In *Proc. PODC’03*. ACM Press, 2003.
- [3] M. K. Aguilera, C. Delporte-Gallet, H. Fauconnier, and S. Toueg. Communication-efficient leader election and consensus with limited link synchrony. In *Proc. PODC’04*, pages 328–337. ACM Press, 2004.
- [4] R. Boichat, P. Dutta, S. Frolund, and R. Guerraoui. Reconstructing Paxos. *ACM SIGACT News*, 34(1):47–67, 2003.
- [5] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of ACM*, 43(2):225–267, 1996.
- [6] B. Charron-Bost and A. Schiper. The “Heard-Of” model: Unifying all benign faults. Technical Report TR, EPFL, June 2006.
- [7] B. Charron-Bost and A. Schiper. Improving fast paxos: being optimistic with no overhead. In *Pacific Rim Dependable Computing, Proceedings*, 2006.
- [8] F. Cristian and C. Fetzer. The timed asynchronous distributed system model. *IEEE Transactions on Parallel and Distributed Systems*, 10(6):642–657, 1999.
- [9] A. Doudou and A. Schiper. Muteness Failure Detectors for Consensus with Byzantine Processes. TR 97/230, EPFL, Dept d’Informatique, October 1997.
- [10] A. Doudou and A. Schiper. Muteness detectors for consensus with byzantine processes (Brief Announcement). In *Proc. PODC’98*, Puerto Vallarta, Mexico, July 1998.
- [11] P. Dutta, R. Guerraoui, and L. Lamport. How fast can eventual synchrony lead to consensus? In *Proc. DSN’05*, pages 22–27, Los Alamitos, CA, USA, 2005.
- [12] C. Dwork, N. Lynch, and L. Stockmeyer. Consensus in the presence of partial synchrony. *Journal of ACM*, 35(2):288–323, April 1988.
- [13] M. Fischer, N. Lynch, and M. Paterson. Impossibility of Distributed Consensus with One Faulty Process. *Journal of ACM*, 32:374–382, April 1985.
- [14] E. Gafni. Round-by-round fault detectors: Unifying synchrony and asynchrony. In *Proc of the 17th ACM Symp. Principles of Distributed Computing (PODC)*, pages 143–152, Puerto Vallarta, Mexico, June-July 1998.
- [15] R. Guerraoui and A. Schiper. “gamma-accurate” failure detectors. In *Proceedings of the 10th International Workshop on Distributed Algorithms (WDAG’96)*, pages 269–286, London, UK, 1996. Springer-Verlag.
- [16] M. Hutle and A. Schiper. Communication predicates: A high-level abstraction for coping with transient and dynamic faults. Technical Report LSR-REPORT-2006-006, EPFL, 2006. <http://infoscience.epfl.ch/search.py?recid=97290>.
- [17] I. Keidar and A. Shraer. Timeliness, failure-detectors, and consensus performance. In *Proc. PODC’06*, pages 169–178, New York, NY, USA, 2006. ACM Press.
- [18] K. P. Kihlstrom, L. E. Moser, and P. M. Melliar-Smith. Solving consensus in a byzantine environment using an unreliable fault detector. In *Proceedings of the International Conference on Principles of Distributed Systems (OPODIS)*, pages 61–75, Chantilly, France, Dec. 1997.
- [19] L. Lamport. The Part-Time Parliament. *ACM Trans. on Computer Systems*, 16(2):133–169, May 1998.
- [20] G. Le Lann and U. Schmid. How to implement a timer-free perfect failure detector in partially synchronous systems. Technical Report 183/1-127, Technische Universität Wien, Department of Automation, Jan. 2003.
- [21] N. Santoro and P. Widmayer. Time is not a healer. In *Proceedings of the 6th Symposium on Theor. Aspects of Computer Science*, pages 304–313, Paderborn, Germany, 1989.
- [22] F. B. Schneider. What Good are Models and What Models are Good. In S. Mullender, editor, *Distributed Systems*, pages 169–197. ACM Press, 1993.
- [23] T. K. Srikanth and S. Toueg. Optimal clock synchronization. *Journal of the ACM*, 34(3):626–645, 1987.
- [24] J. Widder. Booting clock synchronization in partially synchronous systems. In *Proceedings of the 17th International Conference on Distributed Computing (DISC’03)*, pages 121–135, 2003.

# Appendix

## A Solving Consensus with failure detectors: crash-stop vs. crash-recovery

We give the consensus algorithms for the crash-stop (Chandra and Toueg [5]) and the crash-recovery model (Aguilera et al. [1]), as discussed in Section 2.

### A.1 Crash-stop algorithm

---

**Algorithm 5** Chandra and Toueg [5]: Solving consensus in the crash-stop model with  $\diamond\mathcal{S}$

---

```
1: procedure propose( $v_p$ )
2:    $estimate_p \leftarrow v_p$  { $p$ 's estimate of the decision value}
3:    $state_p \leftarrow undecided$ 
4:    $r_p \leftarrow 0$  { $p$ 's current round number}
5:    $ts_p \leftarrow 0$  {last the round in which  $p$  updated  $estimate_p$ , initially 0}

6: while  $state_p = undecided$  do
7:    $r_p \leftarrow r_p + 1$ 
8:    $c_p \leftarrow (r_p \bmod n) + 1$  { $c_p$  is the current coordinator}

9:   Phase 1:
10:  send ( $p, r_p, estimate_p, ts_p$ ) to  $c_p$ 

11:  Phase 2:
12:  if  $p = c_p$  then
13:    wait until [for  $\lceil \frac{(n+1)}{2} \rceil$  processes  $q$ : received ( $q, r_p, estimate_q, ts_q$ ) from  $q$ ]
14:     $msgs_p[r_p] \leftarrow \{(q, r_p, estimate_q, ts_q) \mid p \text{ received } (q, r_p, estimate_q, ts_q) \text{ from } q\}$ 
15:     $t \leftarrow$  largest  $ts_q$  such that  $(q, r_p, estimate_q, ts_q) \in msgs_p[r_p]$ 
16:     $estimate_p \leftarrow$  select one  $estimate_q$  such that  $(q, r_p, estimate_q, t) \in msgs_p[r_p]$ 
17:    send ( $p, r_p, estimate_p$ ) to all

18:  Phase 3:
19:  wait until [received ( $c_p, r_p, estimate_{c_p}$ ) from  $c_p$  or  $c_p \in \mathcal{D}_p$ ] {query the failure detector}
20:  if received ( $c_p, r_p, estimate_{c_p}$ ) from  $c_p$  then
21:     $estimate_p \leftarrow estimate_{c_p}$ 
22:     $ts_p \leftarrow r_p$ 
23:    send ( $p, r_p, ack$ ) to  $c_p$ 
24:  else
25:    send ( $p, r_p, nack$ ) to  $c_p$  { $p$  suspects that  $c_p$  crashed}

26:  Phase 4:
27:  if  $p = c_p$  then
28:    wait until [for  $\lceil \frac{(n+1)}{2} \rceil$  processes  $q$ : received ( $q, r_p, ack$ ) or ( $q, r_p, nack$ )]
29:    if for  $\lceil \frac{(n+1)}{2} \rceil$  processes  $q$ : received ( $q, r_p, ack$ ) then
30:       $R$ -broadcast( $p, r_p, estimate_p, decide$ ) to all

31: when  $R$ -deliver( $q, r_q, estimate_q, decide$ ) for the first time
32: if  $state_p = undecided$  then
33:    $state_p \leftarrow decided$ 
34:    $decide(estimate_q)$ 
```

---

## A.2 Crash-recovery algorithm

---

**Algorithm 6** *Aguilera et al. [1]*: Solving consensus in the crash-recovery model with  $\diamond \mathcal{S}_u$

---

```

1: Initialization:
2:   for all  $q \in \Pi \setminus \{p\}$  do  $xmitmsg[q] \leftarrow \perp$ 

3: To s-send  $m$  to  $q$ :
4:   if  $q \neq p$  then  $xmitmsg[q] \leftarrow m$ ; send  $m$  to  $q$  else simulate receive  $m$  from  $p$ 

5: Task retransmit:
6:   repeat forever
7:     for all  $q \in \Pi \setminus \{p\}$  do if  $xmitmsg[q] \neq \perp$  then send  $xmitmsg[q]$  to  $q$ 

8: upon propose( $v_p$ ):
9:    $(r_p, estimate_p, ts_p) \leftarrow (1, v_p, 0)$ 
10:  fork task  $\{4phases, retransmit\}$ 

11: Task 4phases:
12:  store  $\{r_p\}$ ;  $c_p \leftarrow (r_p \bmod n + 1)$ ; fork task  $\{skip\_round, participant\}$ 
13:  if  $p = c_p$  then fork task coordinator

14: Task coordinator:
15:  {Phase NEWROUND}
16:  if  $ts_p \neq r_p$  then
17:    s-send  $(r_p, NEWROUND)$  to all
18:    wait until [ received  $(r_p, estimate_q, ts_q,$ 
19:      ESTIMATE) from  $\lceil (n+1)/2 \rceil$  processes ]
20:     $t \leftarrow$  largest  $ts_q$  such that  $p$  received
21:     $(r_p, estimate_q, ts_q, ESTIMATE)$ 
22:     $estimate_p \leftarrow$  select one  $estimate_q$  such that
23:     $p$  received  $(r_p, estimate_q, t, ESTIMATE)$ 
24:     $ts_p \leftarrow r_p$ 
25:    store  $\{estimate_p, ts_p\}$ 

26:  {Phase NEWESTIMATE}
27:  s-send  $(r_p, estimate_p, NEWESTIMATE)$  to all
28:  wait until [ received  $(r_p, ACK)$  from
29:    from  $\lceil (n+1)/2 \rceil$  processes ]
30:  s-send  $(estimate_p, DECIDE)$  to all

31: Task participant:
32:  {Phase ESTIMATE}
33:  if  $ts_p \neq r_p$  then
34:    s-send  $(r_p, estimate_q, t, ESTIMATE)$  to  $c_p$ 
35:    wait until [ received  $(r_p, estimate_{c_p},$ 
36:      NEWESTIMATE) from  $c_p$  ]
37:  if  $p \neq c_p$  then
38:     $(estimate_p, ts_p) \leftarrow (estimate_{c_p}, r_p)$ 
39:    store  $\{estimate_p, ts_p\}$ 

40:  {Phase ACK}
41:  s-send  $(r_p, ACK)$  to  $c_p$ 

42: Task skip_round:
43:   $d \leftarrow \mathcal{D}_p$ 
44:  if  $c_p \in d.trustlist$  then {query  $\diamond \mathcal{S}_u$ }
45:    repeat  $d' \leftarrow \mathcal{D}_p$  {query  $\diamond \mathcal{S}_u$ }
46:    until  $[ c_p \notin d'.trustlist$  or  $d'.epoch[c_p] < d'.epoch[c_p]$  or received some message  $(r, \dots)$  such that  $r > r_p ]$ 
47:    terminate task  $\{4phases, participant, coordinator\}$  {abort current round}
48:    repeat  $d \leftarrow \mathcal{D}_p$  until  $d.trustlist \neq \emptyset$  {query  $\diamond \mathcal{S}_u$  to go to a higher round}
49:     $r_p \leftarrow$  the smallest  $r > r_p$  such that  $[(r \bmod n) + 1] \in d.trustlist$  and  $r \geq \max\{r' \mid p \text{ received } (r', \dots)\}$ 
50:    fork task  $4phases$ 

51: upon receive  $m$  from  $q$  do:
52:  if  $m = (estimate, DECIDE)$  and decide( $-$ ) has not yet occurred then {check stable storage about decide}
53:    DECIDE( $estimate$ ) {decide is logged into stable storage}
54:    terminate task  $\{skip\_round, 4phases, participant, coordinator, retransmit\}$ 
55:  if  $m \neq (-, decide)$  and decide( $estimate$ ) has occurred then {check stable storage about decide}
56:    send  $(estimate, DECIDE)$  to  $q$ 

57: upon recovery:
58:  for all  $q \in \Pi \setminus \{p\}$  do  $xmitmsg[q] \leftarrow \perp$ 
59:  if propose( $v_p$ ) has occurred and decide( $-$ ) has not occurred then {check st. storage about propose and decide}
60:    retrieve  $\{r_p, estimate_p, ts_p\}$ 
61:    if  $r_p = \perp$  then  $r_p \leftarrow 1$ ; if  $estimate_p = \perp$  then  $(estimate_p, ts_p) \leftarrow (v_p, 0)$ 
62:    fork task  $\{4phases, retransmit\}$ 

```

---

## B Proofs for the duration of a good period

### B.1 Proofs for Section 4.2.1

Let in the sequel  $\rho_0$  denote the largest round number such that no process has sent a message for round  $\rho_0$  by time  $\tau_G$ , but some process has send a message for round  $\rho_0 - 1$ .

**Proposition B.1.** *For implementing  $\mathcal{P}_{otr}^2(\pi_0)$  with Algorithm 2, we need one “ $\overline{\pi_0}$ -down” good period of length*

$$(6\delta + 3n\phi + 6\phi + 2)\phi + \delta + \phi.$$

*For implementing  $\mathcal{P}_{otr}^{1/1}(\pi_0)$  with Algorithm 2, we need two “ $\overline{\pi_0}$ -down” good period of length*

$$(4\delta + 2n\phi + 4\phi + 1)\phi + \delta + \phi.$$

*Proof.* Follows directly from Theorem 3 below.  $\square$

**Lemma B.2.** *For each round number  $\rho$ , if a message for round  $\rho + 1$  is sent by some process then there is a process that sent previously a message for round  $\rho$ .*

*Proof.* Assume by contradiction that no message for round  $\rho$  is sent, but for round  $\rho + 1$ . The first message for round  $\rho + 1$  is sent by some process  $p$  only if  $p$  received a message for round  $\rho + 1$  (impossible, since it is by assumption the first message), or when  $p$  has taken sufficient receive steps in round  $\rho$ , which starts always with sending a message for  $\rho$ . A contradiction.  $\square$

**Lemma B.3.** *By time  $\tau_G + (2\delta + (n + 2)\phi + 1)\phi$  some process sends a message for round  $\rho_0$ .*

*Proof.* Let  $p$  be the (not necessarily unique) process that has sent a message with the highest round number (equal to  $\rho_0 - 1$ ) at time  $\tau_G$ . By the definition of  $\rho_0$  such a process exists. Because of line 12,  $p$  takes at most  $2\delta + (n + 2)\phi$  receive steps, before sending a new round message (one send step). By the synchrony assumptions,  $p$  takes at least one step in every time interval of length  $\phi$ , so executing  $2\delta + (n + 2)\phi + 1$  steps takes at most  $(2\delta + (n + 2)\phi + 1)\phi$  time units. So  $p$  sends a message for a new round no later than  $\tau_G + (2\delta + (n + 2)\phi + 1)\phi$  (\*). By Lemma B.2 a message is sent for round  $\rho_0$  (\*\*). The result follows from (\*) and (\*\*).  $\square$

**Lemma B.4.** *Let  $p$  be the first (not necessarily unique) process that sends at time  $\tau \geq \tau_G$  a message for a round  $\rho \geq \rho_0$ . Then no process sends a round  $\rho + 1$  message before time  $\tau + 2\delta + (n + 2)\phi + 1$ .*

*Proof.* Assume by contradiction a process  $q$  that is the first to send a round  $\rho + 1$  message before time  $\tau + 2\delta + (n + 2)\phi + 1$ . A process can enter round  $\rho + 1$ , if (a) it receives a round  $\rho + 1$  message, or (b) it makes  $2\delta + (n + 2)\phi$  receive steps

in round  $r$ . (a) cannot be the case, since  $q$  is by assumption the first process in round  $\rho + 1$ . In case (b), making  $2\delta + (n + 2)\phi$  receive steps and a send step for sending the round  $\rho + 1$  message takes at least  $2\delta + (n + 2)\phi + 1$  time, so the send step in round  $\rho + 1$  can be at the earliest at time  $\tau + 2\delta + (n + 2)\phi + 1$ . A contradiction.  $\square$

**Lemma B.5.** *Let  $p$  be the first (not necessarily unique) process that sends a message for a new round  $\rho > \rho_0$  at some time  $\tau > \tau_G$ . Then every process  $q \in \pi_0$  sends a round  $\rho$  message by time  $\tau + \delta + 2\phi$ .*

*Proof.* Since we are in a good period, the message sent by  $p$  in round  $\rho$  is in  $buffer_q$  by time  $\tau + \delta$  (\*). By Lemma B.4, there is no round  $\rho' > \rho$  message in the system before time  $\tau + 2\delta + (n + 2)\phi + 1$  (\*\*). The reception policy “highest round first” together with (\*) and (\*\*) ensure that by time  $\tau + \delta + \phi$  some round  $\rho$  message (not necessarily the one from  $p$ ) is received by  $q$ . Because of lines 17–18,  $q$  switches to round  $\rho$  and executes a send step for round  $\rho$ . The send steps takes at most  $\phi$  time units, i.e.,  $q$  sends a round  $\rho$  message by time  $\tau + \delta + 2\phi$ .  $\square$

**Lemma B.6.** *For all rounds  $\rho \geq \rho_0$ , if some process sends a round  $\rho$  message, every process in  $\pi_0$  executes line 20 with  $R$  containing the round  $\rho$  message from every process in  $\pi_0$ , presuming we are still in the good period at this time.*

*Proof.* Let  $\tau$  be the time where some process  $p \in \pi_0$  sends a round  $\rho$  message. By Lemma B.5, every other process  $q \in \pi_0$  sends a round  $\rho$  message  $m_q^\rho$  by time  $\tau + \delta + 2\phi$ . Message  $m_q^\rho$  is in  $buffer_s$  for all  $s \in \pi_0$  by time  $rt = \tau + 2\delta + 2\phi$ . At time  $rt$ , processes in  $\pi_0$  have taken at most  $2\delta + 2\phi$  receive steps in round  $\rho$ . By Lemma B.4, no messages for higher rounds are in the system before  $\tau + 2\delta + (n + 2)\phi + 1$ . So after  $rt$ , every process in  $\pi_0$  takes at least  $n$  more receive steps in round  $\rho$ , and executes line 20 with  $m_q^\rho$  in  $R$ .  $\square$

**Theorem 3.** *With Algorithm 2, the minimal length of a good period to achieve  $\mathcal{P}_{su}(\pi_0, \rho_0, \rho_0 + x - 1)$  is:*

$$(x + 1)(2\delta + (n + 2)\phi + 1)\phi + \delta + \phi.$$

*Proof.* By Lemma B.6, for all rounds  $\rho \geq \rho_0$ ,  $\mathcal{P}_{su}(\pi_0, \rho, \rho)$  holds. So it remains to compute the length of each of these rounds.

For each round  $\rho > \rho_0$ , every process in  $\pi_0$  takes exactly 1 send step and  $2\delta + (n + 2)\phi$  receive steps. Since each step takes at most  $\phi$  time units, each round  $\rho > \rho_0$  takes at most  $(2\delta + (n + 2)\phi + 1)\phi$  time units.

It remains to compute the time at which processes in  $\pi_0$  ends round  $\rho_0$ . By Lemma B.3, by time  $\tau_G + (2\delta + (n + 2)\phi + 1)\phi$ , a message for round  $\rho_0$  is sent. By Lemma B.5, by time  $\tau_G + (2\delta + (n + 2)\phi + 1)\phi + \delta + 2\phi$  every process

sends a round  $\rho_0$  message. Then every process takes  $2\delta + (n+2)\phi$  receive steps. Each of these receive steps take at most  $\phi$  time units. So all processes have taken their receive steps for round  $\rho_0$  by time  $\tau_G + 2(2\delta + (n+2)\phi + 1)\phi + \delta + \phi$ .  $\square$

**Theorem 5.** *With Algorithm 2, the minimal length of an initial good period to achieve space uniformity for  $x$  consecutive rounds is:*

$$x(2\delta + (n+2)\phi + 1)\phi.$$

*Proof.* The result follows directly from the length of a round in a good period computed in the proof of Theorem 3.  $\square$

## B.2 Proofs for Section 4.2.2

For the proof of correctness, we assume that a good period starts at time  $\tau_G$  and lasts forever. We denote with “a round message for  $\rho$ ” a message  $\langle \text{ROUND}, \rho, - \rangle$ , and with “an init message for  $\rho$ ” a message  $\langle \text{INIT}, \rho, - \rangle$ , respectively. Further, let  $\rho_0$  be the round number of the highest round message that has been sent by any process. The timeout of the algorithm is abbreviated with  $\tau_0 = 2\delta + (2n+1)\phi$ .

**Lemma B.7.** *For any time  $\tau$ , let  $F(\tau)$  be the set of  $f+1$  processes that have the highest round number among all processes, and let  $\rho_{\min}(\tau)$ , resp.  $\rho_{\max}(\tau)$ , be the smallest, resp. largest, round number in  $F(\tau)$ . Then  $\rho_{\max}(\tau) - \rho_{\min}(\tau) \leq 1$ .*

*Proof.* Assume by contradiction that there is a process in  $F(\tau)$  with round number  $\rho < \rho_{\max}(\tau) - 1$ . Let  $p$  be the first process in round  $\rho_{\max}(\tau)$ . To enter this round,  $p$  has to receive either (a) a message  $\langle \text{ROUND}, \rho_{\max}(\tau), - \rangle$  (line 12), (2) a message  $\langle \text{INIT}, \rho_{\max}(\tau)+1, - \rangle$  (line 12), or (3)  $f+1$  messages  $\langle \text{INIT}, \rho_{\max}(\tau), - \rangle$  (line 16). (a) and (b) cannot be the case, since  $p$  is by assumption the first process in round  $\rho_{\max}(\tau)$ . Case (c) is not possible, since by assumption, less than  $f+1$  processes are in a round greater or equal  $\rho_{\max}(\tau) - 1$ . So  $p$  cannot be in round  $\rho_{\max}(\tau)$ , a contradiction.  $\square$

**Lemma B.8.** *Every process in  $\pi_0$  sends message  $\langle \text{ROUND}, \rho_0, - \rangle$  by time  $\tau_G + 3\tau_0\phi + 3\delta + 3n\phi + 6\phi$ .*

*Proof.* We first show that every process in  $\pi_0$  sends message  $\langle \text{ROUND}, \rho_0 - 1, - \rangle$  by time  $\tau_G + 2\tau_0\phi + 2\delta + 2n\phi + 4\phi$ . For this, we distinguish two cases.

For the first case, assume that there is a process  $p \in \pi_0$  with  $r_p = \rho_0 - 1$  at time  $\tau_G$ . Then after at most  $\tau_0$  receive steps,  $p$  sends  $\langle \text{INIT}, \rho_0, - \rangle$ , which is by time  $\tau_G + \tau_0\phi + \phi$ . Because we are in a good period and  $p \in \pi_0$ ,  $\langle \text{INIT}, \rho_0, - \rangle$  is ready for reception at each process  $q \in \pi_0$  at most  $\delta$  time

later; because of the reception policy, and since there are no messages with a higher round number from  $p$ ,  $\langle \text{INIT}, \rho_0, - \rangle$  is received by  $q$  by time  $\tau_G + \tau_0\phi + \delta + n\phi + \phi$ . At most  $\phi$  times later  $q$  sends  $\langle \text{ROUND}, \rho_0 - 1, - \rangle$ , and we are done.

For the second case, where there is no process in  $\pi_0$  with round number  $\rho_0 - 1$ , we can conclude from Lemma B.7, that there is at least some process  $p \in \pi_0$ , with  $r_p = \rho_0 - 2$ . By the same argument as above, every process in  $\pi_0$  receives  $\langle \text{INIT}, \rho_0 - 1, - \rangle$  by time  $\tau_G + \tau_0\phi + \delta + n\phi + \phi$ , then (optionally, since this might have happened already before) sends  $\langle \text{ROUND}, \rho_0 - 2, - \rangle$  and takes at most  $\tau_0$  receive steps before sending  $\langle \text{INIT}, \rho_0 - 1, - \rangle$ . Thus, by time  $\tau_G + 2\tau_0\phi + \delta + n\phi + 3\phi$ , every process in  $\pi_0$  sends  $\langle \text{INIT}, \rho_0 - 1, - \rangle$ . Because of  $|\pi_0| = n - f$ , every process in  $\pi_0$  receives thus more than  $f$  messages  $\langle \text{INIT}, \rho_0 - 1, - \rangle$  and sends thus  $\langle \text{ROUND}, \rho_0 - 1, - \rangle$  by time  $\tau_G + 2\tau_0\phi + 2\delta + 2n\phi + 4\phi$ .

In both cases, since  $\langle \text{ROUND}, \rho_0 - 1, - \rangle$  is sent by every process in  $\pi_0$  by time  $\tau_G + 2\tau_0\phi + 2\delta + 2n\phi + 4\phi$ , but after  $\tau_G$ , after at most  $\tau_0$  receive steps (which takes at most  $\tau_0\phi$  time units), every process sends  $\langle \text{INIT}, \rho_0, - \rangle$ , which takes at most  $\Phi$  time. This message is ready for reception at each process in  $\pi_0$  at most  $\delta$  time later, and because of the reception policy, because there are no messages with a higher round number from  $p$ , this message is received at most  $n\phi$  time later. At most  $\phi$  times later  $\langle \text{ROUND}, \rho_0, - \rangle$  is sent. The Lemma follows.  $\square$

**Lemma B.9.** *Let  $p \in \Pi$  be the first (not necessarily unique) process that sends at time  $\tau \geq \tau_G$  message  $\langle \text{ROUND}, \rho, - \rangle$  with  $\rho \geq \rho_0$ . Then no process sends  $\langle \text{ROUND}, \rho+1, - \rangle$  before time  $\tau + \tau_0 + 1$ . Moreover, no process sends  $\langle \text{INIT}, \rho+2, - \rangle$  before time  $\tau + \tau_0 + 1$ .*

*Proof.* For the ROUND message, assume by contradiction that process  $q \in \Pi$  is the first process that sends  $\langle \text{ROUND}, \rho+1, - \rangle$  before time  $\tau + \tau_0 + 1$ . Process  $q$  can send this message if it (a) receives a message  $\langle \text{ROUND}, \rho+1, - \rangle$ , (b) a message  $\langle \text{INIT}, \rho+2, - \rangle$ , or (c)  $f+1$  messages  $\langle \text{INIT}, \rho+1, - \rangle$ . Cases (a) and (b) are excluded, since  $q$  is by assumption the first process in round  $\rho+1$ . In case (c), since  $|\pi_0| = n - f$ , at least one of the messages  $\langle \text{INIT}, \rho+1, - \rangle$  is from a process in  $\pi_0$ . In order to send  $\langle \text{INIT}, \rho+1, - \rangle$ , a process makes  $\tau_0$  receive steps in round  $\rho$ , which takes for processes in  $\pi_0$  at least  $\tau_0$  time, plus the send step that costs at least one time unit. So message  $\langle \text{ROUND}, \rho+1, - \rangle$  is sent by  $q$  at the earliest at time  $\tau + \tau_0 + 1$ . A contradiction.

For the INIT message the result follows from the fact that a message  $\langle \text{INIT}, \rho+2, - \rangle$  is only sent after a message  $\langle \text{ROUND}, \rho+1, - \rangle$ .  $\square$

**Lemma B.10.** *Let  $p$  be the first (not necessarily unique) process from  $\pi_0$  that sends message  $\langle \text{ROUND}, \rho, - \rangle$  with  $\rho \geq \rho_0$  at some time  $\tau \geq \tau_G$ . Then every other process in  $\pi_0$  sends message  $\langle \text{ROUND}, \rho+1, - \rangle$  by time  $\tau + \delta + n\phi + \phi$ .*

*Proof.* Since we are in a good period, the message sent by  $p$  in round  $\rho$  is in  $buffer_q$  by time  $\tau + \delta$  (\*). By Lemma B.9, there is no message  $\langle \text{ROUND}, \rho', - \rangle$  with  $\rho' > \rho$  in the system before time  $\tau + \tau_0 + 1$ . Since  $\tau_0 = 2\delta + 2n\phi + \phi$  no such message is the system before time  $\tau + 2\delta + 2n\phi + \phi + 1$  (\*\*). The reception policy “highest round first from each process in a round robin fashion” together with (\*) and (\*\*), ensures that by time  $\tau + \delta + n\phi$  some message  $\langle \text{ROUND}, \rho, - \rangle$  (not necessarily the one from  $p$ ) is received by  $q$ . Because of lines 14–15,  $q$  switches to round  $\rho$  and sends message  $\langle \text{ROUND}, \rho, - \rangle$ . The send step takes at most  $\phi$  time units, i.e.,  $q$  sends message  $\langle \text{ROUND}, \rho, - \rangle$  by time  $\tau + \delta + (n + 1)\phi$ .  $\square$

**Lemma B.11.** *For all rounds  $\rho \geq \rho_0$ , if some process  $p \in \pi_0$  sends message  $\langle \text{ROUND}, \rho, - \rangle$ , every process  $q \in \pi_0$  executes line 20 with  $R$  containing the round message  $\langle \text{ROUND}, \rho, - \rangle$  from every process in  $\pi_0$ , presuming we are still in the good period at this time.*

*Proof.* Let  $\tau$  be the time where some process  $p \in \pi_0$  sends message  $\langle \text{ROUND}, \rho, - \rangle$ . By Lemma B.10, every other process  $q \in \pi_0$  sends message  $m_q^\rho = \langle \text{ROUND}, \rho, - \rangle$  by time  $\tau + \delta + n\phi + \phi$ . Message  $m_q^\rho$  is in  $buffer_s$  for all  $s \in \pi_0$  by time  $\tau_\rho = \tau + 2\delta + n\phi + \phi$ . At time  $\tau_\rho$ , processes in  $\pi_0$  have taken at most  $2\delta + n\phi + \phi$  receive steps in round  $\rho$ . By Lemma B.9, no messages for higher rounds are in the system before  $\tau + 2\delta + (2n + 1)\phi + 1$ . So after  $\tau_\rho$ , every process in  $\pi_0$  takes at least  $n$  more receive steps in round  $\rho$ , and executes line 22 with  $m_q^\rho$  in  $R$ .  $\square$

**Theorem 6.** *With Algorithm 3 and  $f < n/2$ , the minimal length of a good period to achieve  $\mathcal{P}_k(\pi_0, \rho_0 + 1, \rho_0 + x)$  is*

$$(x + 2)[\tau_0\phi + \delta + n\phi + 2\phi] + \tau_0\phi = \\ = (x + 2)[(2\delta + 2n\phi + \phi)\phi + \delta + 2n\phi + 2\phi] + (2\delta + 2n\phi + \phi)\phi.$$

*Proof.* By Lemma B.11, for  $x \geq 1$ , the predicate  $\mathcal{P}_k(\pi_0, \rho_0 + 1, \rho_0 + x)$  holds. It remains to compute the length of each of these rounds.

According to Lemma B.8, the last process starts round  $\rho_0$  by time  $\tau_G + 3\tau_0\phi + 3\delta + 3n\phi + 6\phi$ . For each round  $\rho \geq \rho_0$ , every process in  $\pi_0$  takes at most  $\tau_0$  receive steps and one send step for message  $\langle \text{ROUND}, \rho, - \rangle$ . This leads to at most  $\tau_0\phi + \phi$  time units. Then we have the sending and reception of the INIT messages: one send step (at most  $\phi$  time units), the transmission (at most  $\delta$ ), and at most  $n$  receive steps ( $f + 1 < n$ ). We get a total of  $\tau_0\phi + \delta + n\phi + 2\phi$  for a round. However, the INIT messages can be ignored for the last round, resulting in a duration of  $\tau_0\phi + \phi$  for the last round. The result follows.  $\square$

**Theorem 7.** *With Algorithm 3 and  $f < n/2$ , the minimal length of an initial good period to implement  $\mathcal{P}_k(\pi_0, 1, x)$  is:*

$$(x - 1)[\tau_0\phi + \delta + n\phi + 2\phi] + \tau_0\phi + \phi.$$

*Proof.* The result follows directly from the length of a round in a good period computed the proof of Theorem 6, and taking into account, as in the proof of Theorem 6, that the INIT messages can be ignored for the last round.  $\square$

## C Proof of Algorithm 4

Algorithm 4 can more easily presented as a *translation*, as formally defined in [6]. For sake of self-containment and conciseness, we chose for Algorithm 4 a presentation in the system level model instead of introducing the concept of a translation.

For an introduction into translations we refer to [6]. It is easy to see that Algorithm 4 is obtained from the translation in Algorithm 7, by changing  $Known_p$ , which is a set of  $\langle message, process \rangle$  in Algorithm 4, to be a set of processes in Algorithm 7.

**Algorithm 7** Translating  $\mathcal{P}_k(-, -, -)$  to  $\mathcal{P}_{su}(-, -, -)$  (adapted from [6]).

---

```

1: Initialization:
2:    $Listen_p \in 2^\Pi$ , initially  $\Pi$ 
3:    $NewHO_p \in 2^\Pi$ 
4:    $Known_p \in 2^\Pi$ , initially  $\{p\}$ 
5: Round  $r$ :
6:    $S_p^r$  :
7:     send  $\langle Known_p \rangle$  to all processes
8:    $T_p^r$  :
9:      $Listen_p \leftarrow Listen_p \cap HO(p, r)$ 
10:    if  $r \not\equiv 0 \pmod{f+1}$  then
11:       $Known_p \leftarrow Known_p \cup \bigcup_{q \in Listen_p} Known_q$ 
12:    else
13:       $NewHO_p \leftarrow \{s \mid s \in Known_q \text{ for } n - f \text{ processes } q \in Listen_p\}$ 
14:       $Listen_p \leftarrow \Pi$ 
15:       $Known_p \leftarrow \{p\}$ 

```

---

Thus we have to show only the correctness of Algorithm 7, which can be done in a more elegant way.

We extend the proof in Section 4.3 of [6] to show that Algorithm 7 translates  $\mathcal{P}_k(\Pi_0, r_1, r_1 + f)$  to  $\mathcal{P}_{su}(\Pi_0, \mathcal{R}, \mathcal{R})$  in  $f + 1$  rounds, where  $f = n - |\Pi_0|$ . The proof is taken from [6] and modified at a few crucial points. In contrast to [6], we have to assume  $n > 2f$ .

We fix a macro-round  $\mathcal{R}$  and introduce some piece of notation relative to  $\mathcal{R}$ . Let  $r_1, \dots, r_{f+1}$  denote the sequence of the  $f + 1$  rounds that form  $\mathcal{R}$ . Let  $K_{\Pi_0}(\mathcal{R})$  denote the  $\Pi_0$ -restricted kernel of macro-round  $\mathcal{R}$ , i.e.,

$$K_{\Pi_0}(\mathcal{R}) = \bigcap_{r=r_1}^{r_f+1} K_{\Pi_0}(r).$$

We further define the restricted co-kernel  $coK_{\Pi_0}(\mathcal{R}) = \Pi \setminus K_{\Pi_0}(\mathcal{R})$ .

**Lemma C.1.** *If  $\mathcal{P}_k(\Pi_0, r_1, r_f)$  holds, then  $coK_{\Pi_0}(\mathcal{R}) \cap \Pi_0 = \emptyset$ .*

*Proof.* If  $\mathcal{P}_k(\Pi_0, r_1, r_f)$  holds, we have  $\Pi_0 \subseteq K_{\Pi_0}(\mathcal{R})$ . The result follows directly from the definition of  $coK_{\Pi_0}(\mathcal{R})$ .  $\square$

We say that process  $p$  knows process  $s$  at round  $r$  if  $s \in Known_p^{(r)}$ . If  $s \in Known_p^{(r)} \setminus Known_p^{(r-1)}$ ,  $q \in Listen_p^r$ , and  $s \in Known_q^{(r-1)}$ , then we say that  $p$  hears of  $s$  from  $q$  at round  $r$ . Finally, process  $s$  is said to be good (at macro-round  $\mathcal{R}$ ) if  $s$  is known by all processes in  $\Pi_0$  at round  $r_f$ ; otherwise  $s$  is bad. In other words, the set of good processes is defined by

$$Good = \bigcap_{p \in \Pi_0} Known_p^{(r_f)}.$$

**Lemma C.2.**

$$K_{\Pi_0}(\mathcal{R}) \subseteq \bigcap_{r \in \Pi_0, r \in \{r_1, \dots, r_{f+1}\}} Listen_p^{(r)}$$

*Proof.* It is immediately from the definition of  $Listen_p$  that for any process  $p$ ,

$$\bigcap_{r=1}^{r_f} Listen_p^{(r)} = \bigcap_{r=1}^{r_f} HO(p, r)$$

and

$$HO(p, r_{f+1}) \subseteq Listen_p^{r_f+1}.$$

The result follows directly.  $\square$

**Lemma C.3.** *Any process  $p$  in  $K_{\Pi_0}(\mathcal{R})$  is a good process.*

*Proof.* Let  $p$  be any process in  $K_{\Pi_0}(\mathcal{R})$ . By lines 15 and 11, it follows that all processes in  $\Pi_0$  know  $p$  at the end of round  $r_f$ . This shows that  $p$  is a good process.  $\square$

**Lemma C.4.** *If a process  $p \in \Pi_0$  hears of some process  $s$  at round  $r_k$ , then there exist  $k-1$  processes  $p_1, \dots, p_{k-1}$ , each different from  $p$  and  $s$ , such that  $p_1$  hears of  $s$  from  $s$  at round  $r_1$ ,  $p_2$  hears of  $s$  from  $p_1$  at round  $r_2$ ,  $\dots$ ,  $p_{k-1}$  hears of  $s$  from  $p_{k-2}$  at round  $r_{k-1}$  and  $p$  hears of  $s$  from  $p_{k-1}$  at round  $r_k$ . Moreover, processes  $p_1, \dots, p_{k-2}$  and  $s$  belong to  $coK_{\Pi_0}(\mathcal{R})$ .*

*Proof.* Since  $p$  hears of process  $s$  at round  $r_k$ , there exists some process  $p_{k-1}$ , such that  $p_{k-1} \in Listen_p^{(r_k)}$  and  $s \in Known_{p_{k-1}}^{(r_{k-1})}$ . Since  $Listen_p$  is non-increasing,  $p_{k-1} \in Listen_{p_{k-1}}^{(r_{k-1})}$ . This implies that  $p_{k-1}$  hears of  $s$  at round

$r_{k-1}$  since  $p$  does not know  $s$  at this round. In turn, there exists some process  $p_{k-2}$ , such that  $p_{k-2} \in Listen_{p_{k-1}}^{(r_{k-1})}$ , and  $s \in Known_{p_{k-2}}^{(r_{k-2})}$ . From

$$s \in Known_{p_{k-2}}^{(r_{k-2})} \text{ and } s \notin Known_{p_{k-1}}^{(r_{k-1})},$$

we deduce that  $p_{k-2} \notin Listen_p^{(r_{k-1})}$ . By Lemma C.2 and  $p \in \Pi_0$ , we have  $p_{k-2} \in coK_{\Pi_0}(\mathcal{R})$ .

Step by step, we exhibit  $k-1$  processes  $p_1 \dots p_{k-1}$ , such that for any index  $i$ ,  $1 \leq i \leq k-1$ ,

$$s \notin Known_{p_i}^{(r_{i-1})}, s \in Known_{p_i}^{(r_i)}, \text{ and } p_{i-1} \in Listen_{p_i}^{(r_i)}.$$

For any index  $i$  such that  $2 \leq i \leq k-1$ , we have both

$$s \in Known_{p_{i-1}}^{(r_{i-1})} \text{ and } s \notin Known_{p_i}^{(r_i)}.$$

Therefore,  $p_{i-1} \notin Listen_p^{(r_i)}$ . By Lemma C.2 and  $p \in \Pi_0$ , we deduce that  $p_{i-1}$  belongs to  $coK_{\Pi_0}(\mathcal{R})$ . Similarly, we have

$$s \notin Known_p^{(r_1)},$$

and so  $s$  belongs to  $coK_{\Pi_0}(\mathcal{R})$ , too. Thus it follows that all the processes  $p_1, \dots, p_{k-2}$ , and  $s$  are in  $coK_{\Pi_0}(\mathcal{R})$ .  $\square$

**Lemma C.5.** *If process  $p \in \Pi_0$  knows some bad process  $s$  at the end of round  $r_{f+1}$ , then  $p$  has heard of  $s$  by the end of the round  $r_f$ , i.e.,*

$$\forall p \in \Pi_0 : s \in Known_p^{r_f+1} \wedge s \notin Good \Rightarrow s \in Known_p^{(r_f)}.$$

*Proof.* Let  $s$  be a bad process; so there exists some process  $q \in \Pi_0$  such that  $s \notin Known_q^{(r_f)}$ . Suppose for contradiction that  $p$  hears of  $s$  at round  $r_{f+1}$ . By Lemma C.4, there are  $f$  processes  $p_1, \dots, p_f$  each different from both  $p$  and  $s$  such that  $p$  hears of  $s$  from  $p_f$  at round  $r_{f+1}$ , and processes  $p_1, \dots, p_{f-1}$ , and  $s$  are all in  $coK_{\Pi_0}(\mathcal{R})$ . Since  $s \notin Known_q^{(r_f)}$ ,  $Listen_q$  contains neither  $p$  nor  $p_f$  at this round. Therefore  $p$  and  $p_f$  are also in  $coK_{\Pi_0}(\mathcal{R})$ . So  $|coK_{\Pi_0}(\mathcal{R})|$  is larger than  $f$ . By Lemma C.1, we have in this case  $|\Pi_0| < n - f$  which contradicts the fact that  $|\Pi_0| = n - f$ .  $\square$

**Lemma C.6.** *A process is good iff it is known by some process in  $\Pi_0$ , i.e.,*

$$s \in Good \Leftrightarrow \exists p \in \Pi_0 : s \in Known_p^{(r_f)}$$

*Proof.* By definition, a good process is known by all processes in  $\Pi_0$  at round  $r_f$ .

Conversely, let  $s$  be any process known by some process  $p \in \Pi_0$  at round  $r_f$ . Assume, for the sake of contradiction, that  $s$  is bad. Since  $p \in \Pi_0$ , every process  $q \in \Pi_0$  receives a message from  $p$  at round  $r_{f+1}$ , and so  $Known_q$  contains  $s$  at the end of round  $r_{f+1}$ . By Lemma C.5, we deduce that every process in  $\Pi_0$  already knows  $s$  at round  $r_f$ . This contradicts that  $s$  is a bad process.  $\square$



**Lemma C.7.** *For any process  $p \in \Pi_0$ , at the end of round  $r_{f+1}$ ,  $NewHO_p$  is composed of all the good processes, i.e.,  $NewHO_p^{(r_{f+1})} = Good$ .*

*Proof.* We have  $Good \subseteq NewHO_p^{(r_{f+1})}$ , since by definition of  $Good$ , for all processes  $q \in \Pi_0$ ,  $Good \subseteq Known_q^{(r_{f+1})}$ .

Conversely, let  $s$  be any process in  $NewHO_p^{(r_{f+1})}$ ; by line 13,  $s$  is known at round  $r_f$  by  $n - f$  processes in

$Listen_p^{(r_{f+1})}$ , which includes because of  $n > 2f$  also a process from  $\Pi_0$ . By Lemma C.6, it follows that  $s$  is a good process.  $\square$

Therefore we have:

**Theorem 8.** *For  $n > 2f$ , Algorithm 7 is a  $f + 1$  translation of  $\mathcal{P}_k(\Pi_0, -, -)$  into  $\mathcal{P}_{su}(\Pi_0, -, -)$ .*