

# Swift Algorithms for Repeated Consensus

Fatemeh Borran    Martin Hutle    Nuno Santos    André Schiper  
*Ecole Polytechnique Fédérale de Lausanne (EPFL)*  
1015 Lausanne, Switzerland  
{firstname}.{lastname}@epfl.ch

**Abstract**—We introduce the notion of a *swift algorithm*. Informally, an algorithm that solves the repeated consensus is *swift* if, in a partial synchronous run of this algorithm, eventually no timeout expires, *i.e.*, the algorithm execution proceeds with the actual speed of the system. This definition differs from other efficiency criteria for partial synchronous systems.

Furthermore, we show that the notion of *swiftness* explains why failure detector based algorithms are typically more efficient than round-based algorithms, since the former are naturally *swift* while the latter are naturally *non-swift*. We show that this is not an inherent difference between the models, and provide a round implementation that is *swift*, therefore performing similarly to failure detector algorithms while maintaining the advantages of the round model.

## I. INTRODUCTION

Timeouts are often required to solve problems in distributed computing. Due to the FLP impossibility result [1], there is a need of some minimal synchrony assumptions for solving the consensus problem, and timeouts are the dominant mechanism for algorithms to make use of synchrony assumptions.

Timeouts are often chosen conservatively, so that an algorithm is correct for a large number of real-life scenarios. However, timeouts should be used only to cope with faults, and not slow down the execution time in good cases. As an example, when implementing communication-closed synchronous rounds in a synchronous message passing system, after a process sent its messages for a certain round it usually waits for a timeout, before it terminates the round and sends its messages for the next round. However, in many runs of the algorithm, a process might have received all messages from other alive processes already long before that. It would be favorable to start the next round immediately after all messages from correct processes are received. This is, for example, the case for an algorithm that uses a  $\diamond\mathcal{P}$  failure detector (FD). Here, a process waits for a message from some process  $p$  until  $p$  is in the FD output. If  $p$  has crashed, this involves waiting for a timeout, but only once: later rounds profit from the fact that the failure detector “remembers” information about faults. We formally capture such a behavior by the definition of *swift*, which we define in the context of repeated consensus [2]. The main intuition behind our definition is that *swift* algorithms make progress

at the speed of the system, and therefore, are more “efficient” than non-*swift* algorithms. A *swift* algorithm for a repeated problem is thus one in which eventually all instances of the problem are “efficient”.

In more detail, for the definition of *swift* we look at partial synchronous runs, *i.e.*, runs where a bound  $\Delta$  on the transmission delay eventually holds forever.<sup>1</sup> For the good period of such a run, that is the partial run  $R$  in which bound  $\Delta$  holds, we can define the actual transmission delay  $\delta(R)$  as the maximum of all transmission delays in  $R$ . Such an actual transmission delay can be much smaller than the bound  $\Delta$ . If in this case the execution time for each instance of the repeated consensus eventually depends only on  $\delta(R)$  (in contrast to  $\Delta$ ), the algorithm is *swift*.

While intuitively *swift* algorithms progress at the speed of messages in good periods, and non-*swift* algorithms progress sometimes only by the expiration of timeouts, we refrained from calling these two classes of algorithms *message-driven* and *timeout driven*. This is because the term *message-driven* is used in [3], [4] with a different meaning, namely to refer to the way events are generated at a process. If processes are allowed to measure time (*e.g.*, with clocks or step counting), then it is possible to construct *message-driven* algorithms (according to this definition) that are not *swift*. On the other hand, if processes use an adaptive timeout, then the algorithm can be *swift* despite timeout expiration. Thus these terms are not suitable to precisely characterize this class of algorithms.

Other notions of efficiency for distributed algorithms have been considered. The term *fast* has been used to refer to (consensus) algorithms that solve consensus with less communication steps in favorable cases [5]. A favorable case corresponds usually to an execution without faults that is synchronous from the beginning. On the contrary, the definition of *swift* is related to the execution *time* of an algorithm in the context of *repeated* consensus. Furthermore, the definition of *swift* considers also runs with faults. The notion of *fast* is orthogonal to the notion of *swift*: it is possible to design both, *fast* algorithms that are *swift* and *fast* algorithms that are not *swift*. The same argument holds for *early terminating* algorithms [6].

<sup>1</sup>Note that such a run exists also, *e.g.*, in an asynchronous system, and all runs of a synchronous systems are of course also partial synchronous. The definition is thus not limited to partial synchronous systems.

The paper makes the following two contributions. The first contribution is the definition of swift algorithms that we just discussed. The second contribution is a new implementation of a communication-closed rounds in a partial synchronous system with crash faults. This new implementation leads to swift round-based consensus algorithms, while previous round implementations, including those described in [7], [8] are not swift. This result is especially relevant in the context of comparing advantages and drawbacks of the failure detector approach [9] with the round-based approach [7], [10] for solving agreement problems. Indeed, failure detector based algorithms, despite the usage of timeouts in the implementation of the failure detector algorithm, are naturally swift. On the other hand, round implementations in a partial synchronous model have some advantages over FD based implementations [11]. Our new solution thus combines the advantages of both approaches.

The rest of the paper is structured as follows. In the next section, we specify our model and give a formal definition of *swift*. Then, in Section III we show a simple round-based consensus algorithm that is not swift, and in Section IV we show that the same consensus algorithm expressed using a failure detector is swift. In Section V we present our main contribution: we show a new implementation of rounds that is swift. Section VI validates the theoretical analysis with experimental results comparing the swift and non-swift implementations. Section VII concludes the paper.

## II. DEFINITIONS AND MODEL

We consider a system of  $n$  processes connected by a message-passing network. Among these  $n$  processes, at most  $f$  may crash. We attach an in-queue and an out-queue to each process, where for repeated consensus, the in-queue contains the consensus proposals, and the out-queue contains the consensus decisions. Processes execute an algorithm by taking steps, where a step can be either a send step  $\langle p, \text{SEND}, m \rangle$ , in which a process sends a message to another process, a receive step  $\langle p, \text{RECEIVE}, S \rangle$ , in which a (possibly empty) set  $S$  of messages is received, an input step  $\langle p, \text{IN}, I \rangle$ , in which a value is read from  $p$ 's in-queue, or an output step  $\langle p, \text{OUT}, O \rangle$ , in which a value is output to  $p$ 's out-queue. We denote with  $In_p$  (resp.  $Out_p$ ) the in-queue (resp. out-queue) of process  $p$ . In each step a process also performs a state transition.

We assume an abstract global discrete time. Without loss of generality, at each time  $t$  at least one process makes a step. A single process can make at most one step at any time. Processes measure time by counting their own steps.

Channels satisfy validity and integrity.<sup>2</sup> Channels are reliable if additionally the following property holds:

<sup>2</sup>*Validity*: A message  $m$  that is received by  $q$  was previously sent by some process  $p$  to  $q$ ; *Integrity*: A message  $m$  that is sent from  $p$  to  $q$  is received by  $q$  at most once.

*Reliability*: If message  $m$  is sent from  $p$  to  $q$  and  $q$  performs an infinite number of receive steps, then eventually  $m$  is received by  $q$ .

We consider partial synchronous runs, defined by a bound  $\Phi$  on the process relative speeds and a bound  $\Delta$  on the transmission delay of messages [7]. For a run  $R$ , we say that the process speed bound  $\Phi$  holds in  $R$  if, in any partial run of  $R$  that contains  $\Phi$  steps, every non-crashed process makes at least one step. Further, we say that the transmission delay  $\Delta$  holds in  $R$  after some time  $t_0$  if (i) any message sent by  $p$  to  $q$  at time  $t \geq t_0$  is received the latest in the first receive step after  $t + \Delta$ ; and (ii) every message sent before  $t_0$  is received the latest in the first receive step after  $t_0 + \Delta$ .

**Definition 1** (Partial synchrony). *A run  $R$  is  $(\Delta, \Phi)$ -partial synchronous if there is a time  $GST$  (Global Stabilization Time) such that after  $GST$  the transmission delay bound  $\Delta$  holds, the process speed bound  $\Phi$  holds, and no process crashes after  $GST$ .*

We call the time interval  $(GST, \infty)$  the *good period* of  $R$ . We say a *system* is  $(\Delta, \Phi)$ -partial synchronous if every run  $R$  of the system fulfills Definition 1. To simplify the presentation, we assume  $\Phi = 1$ , and write  $\Delta$ -partial synchronous for  $(\Delta, 1)$ -partial synchronous.

**Definition 2** (Actual parameters). *Let  $R'$  be a partial run. Then  $\delta(R')$  denotes the maximum transmission delay of the partial run  $R'$ , i.e., the smallest value  $\bar{\delta}$  such that the transmission delay is bounded by  $\bar{\delta}$  in the partial run  $R'$ .*

If  $R'$  is the good period of a  $\Delta$ -partial synchronous system, then  $\delta(R') \leq \Delta$ . When  $R'$  is clear from the context, we simply write  $\delta$ . The bound  $\Delta$  may be known or unknown. For the algorithms in this paper, we assume that  $\Delta$  is known. However,  $\delta$  is unknown (it represents the performance metric of a single run).

### A. Repeated consensus

We focus on the repeated consensus problem. The in-queue and out-queue are queues of pairs  $\langle i, v \rangle$ , where  $i$  is a consensus instance number and  $v$  a value. In the repeated consensus problem, for each instance  $i$ , the following holds:

- *Validity*: For every process  $p$ , if  $\langle i, v \rangle \in Out_p$  then there exists some process  $q$  such that  $\langle i, v \rangle \in In_q$ .
- *Uniform agreement*: For all processes  $p, q$ , if  $\langle i, v \rangle \in Out_p$  and  $\langle i, v' \rangle \in Out_q$  then  $v = v'$ .
- *Termination*: For every correct process  $p$  there exists  $v$  such that  $\langle i, v \rangle \in Out_p$ .

### B. Swift algorithms

Before giving a formal definition of swift, we need to formalize the notion of execution time of an instance of consensus.

**Definition 3** (Execution time). *Consider a run  $R$  of a repeated consensus algorithm. The execution time  $\tau_i(R)$  of*

**Algorithm 1** OneThirdRule (OTR) (code of process  $p$ )

---

```

1: State:
2:    $x_p \in V$ 
3:    $decision_p \in V$ 

4: Round  $r$ :
5:    $S_p^r$ :
6:   send  $\langle x_p \rangle$  to all processes
7:    $T_p^r$ :
8:   if number of values received  $> 2n/3$  then
9:      $x_p \leftarrow x$  smallest most often received value
10:    if more than  $2n/3$  values received are equal to  $v$  then
11:       $decision_p \leftarrow v$ 

```

---

instance  $i$  of consensus is defined as follows. Let  $t_{in} = \max\{t : \langle i, v \rangle \text{ is taken from } In_i \text{ at some process } p \text{ at time } t\}$ ,  $t_{out} = \max\{t : \langle i, v \rangle \text{ is output to } Out_i \text{ at some process } p \text{ at time } t\}$ . Then  $\tau_i(R) = t_{out} - t_{in}$ .

Let  $A(\Delta)$  denote algorithm  $A$  parametrized with  $\Delta$ .<sup>3</sup>

**Definition 4** (Swift algorithm). An algorithm  $A(\Delta)$  that solves repeated consensus is swift if there are constants  $k, c \in \mathbb{N}$  such that for every run  $R$  of  $A(\Delta)$  that is  $\Delta$ -partial synchronous with good period  $R'$ , and includes an infinite number of instances, there exists  $i'$  such that for all instance  $i \geq i'$ , we have  $\tau_i(R) \leq k\delta(R') + c$ .

Note that this definition does not refer to timeouts. Our definition only depends on the relation between system properties (i.e., transmission delays) and algorithm properties (i.e., execution time), and therefore avoids any reference to timeout expiration.

### III. A NON-SWIFT ROUND-BASED ALGORITHM

We illustrate swiftness and non-swiftness on simple consensus algorithms. The algorithms we consider belong all to the same class of consensus algorithms, i.e., algorithms that require  $f < n/3$ . We consider a round-based algorithm, namely the OneThirdRule (OTR) consensus algorithm from [10], see Algorithm 1. The round-based model has been introduced in [7]. In each round  $r$ , a process sends its estimate  $x_p$  to all processes (line 6) and then, after an implicit receive step where only messages of round  $r$  may be received, performs the state transition function  $T_p^r$  (lines 8 to 11). Algorithm 1 is always safe. For liveness, we need two rounds in which the set  $\Pi_0$  of alive processes (at least  $2n/3$ ) receives all messages from processes in  $\Pi_0$ , and only from these processes. This property is called *space uniformity*. It can be ensured by the round implementation layer during the good period of a partially synchronous system.

The implementation of the round structure is given by Algorithm 2. It is an extension of the implementation given in [11] with support for repeated instances of consensus.

<sup>3</sup>For models with known bounds on transmission delays,  $\Delta$  represent this knowledge. For models with unknown  $\Delta$ , or asynchronous algorithms, we assume  $A(\Delta)$  to be a constant function, i.e.,  $A(\Delta)$  represents one single algorithm.

**Algorithm 2** A non-swift round implementation (code of  $p$ )

---

```

1:  $r_p \leftarrow 1$  /* round number */
2:  $next\_r_p \leftarrow 1$ 
3:  $Rcv_p \leftarrow \emptyset$  /* set of received messages */
4:  $\forall i \in \mathbb{N} : state_p[i] \leftarrow \perp$  /* state of instance  $i$  */

5: while true do
6:    $I \leftarrow input()$ 
7:   for all  $\langle i, v \rangle \in I$  do
8:      $state_p[i] \leftarrow \langle v, \perp \rangle$ 
9:   for all  $i : state_p[i] \neq \perp$  do
10:     $msgs[i] \leftarrow S_p^r(state_p[i])$ 
11:   for all  $q \in \Pi$  do
12:     $M_q \leftarrow \{\langle i, msgs[i][q] \rangle : state_p[i] \neq \perp\}$ 
13:    send( $M_q, r_p, p$ ) to  $q$ 
14:    $i_p \leftarrow 0$ 
15:   while  $next\_r_p = r_p$  do
16:      $i_p \leftarrow i_p + 1$ 
17:     if  $i_p \geq TO$  then
18:        $next\_r_p \leftarrow r_p + 1$ 
19:       receive( $M$ )
20:        $Rcv_p \leftarrow Rcv_p \cup M$ 
21:        $next\_r_p \leftarrow \max\{r : \langle -, r, - \rangle \in Rcv_p\} \cup \{next\_r_p\}$ 
22:    $O \leftarrow \emptyset$ 
23:   for all  $i : state_p[i] \neq \perp$  do
24:     for all  $r \in [r_p, next\_r_p - 1]$  do
25:        $\forall q \in \Pi : M_r[q] \leftarrow m$  if  $\exists M \langle M, r, q \rangle \in Rcv_p$ 
26:          $\wedge \langle i, m \rangle \in M$ , else  $\perp$ 
27:        $state_p[i] \leftarrow T_p^r(state_p[i], M_r)$ 
28:       if the first time  $state_p[i].decision \neq \perp$  then
29:          $O \leftarrow O \cup \langle i, state_p[i].decision \rangle$ 
30:   output( $O$ )
31:    $r_p \leftarrow next\_r_p$ 

```

---

Each iteration of the outermost loop is composed of three parts: *input & send* part, *receive* part and *comp. & output* part. In the *input & send* part, the process queries the input queue for new proposals (line 6), initializes new slots in the *state* vector for each new proposal (line 8), calls the send function of all active consensus instances (line 10), and sends the resulting messages (line 13). The process then starts the *receive* part, where it waits for messages until either the timeout  $TO$  expires (line 17) or it receives a message from a higher round (line 21). Finally, in the *comp. & output* part, the process calls the state transition function of each active instance (line 26), and outputs any new decisions (line 29). Note that some rounds may be partially skipped (no message sent, no message received, only transition function executed): this happens whenever a message from higher round is received.

In Appendix A we prove the correctness of the round implementation for  $TO \geq 2\Delta + 2n + 5$ . We also show that for each instance  $i$  of consensus started after  $GST$ , we have an execution time  $\tau_i \leq 2TO + \delta + 3n + 6$ . This defines the maximum execution time. We now show that the implementation is not swift by computing the minimum execution time for each instance of consensus.

**Lemma 1.** Consider Algorithm 2 with  $TO \geq 2\Delta + 2n + 5$ ,  $n > 3f$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Let  $r_0$  be

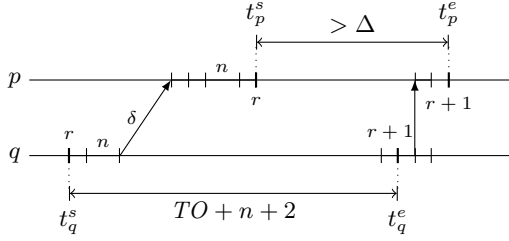


Figure 1. Illustration for Lemma 1

the first new round that is started after GST. Then for all instances  $i$  started in a round  $r \geq r_0$ , we have an execution time  $\tau_i > \Delta$ .

*Proof:* We prove the result by showing that, for every round  $r \geq r_0$ , every process  $p$  stays in round  $r$  for more than  $\Delta$  time.

Let  $t_p^s$  and  $t_p^e$  be the time when  $p$  starts and finishes round  $r$ , respectively. Process  $p$  may finish round  $r$  either (i) by the expiration of its timeout (line 17), or (ii) by receiving a higher round message (line 21).

In case (i) we have  $t_p^e - t_p^s = TO + (n + 2) > \Delta$ , that is the timeout,  $n$  send steps, one input step, and one output step. Thus  $p$  stays in round  $r$  more than  $\Delta$  time.

For case (ii), we calculate the minimum duration of round  $r$  by determining the latest time  $t_p^s$  and the earliest time  $t_p^e$  when  $p$  could have started and ended round  $r$ , respectively (see Figure 1). Let  $q$  be the first process to finish round  $r$  at time  $t_q^e$ . Then the earliest that  $p$  may receive a round  $r + 1$  message is  $t_q^e + 3$  (one input step by  $q$  at the start of round  $r + 1$ , one send step, and one output step by  $p$  to finish round  $r$ ). Hence,  $t_p^e = t_q^e + 3$ .

Let  $t_q^s$  be the time when  $q$  started round  $r$ . Process  $q$  sends a round  $r$  message to  $p$  the latest by  $t_q^s + n + 1$  (if the message to  $p$  is sent in the last send step). By assumption  $r \geq r_0$ , so  $t_q^s$  is after GST. Therefore,  $p$  receives the message at most  $\delta + n + 2$  later ( $\delta$  is the maximum transmission delay in this run and, in the worst case,  $p$  is taking an output step when the message is received, so that in total it takes one output step, one input step, and  $n$  send steps, before the next receive step). After one final output step,  $p$  enters round  $r$ . This happens the latest by  $t_q^s + \delta + 2n + 4$ . Therefore  $t_p^s = t_q^s + \delta + 2n + 4$ .

The minimum duration of round  $r$  at  $p$  is  $t_p^e - t_p^s = (t_q^e + 3) - (t_q^s + \delta + 2n + 4) = (t_q^e - t_q^s) - \delta - 2n - 1$ . To calculate  $t_q^e - t_q^s$ , recall that  $q$  finishes round  $r$  by timeout and not by receiving a higher round message, because by assumption no other process started a round higher than  $r$  before  $q$ . Therefore,  $q$  stays in round  $r$  a total of  $t_q^e - t_q^s = TO + n + 2$ . Substituting  $t_q^e - t_q^s$ , we obtain  $t_p^e - t_p^s = (TO + n + 2) - \delta - 2n - 1 \geq 2\Delta + n + 6 - \delta \geq \Delta$ , which means that  $p$  stays in round  $r$  more than  $\Delta$  time. ■

---

### Algorithm 3 OTR with the failure detector $\diamond\mathcal{P}$ (code of $p$ )

---

```

1: State:
2:    $r_p \leftarrow 1$  /* round number */
3:    $x_p \in V$ 
4:    $decision_p \in V$ 

5: while true do
6:   send  $\langle r_p, x_p \rangle$  to all processes
7:   wait until received values for round  $r_p$  from all processes  $q \notin \diamond\mathcal{P}_p$ 
8:   if number of values received  $> 2n/3$  then
9:      $x_p \leftarrow x$  smallest most often received value
10:  if more than  $2n/3$  values received are equal to  $v$  then
11:     $decision_p \leftarrow v$ 
12:   $r_p \leftarrow r_p + 1$ 

```

---

Since the execution time is proportional to the parameter  $\Delta$  and independent of the effective transmission delay  $\delta$ , the implementation is not swift:

**Theorem 1.** *The round implementation of Algorithm 2 is not swift.*

*Proof:* In case that  $TO < 2\Delta + 2n + 5$ , the algorithm is not live. Therefore we only consider  $TO \geq 2\Delta + 2n + 5$ . Assume by contradiction that the collection of algorithms  $A(\Delta)$  given by Algorithm 2 is swift. Then, there exist  $k, c \in \mathbb{N}$ , such that in every  $\Delta$ -partial synchronous run  $R$  with a good period  $R'$ , there is an  $i_R$  such that, for all instances  $i > i_R$ ,  $\tau_i(R) < k\delta(R') + c$ . For a contradiction, consider  $A(k\delta(R') + c)$ . By Lemma 1, for all instances started after GST, we have  $\tau_i > \Delta = k\delta(R') + c$ . A contradiction. ■

## IV. A FAILURE DETECTOR-BASED ALGORITHM THAT IS SWIFT

We consider now the OTR algorithm expressed with the failure detector  $\diamond\mathcal{P}$  (Algorithm 3). Intuitively it is easy to see that repeated execution of this algorithm is swift. Indeed, some time after GST, the failure detector list contains exactly the faulty processes. At this point, by line 7, all correct processes wait only for messages from correct processes and, since  $f < n/3$ , the condition on line 8 is always true. Note that the failure detector model requires reliable links, contrary to the solution in the previous section.<sup>4</sup> In this section we assume that links are reliable.

Repeated execution of Algorithm 3 is expressed by Algorithm 4. The box in Algorithm 4 corresponds to line 7 of Algorithm 3. For simplicity, we have not shown in Algorithm 4 the (trivial) implementation of  $\diamond\mathcal{P}$ . We assume that both Algorithm 4 and the implementation of  $\diamond\mathcal{P}$  run in the same partial synchronous system in the following way: in every even step Algorithm 4 is executed, in every odd step the implementation of  $\diamond\mathcal{P}$  is executed.

The correctness of Algorithm 4 follows from the following lemma:

<sup>4</sup>Consider two correct processes  $p$  and  $q$  and line 7 executed by  $p$ . If the message sent by  $q$  is lost, and  $p$ 's failure detector never suspects  $q$ , then  $p$  is blocked forever at line 7.

---

**Algorithm 4** Multiple instances of Algorithm 3 (code of  $p$ )

---

```
1: Initialization:
2:    $r_p \leftarrow 1$ 
3:    $\forall i \in \mathbb{N} : x_p[i] \leftarrow \perp$ 
4:    $\forall i \in \mathbb{N} : decision_p[i] \leftarrow \perp$ 

5: while true do
6:    $I \leftarrow input()$ 
7:   for all  $\langle i, v \rangle \in I$  do
8:      $x_p[i] \leftarrow v$ 
9:     send  $\langle r_p, x_p, p \rangle$  to all processes
10:    while not received  $\langle r_p, x_q, q \rangle$  from all processes  $q \notin \diamond\mathcal{P}_p$  do
11:      receive( $M$ )
12:       $Rcv \leftarrow Rcv \cup M$ 
13:     $O \leftarrow \emptyset$ 
14:    for all  $i : x_p[i] \neq \perp$  and  $decision_p[i] = \perp$  do
15:      if number of values received  $\langle r_p, x', - \rangle > 2n/3$  then
16:         $x_p[i] \leftarrow$  smallest most often value  $x'[i]$ 
17:        if more than  $2n/3$  values  $x'[i]$  are equal to  $v$  then
18:           $decision_p[i] \leftarrow v$ 
19:           $O \leftarrow O \cup \{i, v\}$ 
20:    output( $O$ )
21:    $r_p \leftarrow r_p + 1$ 
```

---

**Lemma 2.** For Algorithm 4, there is eventually a round  $GSR$  so that for all rounds  $r \geq GSR$ , every correct process receives a message from every correct process in round  $r$  and receives no message from faulty processes.

*Proof:* By the properties of  $\diamond\mathcal{P}$ , there is a time where the FD is accurate and complete, *i.e.*, a process is suspected if and only if it is faulty. In every round that is started after this time, every correct process waits for a message from every correct process. ■

Theorem 2 proves that Algorithm 4 is swift, by showing that eventually every instance of consensus decides in at most  $3\delta + 6n + 6$ .

**Theorem 2.** For a run of Algorithm 4 with  $n > 3f$  and an infinite number of instances of consensus, there is an instance  $i_0$  such that for all  $i > i_0$ , we have  $\tau_i \leq 3\delta + 6n + 6$ .

*Proof:* Let  $GSR$  be the round defined by Lemma 2. Since in every input step only a finite number of instances are read, there is an input step so that this step and all later input steps are in a round after  $GSR$ . Let  $i_0$  be the largest consensus instance started in a round before  $GSR$  (instance  $i$  is started in the round in which the last process starts instance  $i$ ). Consider an instance  $i > i_0$ . The maximum execution time of  $i$  corresponds to the maximum duration of two rounds. This follows from Lemma 2, which ensures that instance  $i$  decides in at most two rounds. It remains to calculate the maximum time for two rounds after  $GSR$ .

Let  $t$  be the first time a process, say  $p$ , starts round  $r > GSR$ . Since  $r - 1 \geq GSR$ ,  $p$  received round  $r - 1$  messages from all correct processes. This must have happened the latest by time  $t - 2$  in order to allow  $p$  to execute the output

step of round  $r - 1$ , and to enter round  $r$  at time  $t$ .<sup>5</sup> Therefore  $p$  executed the receive step of round  $r - 1$  at latest by time  $t - 4$ , and all correct processes started the send steps for round  $r - 1$  at latest by time  $t - 4$ ; these send steps finished at latest by time  $t - 4 + 2n = t + 2n - 4$ , and messages are received at latest by time  $t + 2n - 4 + \delta$ . Adding the output step, all correct processes started round  $r$  the latest at  $t' = t + 2n - 2 + \delta$ .

By  $t'' = t' + 2n + 2 + \delta$  all round  $r$  messages are thus ready for reception, and received by  $t'' + 2$ . Again by  $t'' + 2 + 2n + 2$  all round  $r + 1$  messages are sent, and thus round  $r + 1$  ends the latest at  $t'' + 2 + 2n + 2 + \delta + 2 = t + 3\delta + 6n + 6$ . ■

*Remark:* Failure detector based solutions require reliable links. This has the following implication. In contrast to partial round implementation of Section III, no round is skipped, *i.e.*, processes send messages for all rounds, and wait for the messages from all unsuspected processes. This implies that, unlike the round implementation in the previous section, it is no more possible to bound the time from  $GST$  until the first decision. To see this, note that at  $GST$ , a process  $p$  might be in a round  $r$  that is arbitrarily smaller than the highest round number  $r_{max}$  at that time. Since other correct processes might wait in any round  $r'$ ,  $r \leq r' \leq r_{max}$ , for the round  $r$  message of process  $p$ ,  $p$  cannot skip the sending step of all rounds between  $r$  and  $r_{max}$ . This takes an unbounded amount of time, as  $r_{max} - r$  can be arbitrarily large. Note that the problem cannot be solved by packing all messages into a single one since, between the sending steps, process  $p$  has to perform receive steps (to receive messages from the other correct processes).

## V. A NEW ROUND IMPLEMENTATION THAT IS SWIFT

We show now that the implementation of the round model can be made swift. Like in the failure detector approach, each process estimates a set of alive processes (the complementary of the set of suspected processes) and uses this set to terminate a round earlier after  $GST$ , namely, as soon as it receives all messages from the alive set. Contrary to the failure detector approach, the algorithm tolerates message loss, by using a timeout which expires only before  $GST$ . Like in the round-based implementation, processes resynchronize after message loss by skipping rounds. Skipping rounds also allows the algorithm to decide in a bounded time after  $GST$ .

### A. Issue to address

Combining the termination of a round upon reception of all messages from alive processes, and the round-skipping mechanism, requires some attention. The problem is illustrated in Figure 2. In this scenario,  $p_3$ 's round  $r$  message is the last message needed by  $p_2$  to have all round  $r$  messages. Let us assume that upon receiving this message,

<sup>5</sup>Note that we have to double the time for a step, since only every second step is of the asynchronous algorithm.

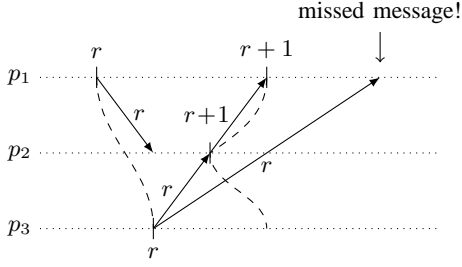


Figure 2. New round implementation: issue to address

$p_2$  immediately sends its round  $r + 1$  message to all. In this case, process  $p_1$  may receive the round  $r + 1$  message of  $p_2$  before the round  $r$  message of  $p_3$ . If  $p_1$  jumps to round  $r + 1$  upon receiving the first round  $r + 1$  message, it will miss  $p_3$ 's round  $r$  message, thereby breaking space uniformity on round  $r$ . This situation may repeat in every round, thus preventing the algorithm from deciding. We show now how we address this problem.

### B. The full algorithm

The ideas described above are used in Algorithm 5, which is a round implementation that is swift. Algorithm 5 enhances Algorithm 2 as follows:

- (i) Each process  $p$  maintains an estimation of the set of alive processes in  $Alive_p$  (see line 13), and updates it every  $TO_A$  steps.  $TO_A$  is thus the timeout used to suspect faulty processes.
- (ii) A process goes directly to the next round if it receives a message from all processes in its alive set (lines 14-15). This is the key point to make the algorithm swift.
- (iii) In any case, a process goes to the next round after  $TO$  time (lines 16-17).  $TO$  is thus the timeout for a round in bad periods.
- (iv) When receiving a round message from the next round for the first time, the process waits for at most  $TO_D$  steps before going into this round (lines 21-22). For this and the last point, each process  $p$  maintains a variable  $timeout_p$ , initially set to  $TO$  (line 8) which is modified when a round  $r + 1$  message is received (line 22). This is used to address the problem described in Section V-A.
- (v) When receiving a message from a round higher than the next round (*i.e.*, larger than  $r_p + 1$ ), the process immediately goes to this round (lines 19-20). This ensures a fast resynchronization of the processes after a bad period.

We now show the correctness of this solution (Section V-C), and that the algorithm is swift (Section V-D).

### C. Correctness

Algorithm 1 together with Algorithm 5 solves repeated consensus in a partial synchronous system. As already

### Algorithm 5 A swift round implementation (code of $p$ )

```

1:  $r_p \leftarrow 1$  /* round number */
2:  $next\_r_p \leftarrow 1$ 
3:  $Rcv_p \leftarrow \emptyset$  /* set of received messages */
4:  $\forall i \in \mathbb{N} : state_p[i] \leftarrow \perp$  /* state for instance  $i$  */

5: while true do
6: input & send /* lines 6-13 of Algorithm 2 */
7:  $i_p \leftarrow 0$ ;
8:  $timeout_p \leftarrow TO$ 
9: while  $next\_r_p = r_p$  do
10:  $i_p \leftarrow i_p + 1$ 
11: receive( $M$ )
12:  $Rcv_p \leftarrow Rcv_p \cup M$ 
13:  $Alive_p \leftarrow \{\text{set of processes from whom}$ 
14:  $\quad \text{there is a message within last } TO_A \text{ steps}\}$ 
15: if  $\forall q \in Alive_p : \exists (M_q, r_p, q) \in Rcv_p$  then
16:  $next\_r_p \leftarrow r_p + 1$ 
17: if  $i_p \geq timeout_p$  then
18:  $next\_r_p \leftarrow r_p + 1$ 
19:  $r \leftarrow \max\{r : \langle -, r, - \rangle \in Rcv_p\}$ 
20: if  $r > r_p + 1$  then
21:  $next\_r_p \leftarrow r$ 
22: if there is a message from round  $r_p + 1$ 
23: for the first time then
24:  $timeout_p \leftarrow \min\{i_p + TO_D, TO\}$ 

25: comp. & output /* lines 22-29 of Algorithm 2 */
26:  $r_p \leftarrow next\_r_p$ 

```

receive

discussed, Algorithm 1 is always safe (with  $n > 3f$ ). Before proving that the round implementation given by Algorithm 5 provides liveness, we show some properties of the algorithm—related to correctness—that hold after  $GST$ .

When the good period starts at  $GST$ , processes will synchronize to the same round using the following two mechanisms: (i) when a process receives a higher round message, it advances rounds either immediately (line 20), or within  $TO_D$  (lines 21-22), or when the original timeout  $TO$  expires; (ii) in any case, processes remain in a round at most  $TO$  time, starting a new round when this timeout expires (lines 16-17 and line 22). Therefore, shortly after  $GST$ , there will be a process  $p$  that starts a new round  $r$  that is higher than any round started by the other alive processes. When the other processes receive the round  $r$  message from  $p$ , they will advance to round  $r$  and send their own messages. These messages are then received by all alive processes, resulting in a space uniform round.

As discussed in Section V-A, a round  $r + 1$  message may be received before all round  $r$  messages (Figure 2). To address this issue, if a process  $p$  in round  $r$  receives a message from round  $r + 1$  for the first time and it has not received all the messages from its alive set, it does not advance immediately. Instead, it waits either for an additional  $TO_D$  or until the end of the original timeout, whichever comes first. During the good period, all the remaining round  $r$  messages will be received before this revised timeout expires. To see why, notice that for a process to send a round  $r + 1$  message, it must have received

all round  $r$  messages from the alive processes, so these messages will also be received by process  $p$  within at most  $TO_D = \Delta + (n - 1)$ , namely  $n - 1$  send steps and  $\Delta$  maximum transmission delay. In any case, all messages will be received before the original round timeout, so the process only has to wait for the minimum of  $TO_D$  or what is left of  $TO$ .

If a process  $p$  in round  $r$  receives a message from round  $r + 2$  or higher, it can conclude that the good period has not yet been started, so  $p$  advances immediately to round  $r + 2$ . This holds for the following reason. Assume that the system is in a good period, and let some process  $q$  send a round  $r + 2$  messages; then either (i)  $q$  received all round  $r + 1$  messages, including  $p$ 's message, which is not possible; or (ii) the timeout for round  $r + 1$  expires, which is not possible as the timeout is chosen in a way that processes have enough time to receive all round messages and messages are not lost in the good period. This shows a contradiction: the system cannot be in a good period.

Thus we can show:

**Theorem 3.** *Consider a run of Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ , and  $TO_A \geq TO + \Delta + (2n + 1)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Then every consensus instance that starts at  $t$  decides the latest at  $\max(t, GST) + TO_A + 2TO + TO_D + 3\Delta + (6n + 15)$ .*

The proof is based on the following two lemmas (for the proof see Appendix B). The first establishes that eventually rounds are space uniform (see Sect. III):

**Lemma 3** (Timeouts  $TO$  and  $TO_D$ ). *Consider Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ . Let  $R$  be a  $\Delta$ -partial synchronous run, and  $t_r$  the time the first process starts a new round  $r$  after  $GST$ , such that all processes have the same *Alive* set after  $t_r$ . Then round  $r$  is space-uniform.*

The previous lemma requires all processes to have the same *Alive* set. Lemma 4 shows that this becomes true shortly after  $GST$ .

**Lemma 4** (Timeout  $TO_A$ ). *Consider Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ , and  $TO_A \geq TO + \Delta + (2n + 1)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Let  $t_r$  be the time the first process starts a new round  $r$  after  $GST$ . Then by time  $t_r + 2 + TO_A$  all processes have the same *Alive* set.*

#### D. Swiftiness

In order to show that Algorithm 1 together with the round implementation provided by Algorithm 5 is swift, we show that the execution time of a consensus instance depends only on  $\delta$  and not on  $\Delta$ .

The main properties of the algorithm related to the swiftiness, which hold after  $GST$ , are the following. First, the

*Alive* set becomes accurate the latest by  $GST + TO + TO_A + n + 4$  (line 13). This follows from Lemma 4, with  $t_r$  being at latest  $GST + TO + n + 2$ . Then, once the *Alive* set is accurate after  $GST$ , it no more changes and therefore no further timeout expires. Finally, all processes finish rounds as soon as all messages from alive processes are received and advance round by lines 14-15, rendering the algorithm swift.

**Theorem 4.** *Consider Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ , and  $TO_A \geq TO + \Delta + (2n + 1)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Then every consensus instance that is started after  $GST + X$  with  $X = TO_A + 3(TO + n + 2) + 2$ , has an execution time of  $\tau_i \leq 3\delta + 3n + 5$ .*

*Proof:* Let  $i_0$  be a new consensus instance started at time  $t_s > GST + X$ . Such an instance exists, because the input queue contains an infinite number of elements and each input step reads a finite number of instances. Let  $p$  be the process that started instance  $i_0$  (last process doing an input step for  $i_0$ ) and  $r_i$  the round where it was started.

We will first show that rounds  $r \geq r_i - 1$  are space uniform. By lines 8, 16 and 22 of Algorithm 5, processes remain on a round for at most  $TO + n + 2$  time. Therefore, the first round started after  $GST$  starts the latest at  $t_0 = GST + TO + n + 2$ . By Lemma 4, the latest at  $t_1 = t_0 + TO_A + 2$  all processes have the same *Alive* set, and by Lemma 3, all rounds started after  $t_1$  are space uniform. Therefore the first space uniform round,  $r'$ , starts the latest at time  $t_2 = t_1 + TO + n + 2$ , and  $r' + 1$  the latest by  $t_3 = t_1 + 2(TO + n + 2)$ . Expanding this expression, we obtain  $t_3 = GST + TO_A + 3(TO + n + 2) + 2 = GST + X$ . Since  $r_i$  started at time  $t_s \geq GST + X$ , rounds  $r \geq r_i - 1$  are space uniform.

Using Theorem 3 we can conclude that instance  $i_0$  is decided by round  $r_i + 1$ . We are now ready to compute the maximum execution time of  $i_0$ . By definition, we have  $\tau_i = t_e - t_s$ , where  $t_e$  is the time when the last process performs an output step for round  $i_0$  (and  $t_s$  is previously defined). To determine the upper bound on  $\tau_i$ , we'll compute the smallest and the largest values for times  $t_s$  and  $t_e$ , respectively. Since by assumption  $t_s$  happens in round  $r_i$ , then  $t_s$  is smallest if  $p$  is the first process starting the round. The largest value for  $t_e$  is the time of the output step of the last process finishing round  $r_i + 1$ . Next we compute  $t_e$ .

Since round  $r_i - 1$  is space uniform, process  $p$  received all round  $r_i - 1$  messages before advancing to round  $r_i$ , hence the latest by  $t_s - 2$  all alive processes had sent their round  $r_i - 1$  message to  $p$ . By  $t_s + n - 2$  all round  $r_i - 1$  messages were sent, and  $\delta$  time later received. Thus, by  $t_s + \delta + 2n$  all processes entered round  $r_i$  and finished sending all messages.  $\delta$  time later all round  $r_i$  messages are received and by time  $t_s + 2\delta + 2n + 1$  all processes started round  $r_i + 1$ . By a similar reasoning, by time  $t_s + 3\delta + 3n + 3$ , all

processes finished round  $r_i + 1$ . Hence, instance  $i_0$  ends at time  $t_e = t_s + 3\delta + 3n + 3$ , and we have  $\tau_i = 3\delta + 3n + 3$ . ■

## VI. EXPERIMENTAL RESULTS

In this section we present the results of an experimental study, comparing the three algorithms presented previously. The main questions we want to answer are (i) how much improvement can be obtained in a round-based algorithm using a swift round implementation, and (ii) are swift round implementations competitive with implementations that use failure detectors.

*Experimental setup:* We performed our experiments both on an emulated network and directly on a physical network (a cluster). The emulated network allowed us to test the behavior of the algorithms with different transmission delays and message loss rates, while the physical network shows what to expect on a cluster environment.

In all experiments, processes were started with 1 second of delay between each other. This prevents initial synchronization and exercises the ability of the algorithms to resynchronize the processes.

The metric considered is the decision time for each consensus instance. Processes run each instance sequentially, starting the next one either when they decide, or when they learn the decision by receiving a message from a higher instance. Each data point shown on the plots below was obtained from a 10 minutes run. We then calculated the average decision time, ignoring the first 10% of the run. For each data point, we show the 95% confidence intervals.

*Implementing  $\diamond\mathcal{P}$  and reliable channels for the failure detector algorithm:* We implemented  $\diamond\mathcal{P}$  by having each process send heartbeats to all every  $\eta$  time. A process  $p$  suspects  $q$  if it does not receive any heartbeat for more than  $\tau$  time. We also implemented reliable channels using message acknowledgments and retransmission. We decided not to use TCP, because our initial experiments using TCP resulted in very poor performance under high message loss conditions. TCP is designed to interpret message loss as an indication of congestion, and therefore it reacts by increasing the retransmission time. On a typical TCP implementation, the interval between retransmissions may reach several minutes, which in practice forces the algorithms running on top of it to stop.

*Notation:* In the following,  $\delta_{net}$  denotes the one-way transmission delay of the physical network,  $\delta_{emu}$  the delay emulated by ModelNet, and  $\delta_{eff}$  the effective one-way transmission delay between two processes. On the experiments run directly on the physical network,  $\delta_{eff} = \delta_{net}$ . However, when using ModelNet,  $\delta_{eff} = 2\delta_{net} + \delta_{emu}$ , since each packet is transmitted two times on the physical network (see Section VI-A). Finally, note that contrary to  $\delta$  defined previously in the paper,  $\delta_{net}$  is not a bound. Instead, it is

a random variable, reflecting the non-deterministic behavior of a physical network.

In the following, NS-OTR, S-OTR, and FD-OTR denote respectively the non-swift OTR (Algorithm 1 + Algorithm 2), the swift OTR (Algorithm 1 + Algorithm 5), and OTR with FD (Algorithm 4 +  $\diamond\mathcal{P}$ ).

### A. Emulated network

We used ModelNet [12] to emulate a network. ModelNet uses two types of nodes: a *core node* that applies the traffic policies, and one or more *edge nodes* that run the application being tested. The edge nodes redirect all traffic sent by the processes to the core node, which applies the traffic policy (e.g., delay, loss and maximum bandwidth) and then transmits the packet to the intended receiver. We varied the emulated delay and loss rate, while leaving the emulated bandwidth set to 1Gbps. We used two physical machines for all experiments run on ModelNet. All 4 replicas were running on a dual Pentium 4 at 3.6GHz with 1GB RAM, while the core node was a Pentium Pro at 200MHz with 70MB of RAM. The machines were connected by a full duplex 100Mbits Ethernet, and had a ping time of approximately 0.3ms. Hence,  $\delta_{eff} \approx 0.3 + \delta_{emu}$ .

1) *Varying the timeout:* In the first set of experiments, we fixed the emulated transmission delay while varying the timeout  $TO$  used by the algorithms. Figure 3 shows the results for  $\delta_{emu} = 0ms$  and Figure 4 for  $\delta_{emu} = 40ms$ . The  $x$  scale indicates the timeout  $TO$  used by the algorithms to terminate a round.<sup>6</sup> For the tests with  $\delta_{emu} = 40ms$ , the failure detector was configured with  $\eta = TO/2$  and  $\tau = TO$ . The rationale is that  $TO$  is the time an algorithm should wait before declaring a failure and taking corrective measures, e.g., advancing rounds or suspecting a process. With  $\delta_{emu} = 0ms$ , following the same policy would result in the network being overloaded with heartbeats, so we opted for  $\eta = TO$  and  $\tau = 2TO$ .

The results clearly validate the main motivation behind this work, in that S-OTR performs at the speed of the network, being independent from the timeout.

With  $\delta_{emu} = 0ms$  (Figure 3-left), FD-OTR performs poorly with low timeouts. This is caused by the additional messages sent by the failure detector and the reliable channels implementation, which slow down the processes and congest the network. For higher timeouts, this overhead becomes less significant and the algorithm starts performing similarly to the other implementations. When looking only at NS-OTR vs S-OTR (Figure 3-right), it is clear that the decision time of NS-OTR increases linearly with the timeout, while S-OTR is constant. Furthermore, even with the optimal timeout of 2ms, NS-OTR performs worse than S-OTR, because no fixed timeout can approximate perfectly the time that it takes for a process to receive all messages (it fluctuates from round to round).

<sup>6</sup>Equivalent to  $2\Delta$  on the NS-OTR and  $3\Delta$  for the S-OTR.



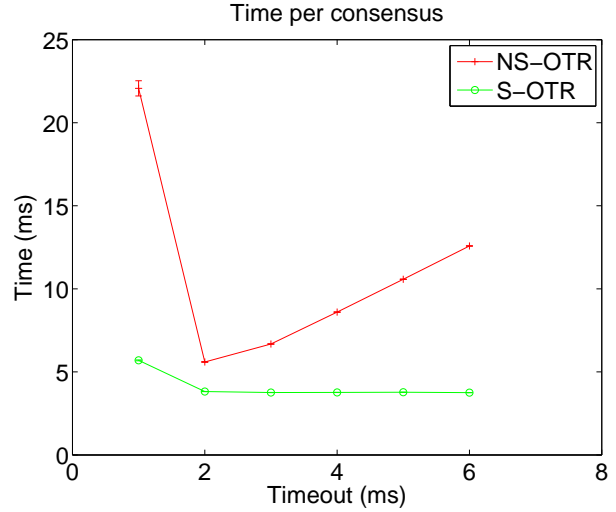
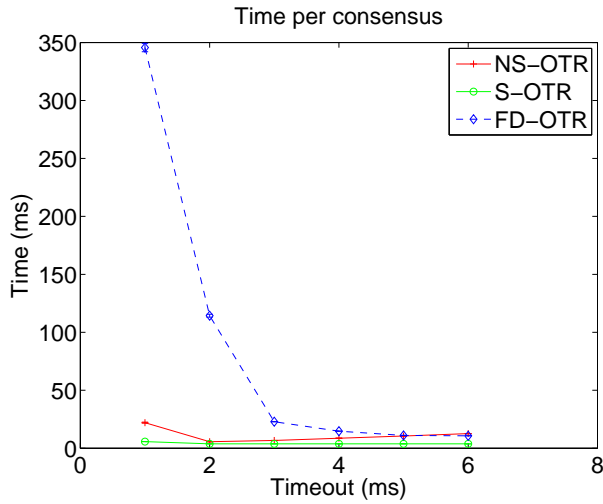


Figure 3. Performance on ModelNet with  $\delta_{eff} \approx 0.3ms$  ( $\delta_{emu} = 0$ ,  $2\delta_{net} \approx 0.3$ ). The figure on the right repeats NS-OTR and S-OTR from the left, with a different time scale.

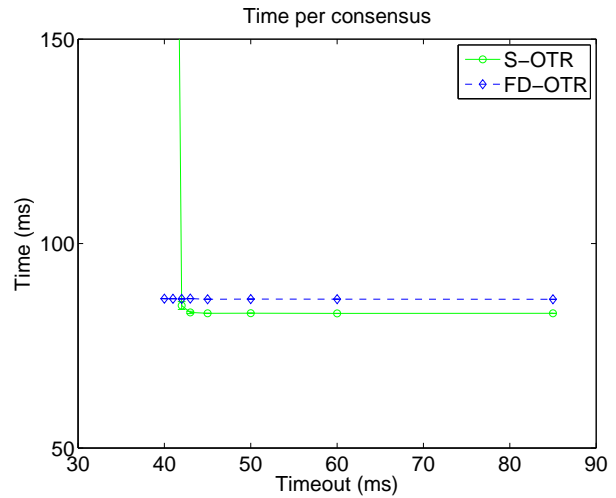
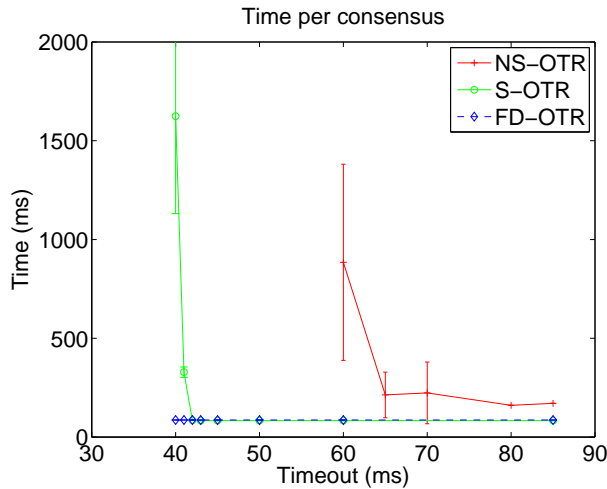


Figure 4. Performance on ModelNet with  $\delta_{eff} \approx 40.3ms$  ( $\delta_{emu} = 40$ ,  $2\delta_{net} \approx 0.3$ ). The figure on the right repeats S-OTR and FD-OTR from the left, with a different time scale.

With  $\delta_{emu} = 40ms$  (Figure 4-left), NS-OTR performs poorly with timeouts lower than  $80ms$ . For timeouts lower than  $60ms$ , the algorithm took hundreds of rounds for each decision, so we did not show the results as they were not statistically significant. Notice that  $80ms \approx 2\delta_{eff}$ , which matches the results from the analytical analysis, where a round must last  $TO = 2\Delta$  in order to ensure decision. The swift version S-OTR is more tolerant to a non-optimal timeout, being able to synchronize even with timeouts slightly above  $40ms$ . This is because processes finish rounds early, after receiving all messages, allowing the processes that are behind to slowly catch-up with the ones in the lead.

FD-OTR is also independent of the timeout, producing the optimal performance regardless of the values used for

the underlying failure detector. Recall that in the absence of message loss, the values chosen for the failure detector (*i.e.*,  $\tau = 2\eta$ ) prevent false suspicions, so FD-OTR can proceed at the speed of the network. The overhead of the implementation of failure detectors and reliable channels is less in this scenario, as shown in Figure 4-right, where FD-OTR performs only slightly worse than S-OTR.

2) *Message loss*: Figure 5 shows the behavior of the algorithms in networks with message loss. The experiment was run on ModelNet with  $\delta_{emu} = 0$ . Both the swift and the non-swift versions were configured with a timeout of  $10ms$ . The failure detector was configured with  $\eta = 10ms$  and  $\tau = 25ms$ , so that it tolerates 2 or 3 lost heartbeats before (wrongly) suspecting a process. The reliable channels

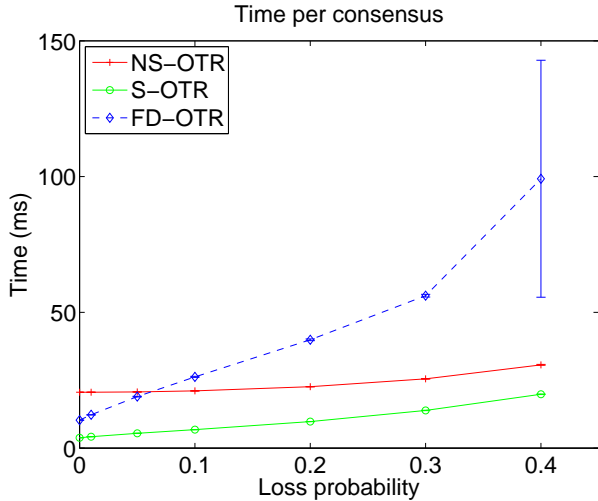


Figure 5. Performance with message loss:  $\delta_{eff} \approx 0.3ms$  ( $2\delta_{net} \approx 0.3$ ,  $\delta_{emu} = 0$ )

implementation retransmits a message every  $25ms$ .

Both NS-OTR and S-OTR are very resilient to message loss. Even with 40% messages loss, the average decision time is only a few milliseconds more than with no message loss. This is because the algorithms make progress as soon as a single process receives three messages ( $2n/3$ ), *i.e.*, two messages from other processes since its own message is always delivered. FT-OTR performs worse because it waits for messages from all processes that are not suspected, so that a single message loss in a round is enough to delay progress (suspecting a process requires more than a single message loss).

S-OTR outperforms both NS-OTR and FD-OTR in the presence of message loss. In particular, the performance of FD-OTR degrades significantly with message loss, caused by the overhead of the retransmissions to simulate reliable links.

### B. Physical network (cluster)

For the tests with the physical network, we used a cluster of Dual Pentium 4 at 3.00GHz with 1GB memory connected by a 1Gbit Ethernet. Each process run on a separate node and the ping time between two nodes was between 0.1 and 0.2ms. The failure detector was configured with  $\eta = TO$  and  $\tau = 2TO$ .

Figure 6 shows that on the cluster even a timeout of  $1ms$  is enough for OTR to terminate. S-OTR always outperforms the two other algorithms. Compared to NS-OTR, even with a  $1ms$  timeout, S-OTR performs better. Lowering the timeout of NS-OTR may improve its performance, but with such small timeouts the algorithm becomes sensible to the normal variability of the system, which is caused by non-deterministic factors like OS scheduling and background activity, either on the hosts or on the network. This will cause

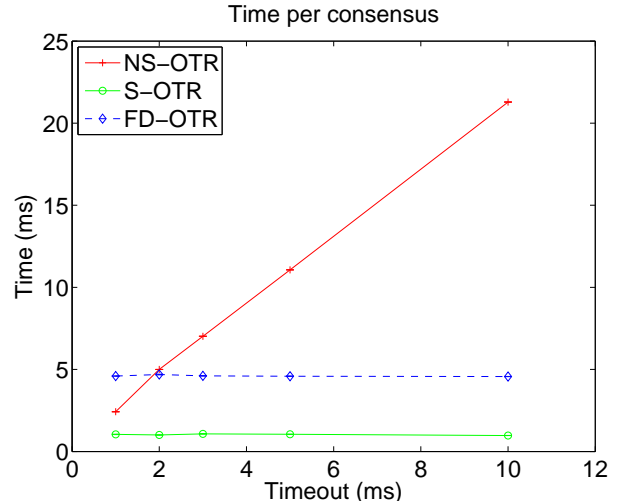


Figure 6. Performance on cluster ( $\delta_{eff} \approx 0.1ms$ )

rounds to finish without receiving all required messages, leading to unstable performance. The timeout of S-OTR can be set to a conservative value, making the algorithm immune to non-deterministic factors, while still providing optimal performance.

FD-OTR suffers again from the overhead of the implementation of failure detectors and reliable channels, resulting in a performance worst than S-OTR.

## VII. DISCUSSION

Table I summarizes the results of the paper. We have analyzed the efficiency of algorithms for solving repeated consensus in two models: the round-based model (which can be implemented on top of a partially synchronous system), and the asynchronous system augmented with failure detectors. Efficiency refers here to *swiftness*, a new notion that captures the fact that an algorithm, once the system has stabilized, progresses at the speed of the messages. Our new round-based implementation combines the advantages of failure detector solutions (swiftness) and round-based model (lossy links). This weak link assumption makes round-based algorithm easy to adapt to the crash-recovery model with stable storage [11].

We have illustrated the new round-based implementation on a specific consensus algorithm (OTR). This does not

	Classical round-based [11] (Algorithm 2)	FD-based [8] (Algorithm 4)	New round-based (Algorithm 5)
Link	lossy	reliable	lossy
Exec. time	$4\Delta + \delta + O(1)$	$3\delta + O(1)$	$3\delta + O(1)$
Swift	no	yes	yes

Table I  
REPEATED CONSENSUS: ALGORITHMS ANALYZED IN THE PAPER

mean that the new solution is limited to OTR. It applies to any consensus algorithm expressed in the round model, in particular to the *LastVoting* algorithm [10], a round-based variant of Paxos [13] that requires only  $n > 2f$ .

#### REFERENCES

- [1] M. J. Fischer, N. A. Lynch, and M. S. Paterson, "Impossibility of distributed consensus with one faulty process," *Journal of the ACM*, vol. 32, no. 2, pp. 374–382, Apr. 1985.
- [2] C. Delporte-Gallet, S. Devismes, H. Fauconnier, F. Petit, and S. Toueg, "With finite memory consensus is easier than reliable broadcast," in *OPODIS*, 2008, pp. 41–57.
- [3] M. Hutle and J. Widder, "On the possibility and the impossibility of message-driven self-stabilizing failure detection," in *Self-Stabilizing Systems*, 2005, pp. 153–170, appeared also as Brief Announcement at PODC'05.
- [4] M. Biely and J. Widder, "Optimal message-driven implementations of omega with mute processes," *ACM Trans. Auton. Adapt. Syst.*, vol. 4, no. 1, pp. 1–22, 2009.
- [5] L. Lamport, "Fast paxos," *Distributed Computing*, vol. 19, no. 2, pp. 79–103, 2006.
- [6] N. Lynch, *Distributed Algorithms*. Morgan Kaufman, 1996.
- [7] C. Dwork, N. Lynch, and L. Stockmeyer, "Consensus in the presence of partial synchrony," *Journal of the ACM*, vol. 35, no. 2, pp. 288–323, Apr. 1988.
- [8] E. Gafni, "Round-by-round fault detectors (extended abstract): unifying synchrony and asynchrony," in *Proceeding of the 16th Annual ACM Symposium on Principles of Distributed Computing (PODC'98)*. Puerto Vallarta, Mexico: ACM Press, 1998, pp. 143–152.
- [9] T. D. Chandra and S. Toueg, "Unreliable failure detectors for reliable distributed systems," *Journal of the ACM*, vol. 43, no. 2, pp. 225–267, Mar. 1996.
- [10] B. Charron-Bost and A. Schiper, "The heard-of model: computing in distributed systems with benign faults," *Distributed Computing*, vol. 22, no. 1, pp. 49–71, 2009.
- [11] M. Hutle and A. Schiper, "Communication predicates: A high-level abstraction for coping with transient and dynamic faults," in *Dependable Systems and Networks (DSN 2007)*. IEEE, Jun. 2007, pp. 92–10.
- [12] A. Vahdat, K. Yocum, K. Walsh, P. Mahadevan, D. Kostić, J. Chase, and D. Becker, "Scalability and accuracy in a large-scale network emulator," *SIGOPS Oper. Syst. Rev.*, vol. 36, no. SI, pp. 271–284, 2002.
- [13] L. Lamport, "The part-time parliament," *ACM Transactions on Computer Systems*, vol. 16, no. 2, pp. 133–169, May 1998.

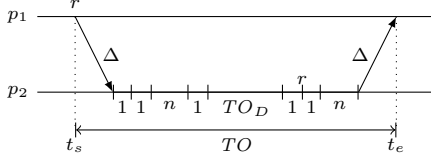


Figure 9. Timeout  $TO \geq TO_D + 2\Delta + (2n + 5)$

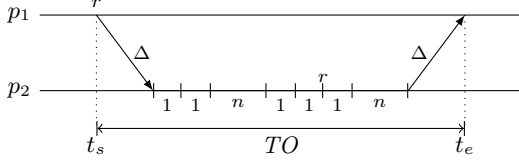


Figure 7. Timeout  $TO \geq 2\Delta + (2n + 5)$

## APPENDIX

### A. Proofs for Section III

**Theorem 5.** Consider a run of Algorithm 2 with  $TO \geq 2\Delta + (2n + 5)$  and  $n > 3f$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Then every consensus instance that starts at  $t$  decides the latest at  $\max(GST, t) + 3TO + \Delta + 4n + 8$ .

Following lemmas together with the results of [10] proves the theorem.

**Lemma 5.** Consider Algorithm 2 with  $TO \geq 2\Delta + (2n + 5)$  and  $n > 3f$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Let  $t_r$  be the time the first process starts a new round  $r$  after  $GST$ . Then round  $r$  is space-uniform.

*Proof:* (See Figure 7 for illustration.) Let  $p$  be the first process to finish the input and send steps for round  $r$ , at time  $t_s$  ( $t_s \leq t_r + (n + 1)$ ). We show that (i) all round  $r$  messages from all alive processes are ready for reception<sup>7</sup> by time  $t_s + TO$ , and (ii) no process expires its round  $r$  timeout before  $t_s + TO$ . This implies that round  $r$  is space-uniform.

(i) By time  $t_s + \Delta$  the round  $r$  message from  $p$  is ready for reception at all processes. Every process  $q$  will make a receive step at most  $(n + 2)$  time later (if at time  $t_s + \Delta$   $q$  was on a output step of a round  $r' < r - 1$ , then it must make one input step and  $n$  send steps before the next receive step). After receiving the round  $r$  message, every process performs an output step for its current round, advances to round  $r$ , performs one input and  $n$  send steps. Therefore, by time  $t_s + \Delta + (2n + 5)$ , all processes have finished sending their round  $r$  messages, and  $\Delta$  time later, by time,  $t_s + 2\Delta + (2n + 5) = t_s + TO$ , all round  $r$  messages are ready for reception at all alive processes. Note that this time is still in the good period, since  $t_s + TO = t_r + TO + (n + 1)$ .

(ii) Since all processes start the timeout for round  $r$  after  $p$ , the timeout of no process will expires before  $t_s + TO$ . Additionally, no process advances to a higher round by receiving a higher round message because for a new round

to start, some the timeout of round  $r$  of some process has to expire. ■

**Lemma 6.** Consider Algorithm 2. Let  $R$  be a  $\Delta$ -partial synchronous run. Then by time  $GST + TO + (n + 2)$  at least one process has started a new round  $r_0$ .

<sup>7</sup>We call a message ready for reception if it must be received with the next receive step of the receiver process.

*Proof:* Let  $p$  be the process with the highest round number  $r$  among all processes. Then the lemma is fulfilled, if  $p$  is at least in round  $r + 1$  by the given time. However, in a good period,  $p$  can be in round  $r$  at most for  $TO + (n + 2)$  time, the timeout and the time for an input, an output, and  $n$  send steps. ■

**Lemma 7.** Consider Algorithm 2 with  $TO \geq 2\Delta + (2n + 5)$ ,  $n > 3f$ , and a  $\Delta$ -partial synchronous run  $R$ . Let  $r_0$  be the first new round that is started after  $GST$ . Then for all instances  $i$  started in a round  $r \geq r_0$ , we have an execution time  $\tau_i \leq 2TO + \delta + (3n + 6)$ .

*Proof:* (See Figure 8 for illustration.) Let  $i$  be an instance started in a round  $r \geq r_0$  by a process  $p$ . Recall that Algorithm 2 needs at most two space-uniform rounds to decide. Since by Lemma 5 rounds  $r$  and  $r + 1$  are space-uniform, all processes decide instance  $i$  by round  $r + 1$  (i.e., they output  $(i, x)$  at line 29, where  $x$  is the decision).

It remains to calculate the maximum time for rounds  $r$  and  $r + 1$ . Let  $p$  be the first process to start round  $r$  at time  $t_r$ . Process  $p$  will finish round  $r$  the latest at  $t_r + TO + (n + 2)$ , and start the send steps for round  $r + 1$ , 1 step later. By time  $t_r + TO + \delta + (2n + 3)$ ,  $p$ 's round  $r + 1$  messages are ready for reception at all processes. At this point, all processes have finished executing the send steps for round  $r$  (process  $p$ 's round  $r$  messages forced them to advance) and are either executing receive steps for round  $r$  or have entered round  $r + 1$ . Therefore, all processes will enter round  $r + 1$  at most 1 step after receiving  $p$ 's round  $r + 1$  message. Round  $r + 1$  will take at most  $TO + (n + 2)$  time, so by time  $t_r + 2TO + \delta + (3n + 6)$  all processes have finished round  $r + 1$ . ■

### B. Proofs for Section V

**Lemma 3.** Consider Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Let  $t_r$  be the time the first process starts a new round  $r$  after  $GST$ . Let all processes have the same Alive set at after  $t_r$ . Then round  $r$  is space-uniform.

*Proof:* Let  $p_1$  be the first process to finish sending its round  $r$  messages at time  $t_s = t_r + (n + 1)$ , and starting the timeout for round  $r$  (see Figure 9). These messages are ready for reception at most  $\Delta$  time later, at  $t_s + \Delta$ . These messages

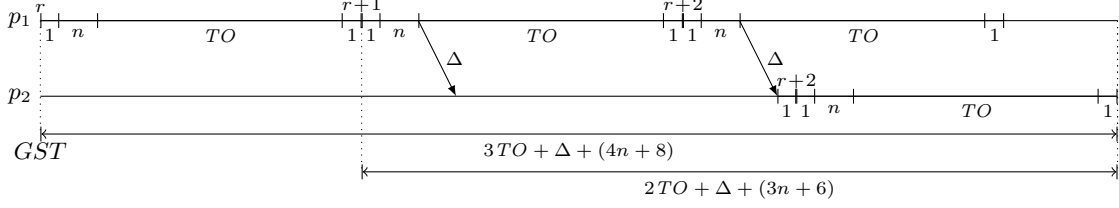


Figure 8. Illustration for Theorem 5 and Lemma 7

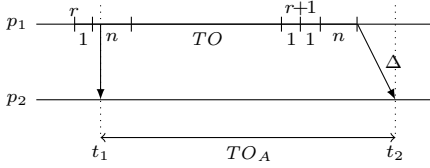


Figure 10. Timeout  $TO_A \geq TO + \Delta + (2n + 1)$

are received in the next receive step, which occurs the latest after  $(n + 2)$  steps (an output step followed by an input step, and  $n$  send steps). This is because some process ( $p_2$  in Figure 9) might be just started executing an output step for some round  $r' < r$ . Therefore,  $p_1$ 's message is received by all processes the latest at time  $t_1 = t_s + \Delta + (n + 3)$ . Any process that receives this message in round  $r - 1$  for the first time, might set its timeout to  $t_1 + TO_D < TO$  (see lines 21-22). And start round  $r$  the latest by time  $t_1 + TO_D + 1$ , after an output step for round  $r - 1$ . By time  $t_2 = t_1 + TO_D + 1 + 1 + n$ , any process (including  $p_2$ ) has performed an input step and  $n$  send steps for round  $r$ . This message is ready for reception the latest at time  $t_e = t_2 + \Delta = t_s + TO_D + 2\Delta + (2n + 5)$ . The timeout  $TO = TO_D + 2\Delta + (2n + 5)$  ensures that no timeout started at time  $t_s$  expires before  $t_e$  (see line 16). So when the timeout expires, all messages for round  $r$  are either received or ready to be received. Before, calling the transition function for round  $r$  (in line 23), a receive step is performed (in line 11); thus every process in round  $r$  receives a message from every process, and round  $r$  is space uniform.

Note that no process in round  $r$  can receive a message from round  $> r + 1$ . We prove this by contradiction. Let  $p$  be a process in round  $r$  that receives a message from round  $r + 2$ . This means that there is some process  $q$  that sent round  $r + 2$  messages. This requires that either (i)  $q$  receives all round  $r + 1$  messages, including  $p$ 's message, which is not possible; or (ii) the timeout for round  $r + 1$  expires, which is not possible inside the given interval.

If a process ends round  $r$  at time  $t$  before the end of timeout  $TO$ , because it has received all round  $r$  messages from its alive set (line 15), any other process does so the latest by time  $t + (n - 1) + \Delta$ . From lines 21-22, a process in round  $r$  that receives a message from round  $r + 1$  for the first time, waits until  $t + TO_D$  time before starting round

$r + 1$ , which is enough to receive all round  $r$  messages. By the assumption, since all processes have the same actual alive set in the given interval, round  $r$  is also space uniform in this case. ■

**Lemma 4.** Consider Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ , and  $TO_A \geq TO + \Delta + (2n + 1)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Let  $t_r$  be the time the first process starts a new round  $r$  after GST. Then by time  $t_r + 2 + TO_A$  all processes have the same Alive set.

*Proof:* By time  $t_1 = t_r + 2$  the first round  $r$  message can be received (see Figure 10). From the code of the algorithm, every process starts a new round the latest every  $TO + (n + 2)$  steps: one input step followed by  $n$  send steps,  $TO$  receive steps followed by an output step. From the fact that a process sends at most one message in each step, every process  $p_1$  sends messages to any process  $p_2$  every  $TO + (n + 2) + (n - 1)$  steps. Since a message can take at least 0 and at most  $\Delta$  time to be received, every process receives a message every  $x = TO + (2n + 1) + \Delta$  time. From the code of the algorithm, process  $p_2$  excludes process  $p_1$  from its alive set, if it does not receive a message within  $TO_A$  steps (see line 13). Comparing  $TO_A$  with  $x$  we have  $TO_A = x$ , which is sufficient to receive a message from any alive process. ■

**Theorem 3.** Consider a run of Algorithm 5 with  $n > 3f$  and the following timeouts:  $TO_D \geq \Delta + (n - 1)$ ,  $TO \geq TO_D + 2\Delta + (2n + 5)$ , and  $TO_A \geq TO + \Delta + (2n + 1)$ . Let  $R$  be a  $\Delta$ -partial synchronous run. Then every consensus instance that starts at  $t$  decides the latest at  $\max(GST, t) + TO_A + 2TO + TO_D + 3\Delta + 6n + 15$ .

*Proof:* See Figure 11 for illustration. We distinguish two cases (1)  $t < GST$ , (2)  $t \geq GST$ . In case (1) by Lemma 3 a new round is started after GST the latest by time  $GST + TO + n + 2$ . From Lemma 4 all processes have the same alive set by time  $t_0 = GST + TO + TO_A + n + 4$ . In case (2) by Lemma 4 all processes have the same alive set by time  $t_0 = t + TO_A + 2$ , which is strictly smaller than  $GST + TO + TO_A + n + 4$ . From the code of the algorithm a process, e.g.,  $p_1$ , starts a new round  $r$  every  $TO + (n + 2)$  steps, i.e., the latest by time  $t_1 = t_0 + TO + n + 2$ . All processes do so by time  $t_2 = t_1 + TO_D + \Delta + (2n + 5)$ . This means that all processes start round  $r$  with the same alive

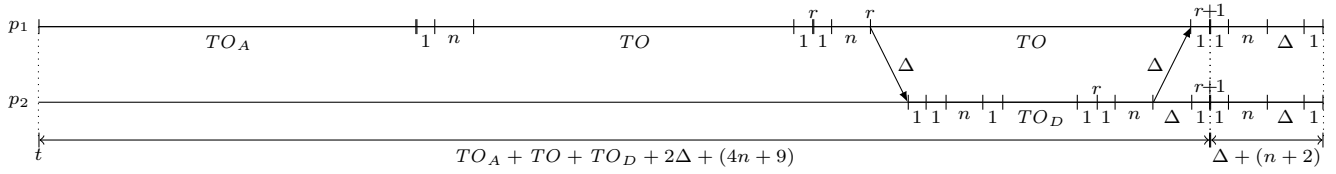


Figure 11. Illustration for Theorem 3

set. From Lemma 3, round  $r$  is space uniform. Furthermore, all processes receive all round  $r$  messages from their alive set, and end round  $r$  the latest by time  $t_3 = t_2 + \Delta + (n+2)$  and start round  $r + 1$  at this time.

From the assumption, no process crashes after  $GST$ , therefore, the alive set remains the same. In round  $r + 1$ , all processes send their messages to all the latest by time

$t_3 + (n + 1)$ . These messages can be received by all processes the latest by time  $t_3 + (n + 1) + \Delta$ . From lines 14-15, all processes end round  $r + 1$  the latest by time  $t_3 + (n + 1) + \Delta + 1$  after an output step. This means that all processes decide the latest by this time which is equal to  $t_0 + TO + (TO_D + 2\Delta + (4n + 9)) + (\Delta + (n + 2))$ , or  $\max(GST, t) + TO_A + 2TO + TO_D + 3\Delta + (6n + 15)$ . ■