

# Relaxed Atomic Broadcast: State-Machine Replication Using Bounded Memory

Omid Shahmirzadi, Sergio Mena, André Schiper  
École Polytechnique Fédérale de Lausanne (EPFL)  
CH-1015 Lausanne, Switzerland  
Email: {first.last}@epfl.ch

**Abstract**—Atomic broadcast is a useful abstraction for implementing fault-tolerant distributed applications such as state-machine replication. Although a number of algorithms solving atomic broadcast have been published, the problem of bounding the memory used by these algorithms has not been given the attention it deserves. It is indeed impossible to solve repeated atomic broadcast with bounded memory in a system (non-synchronous or not equipped with a perfect failure detector) in which consensus is solvable with bounded memory. The intuition behind this impossibility is the inability to safely garbage-collect unacknowledged messages, since a sender process cannot tell whether the destination process has crashed or is just slow.

The usual technique to cope with this problem is to introduce a membership service, allowing the exclusion of a slow or silent process from the group and safely discarding unacknowledged messages sent to this process. In this paper, we present a novel solution that does not rely on a membership service. We *relax* the specification of atomic broadcast so that it can be implemented with bounded memory, while being strong enough to still be useful for applications that use atomic broadcast, e.g., state-machine replication.

## I. INTRODUCTION AND RELATED WORK

Atomic broadcast has been proposed as the key abstraction to implement fault-tolerant distributed services [?] using the state-machine approach [?]. A number of different implementations of atomic broadcast have been proposed in the literature for a variety of system models [?]. However, they rarely tackle the problem of bounding the use of memory. The fact that an algorithm needs a potentially unbounded amount of buffers is often considered as a minor (implementation) issue. Bounding memory might not be a very exciting theoretical issue, it is nevertheless important from a practical point of view, since inability to bound (or garbage-collect) the memory used may lead to serious instability of the application, with effects similar to those of memory leaks. This is definitely not the best feature for algorithms that are supposed to increase availability. As Parnas argues in [?], a model should be simple, but if it becomes too simple it risks being a lie, i.e., not representing reality. No real system can assume it has access to unbounded memory.

Implementing atomic broadcast with bounded memory in a synchronous system is trivial [?]. However, if the system model does not allow us to distinguish a slow process from a crashed process, the ability of atomic broadcast algorithms to bound their memory – without

affecting correctness – becomes challenging. Ricciardi [?] proved that a primitive as basic as (repeated) reliable broadcast cannot be implemented in a system with message losses in which slow processes are indistinguishable from crashed processes. Trivially, Ricciardi’s impossibility result also applies to (repeated) atomic broadcast, since it is strictly stronger than (repeated) reliable broadcast. In this paper, we address the problem of bounded memory in the context of repeated atomic broadcast by weakening the specification of atomic broadcast. Note that (one instance of) consensus has been shown to be solvable with bounded memory [?] in an asynchronous system with the  $\diamond S$  failure detector, and in [?] Delporte-Gallet *et al.* show that solving (repeated) reliable broadcast requires indeed a stronger failure detector than solving (one instance of) consensus.

Group communication prototypes built in the last 20 years have addressed the problem of bounding memory thanks to group membership [?], [?], [?], [?]: slow or irresponsive processes are excluded from the group so that messages sent to them can be safely garbage-collected before buffers at other processes overflow. However, this solution has its own drawbacks. First, the dynamic group model is more complex than the static one. Second, the dynamic model requires the introduction of a group membership service, which adds a performance overhead. Finally, excluding a destination process just because the sender is unable to garbage-collect its output buffers<sup>1</sup> may not always be desirable.

The paper presents *relaxed atomic broadcast*, a novel broadcast primitive defined in the static group model (i.e., no membership service), whose repeated invocation can be implemented using bounded memory. Relaxed atomic broadcast is weak enough so that it can be implemented with bounded memory, yet strong enough to be useful for applications that typically use atomic broadcast, such as state-machine replication. Note that repeated relaxed *atomic* broadcast is implementable with bounded memory in systems where repeated *reliable* broadcast is not. The intuition behind relaxed atomic broadcast is the following. As long as no process lags behind in the execution, relaxed atomic broadcast ensures the same properties as (classic) atomic broadcast. When some process  $p$  appears to be slow, other processes, instead of keeping on buffering messages for  $p$ , discard these messages. As a result,  $p$  will not be able to deliver all the messages that were atomically

broadcast. Missing messages are replaced at  $p$  with the special  $\perp$  message (void), which signals that a message could not be delivered.

At first sight it may seem complicated, when using relaxed atomic broadcast for state-machine replication, to recover from the delivery of  $\perp$ . However, whenever some process  $p$  delivers  $\perp$ , the specification of relaxed atomic broadcast ensures that there exists some correct process that has delivered the missing message and applied it to its state. Thus state transfer, as in the case of dynamic groups, will allow  $p$  to recover from the delivery of  $\perp$ .<sup>2</sup>

The paper is organized as follows. The system model is presented in Section II. Section III discusses atomic broadcast and the problem of implementing repeated atomic broadcast with bounded memory. Approaches to address this are discussed in Section IV. Section V presents our novel approach. In Section VI, we present the implementation of relaxed atomic broadcast and its memory bounds. Section VII compares relaxed atomic broadcast over the solution that uses dynamic groups. Section VIII concludes the paper.

## II. SYSTEM MODEL

We consider a system with a finite set of processes  $\Pi = \{p_1, p_2, \dots, p_n\}$  that communicate by message exchange. Set  $\Pi$  has cardinality of  $n$ . We assume a partially synchronous system [?], where after some unknown time *GST* (*Global Stabilization Time*) the system (both processes and channels) becomes synchronous and channels become reliable.<sup>3</sup> Before GST the system is asynchronous and channels are lossy. Processes can only fail by crashing. A process that crashes stops its operation permanently and never recovers. A process is *faulty* in a run if it crashes in that run. A process is *correct* in a run if it is not faulty in that run. We only consider runs where up to  $f$  processes are faulty ( $f$  is a system parameter). Since processes do not know whether they are before or after GST, a slow process (or a process connected through a slow link) is indistinguishable from a crashed process.

Every pair of processes is connected by a bidirectional communication channel, which provides two communication primitives:  $send(m, q)$  and  $receive(m, q)$ , where  $m \in \mathcal{M}$  (the set of messages) and  $q \in \Pi$ . Channels satisfy the properties mentioned above. All messages in  $\mathcal{M}$  are unique: they are broadcast at most once in a given run.

## III. REPEATED ATOMIC BROADCAST AND FINITE MEMORY

We recall the definition of atomic broadcast. We say that a process  $p$  atomically broadcasts (or simply *abcasts*)

<sup>2</sup>The size of the application state is controlled (and bounded) by the application. This is different from the state required for the implementation of atomic broadcast, which cannot be controlled by the application.

<sup>3</sup>We could also consider a system that alternates between sufficiently long *good* periods (system is synchronous and channels are reliable) and *bad* periods (system is asynchronous and channels are lossy). The algorithms would be the same.

message  $m$  if  $p$  executes  $abcast(m)$ . Likewise, we say that a process  $p$  atomically delivers (or simply *adeli*vers) message  $m$  if  $p$  executes  $adeli$ ver( $m$ ). Atomic broadcast is defined by the following properties:

*Property 3.1: VALIDITY.* If a correct process  $p$  abcasts message  $m$ , then some correct process eventually adelivers  $m$ .

*Property 3.2: UNIFORM INTEGRITY.* Every process adelivers a message  $m$  at most once and only if  $m$  was previously abcast by some process.

*Property 3.3: UNIFORM AGREEMENT.* If a process adelivers a message  $m$  then every correct process also adelivers  $m$ .

*Property 3.4: UNIFORM TOTAL ORDER.* For any two processes  $p$  and  $q$  and any two messages  $m$  and  $m'$ , if  $p$  adelivers  $m$  before  $m'$ , then  $q$  adelivers  $m'$  only after having adelivered  $m$ .

Repeated atomic broadcast is the case where at least one process executes atomic broadcast infinitely often.

Reliable broadcast is defined by properties 3.1, 3.2, and the non-uniform version of 3.3. As shown by Ricciardi, repeated reliable broadcast cannot be implemented in a system with message losses in which slow processes are indistinguishable from crashed processes [?]. The intuition behind this impossibility result is the following. Consider a sender process  $p$ , and its output buffer to  $q$  that contains unacknowledged messages sent to  $q$ . If  $p$  is unable to distinguish whether  $q$  has crashed or is just slow (or connected through a slow link), then  $p$  cannot safely dispose of unacknowledged messages sent to  $q$ . However, if  $q$  has actually crashed, the set of unacknowledged messages will grow forever [?].

The impossibility of repeated reliable broadcast also applies to repeated atomic broadcast, since atomic broadcast is strictly stronger than reliable broadcast.

## IV. HOW TO DEAL WITH FINITE MEMORY

Consider atomic broadcast used to implement state-machine replication [?] in a system with three processes ( $n = 3$ ). Process  $p_1$ , which receives clients' requests, issues abcasts. Assume that the adelivery of these messages requires the cooperation of  $p_1$  with only  $p_2$  or with only  $p_3$ . Consider the former case, and assume  $p_3$  is slow (or connected to  $p_1$  and  $p_2$  through slow channels). Since  $p_1$  and  $p_2$  do not know whether  $p_3$  has crashed or not, they cannot safely dispose of unacknowledged messages sent to  $p_3$ , and their buffer to  $p_3$  may grow infinitely.

We now present two approaches to deal with this problem.

*The dynamic model:* The traditional solution to bound memory consists in switching to the dynamic system (or dynamic group) model [?], [?], [?], [?], [?].<sup>4</sup> In such a model processes can be added/removed to/from

<sup>4</sup>Note that this argument is not always explicit in these papers.

the system (or group) on the fly. In a dynamic model, a *view* describes the set of processes that are currently part of the system (or group). Views are maintained by a *membership* service, which adds and removes processes. Let us consider again state-machine replication with three replicas  $p_1$ ,  $p_2$  and  $p_3$ . If the buffer from  $p_1$  to  $p_3$  is full,  $p_1$  may ask to remove  $p_3$  from the view. Once this is done, all unacknowledged messages to  $p_3$  can be discarded. However, the dynamic model is not so straightforward as the static one: protocol specifications and implementations have to be revised [?] and are more complex. Besides, a membership service is needed, and the application logic needs to become aware of view changes and state transfers (which are needed when an excluded process re-joins the group).

*Relaxing the specification of atomic broadcast:*

The paper proposes another – novel – way to deal with bounded memory. Instead of switching to the dynamic model, we propose to *relax* the specification of atomic broadcast. This is done while keeping the specification strong enough to be useful for practical systems, and ensuring that repeated *relaxed* atomic broadcast is solvable with bounded memory.

## V. RELAXED ATOMIC BROADCAST

We start by defining relaxed atomic broadcast, and then we show how state-machine replication can be implemented using this new primitive.

### A. Specification of relaxed atomic broadcast

We start by extending the set of messages that are delivered with the special void message  $\perp$ , which is not in set  $\mathcal{M}$ . Unlike any other message, this message is not unique, i.e., there may be more than one occurrence of this message in one run. A message  $m$  is called *normal* if it is not the void message  $\perp$  (i.e., if  $m \in \mathcal{M}$ ). The void message  $\perp$  is never broadcast by the application, but might be delivered in substitution of a normal message in certain scenarios. The delivery of  $\perp$  warns the application that a message is missing in its delivery sequence.

We define relaxed atomic broadcast with the primitives  $xbcast(m)$  and  $xdeliver(m')$ , where  $m \in \mathcal{M}$ , and  $m' \in \mathcal{M} \cup \{\perp\}$ . Relaxed atomic broadcast is also called *x-atomic* broadcast. For  $k$  a positive integer, we say that a process  $p$  *xdelivers@k* message  $m$  if  $m \in \mathcal{M} \cup \{\perp\}$  and  $m$  is the  $k^{th}$  message *xdelivered* by  $p$  since system start-up time. If  $k$  is irrelevant then the suffix *@k* is omitted, i.e., *xdeliver@k* simply becomes *xdeliver*. Relaxed atomic broadcast satisfies the following properties:

*Property 5.1: VALIDITY.* If a correct process  $p \in \Pi$  *xbcasts* message  $m$ , then some correct process  $q \in \Pi$  eventually *xdelivers*  $m$ .

This property does not change with respect to classic atomic broadcast (see Sect. III).

*Property 5.2: UNIFORM AGREEMENT.* For all  $k \geq 1$ , if some process *xdelivers@k* a normal message or  $\perp$ , then every correct process *xdelivers@k* a normal message or  $\perp$ .

The uniform agreement property is usually stated in terms of a given message  $m$ . In contrast, this weaker form only forces correct processes to *xdeliver* (at least) as many messages (normal or  $\perp$ ) as any other process.

*Property 5.3: UNIFORM TOTAL ORDER.* For all  $k, k' \geq 1$ , if process  $p$  *xdelivers@k* normal message  $m$  and process  $q$  *xdelivers@k'* normal message  $m'$ , then  $k = k' \Leftrightarrow m = m'$ .

The simplicity of the definition of uniform total order benefits from the definition of *xdelivery@k*. Property 3.4 could also benefit from this definition, thus becoming simpler.

*Property 5.4: UNIFORM INTEGRITY.* A process *xdelivers* a normal message  $m$  only if  $m$  was previously *xbcast*.

This property is simplified with respect to classic atomic broadcast for two reasons: (1) to allow the void message  $\perp$  to be *xdelivered* more than once, and (2) because Property 5.3 already forbids *xdelivering* a normal message more than once.

*Property 5.5: CONTINUITY.* For all  $k \geq 1$ , a process *xdelivers@k* the void message  $\perp$  only if at least one correct process *xdelivers@k* a normal message.

This safety property forbids runs where no correct process *xdelivers* a normal message at some position in the delivery sequence. Examples of such runs are (1) all processes *xdeliver@k* message  $\perp$ , or (2) correct processes *xdeliver@k* message  $\perp$  and faulty processes *xdeliver@k* a normal message (and crash immediately after). In both cases, the application at surviving processes may not be able to reconstruct a complete delivery sequence of normal messages (i.e., without gaps).

The specification of relaxed atomic broadcast reduces to that of classic atomic broadcast in runs where no void message  $\perp$  is ever *xdelivered*. Relaxed atomic broadcast is thus strictly weaker: any algorithm solving atomic broadcast also solves relaxed atomic broadcast.

### B. Is the new specification useful?

We illustrate now the usefulness of relaxed atomic broadcast in the context of state-machine replication, see Algorithm 1. Basically, the algorithm works as though it was using classic atomic broadcast, but in addition it needs to implement a state transfer in order to recover from gaps in the sequence (when  $\perp$  is *xdelivered*).

The (simple) algorithm works as follows. Two counters keep track of (1) the number of messages *xdelivered*,  $n\text{-}xdel_p$ ; and (2) the number of (normal) messages that have been applied to the application's state,  $n\text{-}st_p$  (i.e.,  $n\text{-}st_p$  messages, in sequence, have updated the application state). Initially these two counters match, and when a normal message is *xdelivered* both are incremented (lines 10 and 15).

If the void message  $\perp$  is *xdelivered*, only  $n\text{-}xdel_p$  is incremented to reflect the *xdelivery*, and process  $p$  halts its execution (line 13) until it receives a (more recent) state

---

**Algorithm 1** State machine replication using relaxed atomic broadcast. Code for process  $p$ .

---

```

1: Initialization:
2:    $n\text{-}xdel_p \leftarrow 0$            {Number of messages xdelivered}
3:    $n\text{-}st_p \leftarrow 0$         {Number of messages applied to current state}
4:    $state_p \leftarrow \text{initial state}$    {Replicated state}

5: task Main Thread
6:   repeat forever
7:     wait until received request  $m$  from user
8:      $\text{xbroadcast}(m)$ 

9: upon  $xdeliver(m)$  do
10:   $n\text{-}xdel_p \leftarrow n\text{-}xdel_p + 1$ 
11:  if  $n\text{-}xdel_p = n\text{-}st_p + 1$  then           {Any gaps so far?}
12:    if  $m = \perp$  then
13:      wait until  $n\text{-}xdel_p \leq n\text{-}st_p$ 
14:        {Halt  $xdelivery$  of  $\perp$  until a useful state received}
15:    else
16:       $n\text{-}st_p \leftarrow n\text{-}st_p + 1$ 
17:       $state_p \leftarrow \text{apply } m \text{ to } state_p$ 

18: task Resend
19:   repeat forever
20:     if  $n\text{-}xdel_p > n\text{-}st_p$  then
21:       send  $\langle \text{STATE-REQ}, n\text{-}xdel_p \rangle$  to all

22: upon receive  $\langle \text{STATE-REQ}, n \rangle$  from  $q$  do
23:   if  $n \leq n\text{-}st_p$  then
24:     send  $\langle \text{STATE-REP}, n\text{-}st_p, state_p \rangle$  to  $q$ 

25: upon receive  $\langle \text{STATE-REP}, n, st \rangle$  from  $q$  do
26:   if  $n \geq n\text{-}xdel_p$  then
27:      $n\text{-}st_p \leftarrow n$ 
28:      $state_p \leftarrow st$ 

```

---

from another process  $q$  whose state has been updated by applying the message missing at  $p$ . To do so, if process  $p$  detects that the number of messages applied to its state ( $n\text{-}st_p$ ) lags behind with respect to the number of messages  $xdelivered$  ( $n\text{-}xdel_p$ ) due to the  $xdelivery$  of  $\perp$ , then  $p$  starts sending out state request messages repeatedly (line 20). When another process  $q$  receives the state request message (line 21), it checks whether its current state would be useful to the requesting process (the state is useful if it has been updated with at least as many messages as specified in the state request). If so,  $q$  sends back a state reply with its state and  $n\text{-}st_p$ . Finally, when the sender of the request receives a state reply (line 24) it checks whether that state is recent enough to fill the gaps in its  $xdelivery$  sequence. If it is the case, it replaces its state by the one it has just received, and updates  $n\text{-}st_p$  accordingly. Note that the state received by  $p$  might have been updated with messages that have not (yet) been  $xdelivered$  at  $p$ . In this case, the algorithm ignores those messages when they are finally  $xdelivered$  (line 11).

If the application state is large, state transfer may be costly. However, this cost is the same as with the dynamic

group solution.

*Concurrency control:* The state updates when an *upon* clause is executed should be atomic to avoid inconsistencies. A simple approach is to assume that the algorithm behaves like a monitor: *upon* clauses and tasks are executed in mutual exclusion, except when a *wait until* statement is reached, where another task or *upon* clause can take over the execution. Task *Resend* is an exception: it executes in mutual exclusion only within its loop: mutual exclusion is not preserved across consecutive executions of lines 19-20.

*Memory bounds:* The memory required by Algorithm 1 is bounded if we can bound the memory usage of relaxed atomic broadcast. Indeed, Algorithm 1 uses (1) two integers ( $M_{int}$  bits for each, see discussion in Section VI-B), (2) needs to store the application state that we assume to be bounded by  $M_{state}$  and a client request  $m$  that we assume to be bounded by  $M_{req}$ , and (3) needs memory space for the interaction between Algorithm 1 and the communication channels, and between Algorithm 1 and the relaxed atomic broadcast implementation (see Figure 1).

The interaction between Algorithm 1 and the communication channels is modeled by input and output buffers. Only one of each is represented in Figure 1, although we assume one pair for each channel (total of  $n$  pairs). Sending a message  $m$  is modeled by writing  $m$  into the output buffer. Receiving a message is modeled by an up-call that reads the input buffer and hands it over to Algorithm 1 (lines 21 and 24). These two buffers are bounded by the size of the longest message, the one with tag STATE-REP. The bound is  $1 + M_{int} + M_{state}$  bits. The interaction between Algorithm 1 and the relaxed atomic broadcast implementation is modeled by function calls ( $\text{xbroadcast}$  is a down-call,  $\text{xdeliver}$  is an up-call). This interaction model does not add anything to the memory requirements of both components.

## VI. IMPLEMENTING REPEATED RELAXED ATOMIC BROADCAST WITH BOUNDED MEMORY

In this section, we present an algorithm that implements repeated relaxed atomic broadcast with bounded memory. For the sake of simplicity, from now on whenever we use the term relaxed atomic broadcast, we mean *repeated* relaxed atomic broadcast.

We first introduce the building blocks that our application (state machine replication) uses along with their interaction model, then we present the implementation of each building block followed by an analysis of the amount of memory needed. We will also have a short discussion regarding integers.

### A. Building blocks and interaction model

Figure 1 depicts the building blocks of our implementation, as well as their interactions. Relaxed atomic broadcast uses consensus, and consensus is expressed in a round-based model implemented by the corresponding building block. The round-based model block interacts

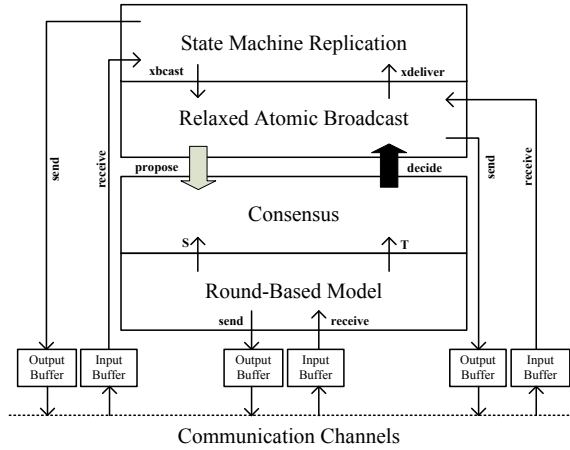


Figure 1. Building blocks. Small arrows represent function calls; large arrows represent spawning/killing (propose), and decision delivery (decide) of consensus instances.

with consensus by calling functions  $S$  and  $T$ . Likewise, state machine replication and relaxed atomic broadcast interact by calling functions  $xbroadcast$  and  $xdeliver$ , which are called in opposite direction. The interaction between relaxed atomic broadcast and consensus is different: when relaxed atomic broadcast calls  $propose$  a new instance of the consensus and round-based blocks (as well as their input/output buffers) is spawned, and any previous instance of these created blocks is immediately killed and garbage-collected. When consensus calls  $decide$ , a task within relaxed atomic broadcast is already waiting for it, so the call simply unblocks the task (and passes  $decide$ 's parameters) as we will see later. As explained above for state machine replication, the interaction with the channels is represented by input buffers and output buffers, one pair of buffers for the “relaxed atomic broadcast” block (i.e., one pair per channel), and one pair for the implementation of the round model (one pair per channel).

### B. The issue of integers

Integer variables are used by all layers of our implementation. Some of these integers, such as message ids or round numbers are constantly increasing during system lifetime. This means that, at least theoretically, the number of bits needed by these variables cannot be bounded. However, this is not a problem from a practical point of view. Indeed, if we use 64 bits to represent some integer variable  $i$ , and we assume that  $i$  is increased by 1 every micro-second, then the largest integer is reached only after 584'000 years. This is long enough from a practical point of view (see also related discussion in [?]).

### C. Relaxed atomic broadcast

1) *Algorithm*: Algorithm 2 solves relaxed atomic broadcast, for  $f < n/2$ , by reduction to a sequence of consensus [?]. However, contrary to [?], each consensus decides only on one single message (in order to bound memory) rather than on a batch of messages. Although a number of optimizations can be performed, we have kept

the algorithm as simple as possible, while preserving its correctness (see [?] for the proofs).

The algorithm is structured in two tasks, *Sequencer* and *Gossip*, and works as follows. When  $p$ 's application  $xbroadcast$ s a message  $m$ , a new identifier is attached to  $m$ . Then,  $m$  is stored in  $Rcv_p[p]$  (line 9). Vector  $Rcv_p$  contains messages that  $p$  knows of but has not yet  $xdelivered$ . If  $p$  has previously  $xbroadcast$  another message  $m'$  not yet  $xdelivered$ , then  $p$ 's application is blocked (i.e.,  $p$  cannot  $xbroadcast$  any further message), since  $Rcv_p[p]$  can only store one message at a time. This is a simple flow-control technique that can be optimized. The elements of vector  $Rcv_p$  will later become proposed values for consensus. This is the mission of task *Sequencer* (line 37), which executes a sequence of consensus instances. The *Sequencer* task waits until there are undelivered messages in vector  $Rcv_p$  (lines 39-40). Then, it starts a new consensus instance. For each instance  $\#k_p$  a sender  $c_p$  is designated in a round-robin manner, with the goal to propose  $Rcv_p[c_p]$  as the initial value for consensus (line 42). This initial value could be optimized to be the whole  $Rcv_p$  vector [?], but the rotating sender approach makes it easier to present both our algorithm and its memory bounds. When consensus  $\#k_p$  decides,  $p$  waits for evidence that at least  $f + 1$  other processes have also decided for consensus  $\#k_p$  (line 45). This mechanism enforces the continuity property of relaxed atomic broadcast, since it ensures that at least one correct process (that can be queried later) has decided. Then,  $p$   $xdelivers$  the message in  $decision_p$  only if its identifier matches the value of  $NextId_p[c_p]$ , otherwise the decision is discarded (lines 46 and 29-32). This simple method demonstrates how to avoid  $xdelivering$  duplicates using bounded memory. Its side effect is that it enforces FIFO order amongst messages  $xbroadcast$  by process  $c_p$ . This may affect performance, but the algorithm can be optimized to relax this condition. Finally, variable  $k_p$  is incremented (line 33) and the loop starts over with a new iteration.

The *Gossip* task (line 34) sends periodically GOSSIP messages to all processes in order to disseminate (1) recently  $xbroadcast$  messages (vector  $Rcv_p$ ), and (2) the status of the sender's current consensus instance ( $k_p$  and  $decided_p$ ). When process  $p$  receives a GOSSIP message from process  $q$  (line 10), it checks whether  $q$  is either ahead or lagging behind. If  $q$  is ahead (or at the same consensus instance as  $p$  but has already decided),  $p$  adds  $q$  to its set  $Finished_p$  (line 12), which contains processes that already finished  $p$ 's current consensus. When the size of this set reaches  $f + 1$ ,  $p$  can infer that at least one correct process has decided; so  $p$  can proceed to consensus  $k_p + 1$  as soon as it is done with consensus  $k_p$  (line 45 is no longer blocking). If  $q$  is lagging behind (line 13), then  $p$  simply sends  $q$  a SLOW message containing part of its current state. Additionally, if both  $p$  and  $q$  are at the same consensus instance (line 14), then  $p$  copies to its  $Rcv_p$  vector all messages received from  $q$  that  $p$  has not yet  $xdelivered$ .

A SLOW message conveys the part of the sender's state that a slow process needs in order to catch up. Upon

---

**Algorithm 2** Solving relaxed atomic broadcast. Code for process  $p$ .

---

```

1: Initialization:
2:    $id_p \leftarrow 0; c_p \in \Pi; decision_p \in \mathcal{M}$ 
3:    $k_p \leftarrow 0; Finished_p \leftarrow \emptyset; decided_p \leftarrow \mathbf{false}$ 
4:   for all  $r \in \Pi$  do  $Rcv_p[r] \leftarrow \perp; NextId_p[r] \leftarrow 0$ 
5:   fork_task(Gossip, Sequencer)

6: upon  $xbcast(m)$  do
7:    $m.id \leftarrow id_p; id_p \leftarrow id_p + 1$ 
8:   wait until  $NextId_p[p] = m.id$ 
9:    $Rcv_p[p] \leftarrow m$ 

10: upon  $receive(GOSSIP, k_q, d_q, Rcv_q)$  from  $q$  do
11:   if  $k_q > k_p$  or  $k_q = k_p$  and  $d_q$  then
12:      $Finished_p \leftarrow Finished_p \cup \{q\}$ 
13:   if  $k_q < k_p$  then  $send(SLOW, k_p, NextId_p)$  to  $q$ 
14:   if  $k_q = k_p$  then {Message dispersal}
15:     for all  $r \in \Pi$  do
16:       if  $Rcv_q[r] \neq \perp$ 
17:         and  $Rcv_q[r].id = NextId_p[r]$  then
18:            $Rcv_p[r] \leftarrow Rcv_q[r]$ 

19: upon  $receive(SLOW, k_q, N_q)$  from  $q$  do
20:   if  $k_q > k_p$  then {p is late}
21:     kill_task(Sequencer)
22:     if  $decided_p$  then  $deliver()$ 
23:      $msgs\_skipped \leftarrow \sum_{r \in \Pi} (N_q[r] - NextId_p[r])$ 
24:     repeat  $msgs\_skipped$  do  $xdeliver(\perp)$ 
25:      $NextId_p \leftarrow N_q; k_p \leftarrow k_q$ 
26:      $Finished_p \leftarrow \emptyset; decided_p \leftarrow \mathbf{false}$ 
27:     fork_task(Sequencer)

28: procedure  $deliver()$ 
29:   if  $decision_p \neq \perp$ 
30:     and  $decision_p.id = NextId_p[c_p]$  then
31:        $xdeliver(decision_p)$ 
32:        $NextId_p[c_p] \leftarrow NextId_p[c_p] + 1$ 
33:    $k_p \leftarrow k_p + 1$ 

34: task Gossip
35:   repeat forever
36:      $send(GOSSIP, k_p, decided_p, Rcv_p)$  to all

37: task Sequencer
38:   repeat forever
39:     wait until  $\exists r : (Rcv_p[r] \neq \perp$ 
40:       and  $Rcv_p[r].id = NextId_p[r])$ 
41:      $c_p \leftarrow k_p \bmod |\Pi|$  { $c_p$  is a rotating sender}
42:      $propose(k_p, Rcv_p[c_p])$  {Delete previous instance}
43:     wait until  $decide(k_p, decision_p)$ 
44:      $decided_p \leftarrow \mathbf{true}$ 
45:     wait until  $|Finished_p| > f$ 
46:      $deliver()$ 
47:      $Finished_p \leftarrow \emptyset; decided_p \leftarrow \mathbf{false}$ 

```

---

reception of such a message (line 19), process  $p$  checks whether the sender is ahead. If that is the case,  $p$  has been lagging behind, so termination of its current consensus instance is not guaranteed because other processes have already moved on to a later instance and disposed of  $p$ 's current consensus (see Sect. VI-D). Therefore,  $p$  stops task *Sequencer* (line 21) and checks whether its current consensus had already finished. If so, the decision is xdelivered (line 22) and  $p$  advances to the next consensus. At this point, if  $p$  is still lagging behind with respect to  $q$ , the following catch-up mechanism is used. Process  $p$  calculates the number of messages it is going to skip when catching up: for each process  $r$ ,  $p$ 's next message id for process  $r$  is subtracted from  $q$ 's (possibly greater) value (line 23). The result of this subtraction is the number of messages sent by  $r$  that were xdelivered between the consensus instances in which  $p$  and  $q$  are. The sum of all these subtractions yields the total amount of messages  $p$  will skip, so it xdelivers as many  $\perp$  messages (line 24). Finally,  $p$  updates  $k_p$  and  $NextId_p$  with the values received from  $q$  and spawns task *Sequencer* again. Note that additional garbage collection can be performed on  $Rcv_p$ , but does not affect correctness.

2) *Concurrency control*: The state of the protocol, in particular variables  $k_p$ ,  $NextId_p$ ,  $Finished_p$ , and  $decided_p$ , should all be updated atomically every time a new consensus instance starts. As in Sect. V-B, a simple approach is to assume that the algorithm behaves like a monitor: *upon* clauses and tasks are executed in mutual exclusion, except when a *wait until* statement is reached. Finally, task *Gossip* executes in mutual exclusion only within its loop (i.e., mutual exclusion is not preserved across consecutive executions of line 36).

3) *Memory bounds*: We show now that our algorithm requires only bounded memory as long as the size of the application payload is bounded to constant  $M_{req}$  (see Section V-B) and consensus requires a maximum of  $M_{cons}$  bits (see Section VI-D).

*State size*: To avoid a boring enumeration, let us assume that the space required for all variables except  $decision_p$  and the vector  $Rcv_p$  amounts to some constant  $c(n)$  (that depends on  $n$ ). Moreover,  $decision_p$  may contain an application message with an attached message  $id$  and vector  $Rcv_p$  is a vector of at most  $n$  application messages with added ids. Together this leads to  $(n+1) \cdot (M_{req} + M_{int})$  bits. Since at most one consensus instance is running at each process, summing everything up, the state space needed by relaxed atomic broadcast is bounded by

$$M_{xbcast} = M_{cons} + (n+1) \cdot (M_{req} + M_{int}) + c(n).$$

*Buffer size*: The algorithm sends/receives two types of messages: GOSSIP and SLOW, with respectively four and three parameters. The former conveys the GOSSIP tag, one integer  $k_q$ , boolean  $d_q$ , and set  $Rcv_q$  of messages with attached ids. The latter contains the SLOW tag, one integer  $k_q$ , and set  $N_q$  of message ids. One bit is enough to represent the message tags. If we use again  $c(n)$  to

represent a constant depending on  $n$ , we get the following bounds:

$$M_{gossip} = n \cdot M_{req} + c(n),$$

$$M_{slow} = c(n).$$

#### D. Consensus

The relaxed atomic broadcast algorithm relies on a consensus algorithm, which ensures the following usual properties:

- *Validity*: If process  $p$  decides  $v$ , then  $v$  has been proposed by some process.
- *Uniform agreement*: No two processes decide differently.
- *Termination*: All correct processes eventually decide.

An unbounded number of consensus instances may be spawned in every run. Every instance of consensus uses its own memory resources. However, each process maintains only one single instance of consensus at a given time. When a process executes *propose*, its current consensus instance (if any) is immediately killed and garbage-collected. Therefore, the termination property of consensus is not guaranteed for all correct processes; rather, only  $f + 1$  processes (whether correct or not) are guaranteed to terminate a consensus instance. Nevertheless, once  $f + 1$  processes have decided for consensus  $\#k$  (i.e., at least one correct process), Algorithm 2 guarantees that all correct processes will eventually stop consensus  $\#k$  and move on to  $\#k + 1$ .

##### 1) Algorithm:

*Round-based model*: We consider a consensus algorithm for a partially synchronous system (see Section II). As in [?], we consider an abstraction on top of the system model, namely a round model. Using this abstraction, rather than the raw system model, improves the clarity of the algorithms and simplifies the proofs. In the round model, processing is divided into rounds of message exchange. Each round  $r$  consists of a *sending step* denoted by  $S_p^r$  (sending step of  $p$  for round  $r$ ), and of a *state transition step* denoted by  $T_p^r$ . In a sending step, each process sends a message to all. A subset of the messages sent is received at the beginning of the state transition step: messages can get lost, and a message sent in round  $r$  can only be received in round  $r$ . We denote by  $\sigma_p^r$  the message sent by  $p$  in round  $r$ , and by  $\vec{\mu}_p^r$  the messages received by process  $p$  in round  $r$  ( $\vec{\mu}_p^r$  is a vector of size  $n$ , where  $\vec{\mu}_p^r[q]$  is the message received from  $q$  or *null* if the message was lost). Based on  $\vec{\mu}_p^r$ , process  $p$  updates its state in the state transition step.

In all rounds executed before *GST* messages can be lost. However, after *GST*, there exists a round *GSR* (*Global Stabilization Round*) such that the message sent in round  $r \geq \text{GSR}$  by a correct process  $q$  to a correct process  $p$  is received by  $p$  in round  $r$ . This is formally expressed by the following predicate (where  $\mathcal{C}$  denotes the set of correct processes):

$$\forall r \geq \text{GSR} : \mathcal{P}_{good}(r),$$

where

$$\mathcal{P}_{good}(r) \equiv \forall p, q \in \mathcal{C} : \vec{\mu}_p^r[q] = \sigma_q^r.$$

An algorithm that ensures this predicate in a partially synchronous system is given in Section VI-E.

---

**Algorithm 3** The *OneThirdRule* (OTR) algorithm [?] ( $f < n/3$ ). Code for process  $p$ .

---

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** (number of messages sent in round  $r$  and received by  $p$  in round  $r$ )  $> 2n/3$  **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  $x$  received

12: **if** more than  $2n/3$  values received are equal to  $\bar{x}$  **then**

13: DECIDE( $\bar{x}$ )

---

*The OTR consensus algorithm*: Algorithm 3 is the consensus algorithm we consider [?]. The algorithm requires  $f < n/3$ . We have chosen this algorithm because of its simplicity. The analysis of *Paxos/LastVoting* [?], [?], which requires only  $f < n/2$  could be used instead, but would require more space.

Algorithm 3 works as follows. As soon as more than  $2n/3$  processes have  $x_p = v$ , then decision  $v$  is locked, i.e., in any future update, variable  $x_p$ , is updated to  $v$ . Termination is ensured by the following observation. Let  $r_0$  be the smallest round after *GSR* such that all faulty processes have crashed before round  $r_0$ . In round  $r_0$  the condition of line 7 is true. Moreover,  $\mathcal{P}_{good}(r_0)$  ensures that all processes that execute round  $r_0$  receive the same set of messages. Therefore, in round  $r_0$ , either all processes execute line 9, or all processes execute line 11. It follows that at the end of round  $r_0$  all processes have  $x_p$  equal to some common value  $v$ , and all processes decide in round  $r_0 + 1$ .

2) *Memory bounds*: As we explain in Section VI-E, the memory required by Algorithm 3 is managed by the implementation of the round-based model. Thus we refer to the next section for the consensus memory bounds.

#### E. Implementation of the round-based model

We describe now the implementation of the round-based model (see Algorithm 4), which is almost identical to the one appearing in [?] (we made small extensions to prevent *msgsRcv<sub>p</sub>* from growing forever). The interaction between Algorithm 4 and Algorithm 3 is by function call: in other words, the execution thread is within Algorithm 4,

and this thread calls functions  $S_p^r$  and  $T_p^r$  defined by Algorithm 3:

- $S_p^r$  is called at line 9 of Algorithm 4 and returns  $x_p$ , see line 5 of Algorithm 3.<sup>5</sup>
- $T_p^r$  is called at line 22 of Algorithm 4 and returns the new state of process  $p$ , see lines 7 to 13 of Algorithm 3.<sup>6</sup>

The state of Algorithm 3 is represented as  $s_p$  in Algorithm 4 (line 3). Moreover, in Algorithm 4,  $\phi$  represents the bound on process relative speed after *GSR*, and  $\delta$  represents the bound on message transmission delay after *GSR*. After *GSR* one send step (line 10) and one receive step (line 16) take each 1 time unit on the fastest process (i.e., at most  $\phi$  time units on the slowest process). If no message is available for reception, then an empty message is received. In one send step a process can send messages to multiple processes, while  $n$  receive steps are needed to receive messages from  $n$  processes.

1) *Algorithm*: Algorithm 4 consists of an infinite loop (see line 8), which includes an inner loop (lines 12 to 21). Each iteration of the outer loop corresponds to one round. The message to send is obtained in line 9, and sent to all in line 10. Each iteration of the inner loop is for the reception of one message for the current round  $r_p$ . The inner loop ends when (i) at least  $2\delta + (n+2)\phi$  time units have elapsed, see lines 14-15 (time is measured by the execution of receive steps: 1 receive step = 1 time unit), or (ii) whenever a message of a round larger than  $r_p$  is received, see lines 20-21. The reader is referred to [?] for a proof that this ensures  $\mathcal{P}_{good}$  after GST. When the inner loop ends, the function  $T_p^r$  is called with the set of messages received in the current round  $r_p$  (line 22). Finally, messages for the current round are garbage collected (line 24).

2) *Memory bounds*: We compute now  $M_{cons}$  – the memory bound for consensus including the implementation of the round-based model – that was referenced in Section VI-C3.

*State size*: Algorithm 4 needs to store three integers ( $r_p$ ,  $next\_r_p$ ,  $i_p$ ), which require  $3M_{int}$ , and variables  $s_p$  and  $msg_p$ , which require  $2M_{req}$  bits. In addition the algorithm needs memory for  $msgsRcv_p$  and  $temp_p$ , which amounts to  $(n+1) \cdot (M_{req} + 2M_{int})$  bits, since  $msgsRcv_p$  stores at most  $n$  messages.

*Buffer size*: All messages sent/received are of the same type and require at most  $M_{req} + M_{int}$  bits each. The algorithm needs only one single output buffer (the same message sent to all) and  $n$  input buffers (one per process). This amounts to  $(n+1) \cdot (M_{req} + M_{int})$  bits.

### F. Summary

Putting everything together, we have shown that all components that appear in Figure 1, including state-machine replication, require only bounded memory. There-

<sup>5</sup>To be consistent, line 4 of Algorithm 3 should be expressed as a function. However, we decided to keep the usual round-based expression for Algorithm 3.

<sup>6</sup>Same comment as for  $S_p^r$ , see footnote 5.

---

### Algorithm 4 Ensuring $\mathcal{P}_{good}$ after GST.

---

```

1:  $r_p \leftarrow 1$  {round number}
2:  $next\_r_p \leftarrow 1$  {next round number}
3:  $s_p \leftarrow init_p$  {state of the consensus algorithm}
4:  $i_p \leftarrow 0$  {counts send/receive steps}
5:  $msg_p$  {message to send in the current round}
6:  $msgsRcv_p \leftarrow \emptyset$  {set of msgs received for the current round}
7:  $temp_p \leftarrow \emptyset$  {contains at most one message received for a round  $> r_p$ }
8: while true do
9:    $msg_p \leftarrow S_p^r(s_p)$ 
10:  send  $\langle msg_p, r_p \rangle$  to all
11:   $i_p \leftarrow 0$ 
12:  while next_r_p = r_p do
13:     $i_p \leftarrow i_p + 1$ 
14:    if  $i_p \geq 2\delta + (n+2)\phi$  then
15:       $next\_r_p \leftarrow r_p + 1$ 
16:      receive a message with highest round number
17:      if received  $\langle msg, r' \rangle$  from  $q$  then
18:        if  $r' = r_p$  then
19:           $msgsRcv_p \leftarrow msgsRcv_p \cup \{\langle msg, r', q \rangle\}$ 
{Messages from old rounds are discarded}
20:        if  $r' > r_p$  then
21:           $next\_r_p \leftarrow r'$ ;  $temp_p \leftarrow \{\langle msg, r', q \rangle\}$ 
22:         $s_p \leftarrow T_p^r(msgsRcv_p, s_p)$ 
23:         $r_p \leftarrow next\_r_p$ 
24:         $msgRcv_p \leftarrow temp_p$  {Garbage collection}
25:         $temp_p \leftarrow \emptyset$ 

```

---

fore, relaxed atomic broadcast has allowed us to implement state-machine replication using bounded memory.

## VII. COMPARISON OF APPROACHES

In Section IV, we have presented two different approaches for implementing state machine replication with bounded memory. Namely, (1) our novel relaxed atomic broadcast algorithm, which was described in detail in Sections V and VI, and (2) atomic broadcast in the dynamic model, i.e., relying on membership [?]. Both approaches rely on state transfer: approach (2) requires a state transfer whenever a new process is added to the dynamic group; approach (1) performs a state transfer whenever a slow process catches up.

Solution (2) is more complex than solution (1). First, solution (2) needs to define a policy for process exclusion [?]. This is simply not needed in (1). Second, static group communication is simpler and easier to understand than dynamic group communication, from a specification as well as from an implementation point of view. Moreover, the complexity added by relaxed atomic broadcast (i.e., the need for state transfer) is also needed in dynamic group communication, as stated above.

If an application is happy with the static group model, and dynamism is introduced only to bound the memory usage, then the solution using relaxed atomic broadcast is a



better choice. If an application requires the dynamic group model, the solution using relaxed atomic broadcast may still be used: it makes sense to combine both approaches, where changes in membership are decoupled from the bounded memory issue.

## VIII. CONCLUSION

We have presented relaxed atomic broadcast, a variant of atomic broadcast that it is weak enough to be solved with bounded memory, yet strong enough to be useful for typical applications like state machine replication. Note that the analysis of the memory requirements forced us to consider the complete protocol stack (i.e., nothing has been swept under the carpet). We have also discussed the advantages of our approach as compared to the solution with group membership.

The solution presented shows an interesting trade-off between the memory allocated and the number of  $\perp$  messages delivered: if a process becomes slow, the more memory we allocate, the longer it will take to run out of buffers. We plan to experimentally analyze this trade-off in the future.

Finally, we recall that the goal when presenting our solution was simplicity. The algorithm can be optimized in a number of ways in order to improve its performance.

## IX. APPENDIX: PROOFS FOR THE ALGORITHM 2

*Definition 9.1:* A process  $p$  arrives at iteration  $k$  at time  $t$  if its  $k_p$  value becomes equal to  $k$  using line 33 or 25 at time  $t$ . If it happens at line 33 we say process  $p$  arrives at iteration  $k$  normally, otherwise we say process  $p$  arrives at iteration  $k$  abnormally.

*Lemma 9.2:* For any process  $p$ , the value of iteration  $k$  is always non-decreasing.

*Proof:* The value of  $k$  increases either at line 25 or 33. at line 25 it increases  $k$  because of line 20. at line 33 also  $k$  will be incremented by one. In both cases  $k$  is increasing and lemma is proved. From now on we can prove lemmas using induction over  $k$ . ■

*Definition 9.3:* A process  $p$  is in iteration  $k$  at time  $t$  if it arrived at iteration  $k$  at time  $t'$  s.t.  $t' \leq t$  and did not arrive at iteration  $k' > k$  by time  $t$  and have not modified its *NextId* vector after arriving at iteration  $k$ .

*Definition 9.4:* Process  $p$  decides  $m$  in iteration  $k$  if  $p$  decides  $m$  while it is in iteration  $k$ .

*Lemma 9.5:* The first process that arrives at iteration  $k$ , arrives normally.

*Proof:* By contradiction suppose the first process, called  $p$ , arrives at iteration  $k$  at time  $t$  abnormally using line 25. It means  $p$  has received a SLOW message (line 13) from some other process who was already in iteration  $k$  at time  $t' < t$ . In this case  $p$  is not the first process who arrives at iteration  $k$ , which is a contradiction. ■

*Lemma 9.6:* If no process arrives at iteration  $k$  at time  $t$ , it is not possible for any process to arrive at iteration  $k'$  before time  $t$  where  $k' = k + 1$ .

*Proof:* By contradiction suppose a process  $p$  arrives at  $k'$  either case 1 )normally or case 2 )abnormally :

Case 1) Suppose  $p$  arrives normally at  $k'$  which means it executed line 33. This case is not possible because in this case  $p$  already was in iteration  $k$  before calling line 33 which is a contradiction.

Case 2) Suppose  $p$  arrives abnormally at  $k'$  which means it executed line 25. In this case the first process which arrived at  $k'$ , arrived normally (according to lemma 9.5). Let us call this first process  $q$  which means it executed line 33 to arrive at iteration  $k'$ . Therefore  $q$  was in iteration  $k$  before executing line 33, which is a contradiction. ■

*Corollary 9.7:* If no process arrives at iteration  $k$  by time  $t$  it is not possible for any process to arrive at any iteration  $k'$  before time  $t$ , where  $k' > k$ .

*Proof:* Suppose that some process arrives at iteration  $k' > k$  while no process arrives at iteration  $k$ . According to Lemma 9.6 some process should arrive at iteration  $k' - 1$  and the same for iteration  $k' - 2$  ( some process should arrive at iteration  $k' - 2$  ) and so on. So some process should also arrived already at iteration  $k$  also which is a contradiction with what we suppose. So it is not possible for any process to arrive at any iteration  $k'$  before time  $t$  while no process arrived at iteration  $k$  before time  $t$ , where  $k' > k$ . ■

*Lemma 9.8:* At any iteration  $k$  at least  $f + 1$  processes decide.

*Proof:* Since the first process  $p$ , arrives normally to iteration  $k$ , at line 33 according to lemma 9.5, it executed function *Deliver()* at line 46. It means the set of *Finished<sub>p</sub>* has cardinality of more than  $f$ , according to line 45. This set can increase its cardinality using only line 12. In this case it means  $p$  received *GOSSIP* messages either a) from processes in iteration more than  $k$  OR b) from processes in iteration  $k$  who have decided (by condition of line 11 variable  $d_p$  would be true). The first case is not possible because of corollary 9.7, because it means while no process arrives at iteration  $k$  at time  $t$ , some processes arrive at iteration  $k' > k$  before time  $t$ . So the only possible case is the case 2, which means at least  $f + 1$  processes decided at iteration  $k$ . ■

*Lemma 9.9:* At least one correct process decides at each iteration  $k$ .

*Proof:* It follows from lemma 9.8. ■

*Lemma 9.10:* All processes have the same *NextId* vector as long as they are in the same iteration  $k$ .

*Proof:* By induction

Base :  $k = 0$  : Trivial.

Induction step : Supposing all processes that are in iteration  $k$ , have the same *NextId* vector. We will show that all processes in iteration  $k + 1$ , also have the same *NextId* vector.

All the processes that are in iteration  $k + 1$ , already arrived at iteration  $k + 1$  in order to be able to propose for consensus at this iteration (at line 42). All processes who arrive at iteration  $k + 1$ , arrived either normally (line 33) or abnormally (line 25). So these processes can be divided

into three categories:

1 ) All those who arrived normally to iteration  $k + 1$ : It means they increased their  $k$  value at line 33. Because of agreement property of consensus we have  $k_p = k_q \rightarrow decision_p = decision_q$ . Now considering agreement property of consensus and what we suppose in induction step (all processes that decide at iteration  $k$  will have the same  $NextId$  vector immediately after deciding at line 43), we can conclude that all processes who arrive at iteration  $k + 1$  normally will update their  $NextId$  vector in the same way (either executing line 32 or not). It means all processes that arrive at iteration  $k + 1$  normally, will have the same  $NextId$  vector upon executing task sequencer at iteration  $k + 1$ .

2) All those who arrived abnormally using a SLOW message from a process who arrived normally at iteration  $k + 1$ : such a process will adopt its  $NextId$  vector from  $NextId$  vector of a process who arrive at this iteration normally using line 25. As we discussed in previous category all processes who arrive at iteration  $k + 1$  normally will have the same  $NextId$  vector upon arrival. So all processes who arrive at iteration  $k + 1$  abnormally using a SLOW message from a process who arrived normally at iteration  $k + 1$  using line 25 (such a process for sure exists as argued in the previous case.), will have the same  $NextId$  vector as those who arrive at iteration  $k + 1$  normally, upon executing task sequencer at iteration  $k + 1$ .

3) All those who come abnormally by receiving a SLOW message from a process who arrived abnormally at iteration  $k + 1$ : The sender of the SLOW message at iteration  $k + 1$ , belongs to the second category. So it have the same  $NextId$  vector as those who arrive at iteration  $k + 1$  normally. So all those who come abnormally using a SLOW message from such a process, will have the same  $NextId$  vector as those who arrive at iteration  $k + 1$  normally. If any other process receives a SLOW message from such processes, It can be said that it indirectly received a SLOW message from a process who arrived normally at iteration  $k + 1$ . So as the second category, they all will have the same  $NextId$  vector as those who arrive at iteration  $k + 1$  normally, upon executing task sequencer at iteration  $k + 1$ .

So all the processes who arrive at iteration  $k + 1$ , either they arrived normally or abnormally, have the same  $NextId$  vector upon executing task sequencer at iteration  $k + 1$ . Since processes in iteration  $k$  never manipulate their  $NextId$  vectors (definition 9.3), all processes will keep their  $NextId$  vectors the same while they are in iteration  $k$ . ■

**Lemma 9.11:** If for process  $p$  its  $\sum_{r \in \Pi} NextId_p[r]$  value becomes equal to  $\kappa$  at time  $t$  then it previously x-delivered  $\kappa - 1$  messages.

*Proof:*  $S = \sum_{r \in \Pi} NextId_p[r]$  can be increased in two lines : either line 25 or line 32. In Line 32,  $S$  will be incremented by 1 (line 32) while exactly one  $x - delivery$  happened before incrementing (line 33). In line 25,  $S$  will be incremented by  $K = \sum_{r \in \Pi} (N_q[r] -$

$NextId_p[r])$  (line 25), you can simply see  $\sum_{r \in \Pi} N_q[r] = \sum_{r \in \Pi} NextId_p[r] + \sum_{r \in \Pi} (N_q[r] - NextId_p[r])$  while exactly  $K x - delivery$  happened before incrementing. So with any increment in  $\sum_{r \in \Pi} NextId_p[r]$ , exactly the same number of  $x$ -delivery will happen for  $p$ . Since the initial value of  $\sum_{r \in \Pi} NextId_p[r]$  is 0 (line 4), whenever process  $p$  is in  $x$ -delivery position  $\kappa$ , it means  $\kappa - 1 = \sum_{r \in \Pi} NextId_p[r]$  (according to definition 9.13), and as explained above it means  $p$  x-delivered already exactly  $\kappa - 1$  messages. ■

**Definition 9.12:** Process  $p$  arrives at  $x$ -delivery position  $\kappa$  at time  $t$  if its  $\sum_{r \in \Pi} NextId_p[r]$  value becomes equal to  $\kappa$  using line 32 or 25 at time  $t$ . If it happens at line 32 we say process  $p$  arrives at  $x$ -delivery position  $\kappa$  normally, otherwise we say process  $p$  arrives at  $x$ -delivery position  $\kappa$  abnormally.

**Definition 9.13:** A process  $p$  is in  $x$ -delivery position  $\kappa$  at time  $t$  if it arrived at  $x$ -delivery position  $\kappa$  at time  $t'$  s.t.  $t' \leq t$  and did not arrive at any  $x$ -delivery position  $\kappa' > \kappa$  by time  $t$ .

**Definition 9.14:** Process  $p$  x-delivers  $m$  in  $x$ -delivery position  $\kappa$  if  $p$  x-delivers  $m$  while it is in  $x$ -delivery position  $\kappa$ .

**Definition 9.15:** For process  $p$ ,  $k$  is corresponding iteration of  $x$ -delivery position  $\kappa$  if  $p$  is in iteration  $k - 1$  when it arrives at  $x$ -delivery position  $\kappa$ .

**Lemma 9.16:** The first process that arrives at  $x$ -delivery position  $\kappa$ , arrives normally.

*Proof:* By contradiction suppose the first process, called  $p$ , arrives at  $x$ -delivery position  $\kappa$  abnormally using line 24 and 25. It means  $p$  has received a SLOW message (line 13) from some other process who was already in  $x$ -delivery position  $\kappa$  (line 25). In this we should have  $\kappa' = \sum_{r \in \Pi} N_q[r] > \sum_{r \in \Pi} NextId_p[r]$  because otherwise  $msgs - skipped$  will be equal to 0 and so no  $x$ -delivery will happen for  $p$  at line 24. It means the sending process already x-delivered  $\kappa' > \kappa$  (based on lemma 9.11). In this case  $p$  is not the first process who arrives at  $x$ -delivery position  $\kappa$ , which is a contradiction. ■

**Lemma 9.17:** For all  $\kappa \geq 1$ , if a process x-delivers  $m \neq \perp$  at  $x$ -delivery position  $\kappa$  then at least one correct process x-delivers at  $\kappa$  a normal message.

*Proof:* Based on lemma 9.16 the first process who arrived at  $x$ -delivery position  $\kappa$ , arrived normally (executed line 32). In corresponding iteration of  $x$ -delivery position  $\kappa$  (look at definition 9.15), which we will call it  $k$ , based on lemma 9.8 at least  $f + 1$  processes decided and having their *decided* variable equal to true (look at line 44). At least one of these processes is correct which we call it  $p$ . Now two cases are possible regarding  $p$ :

Case 1 ) Process  $p$  is the first process who arrived at  $x$ -delivery position  $\kappa$ , hence arrived normally (according to lemma 9.16). In this case lemma is proved because at least one correct process x-deliver at  $\kappa$  a normal message (at line 31). Note that if a process arrives normally at  $x$ -delivery position  $\kappa$ , it delivered at that  $x$ -delivery position a normal message (look at line 29).

Case 2 ) Process  $p$  is not the first process who arrived

at x-delivery position  $\kappa$ . So while  $p$  waiting at line 45 to gather enough messages in order to proceed to  $Deliver()$  function, some process  $q$ ,  $x - delivered(m \neq \perp)$  at x-delivery position  $\kappa$  already. Now consider 2 possible cases while  $p$  is in such a state:

Case 2-1) If  $p$  receives a SLOW message from a process in a higher iteration: In this case the Sequencer task of process  $p$  will be killed (at line 21), but since the  $Decided_p$  is true, it will execute  $Deliver()$  at line 22. Since  $q$  already  $x - delivered(m \neq \perp)$  at line 31,  $p$  will also  $x - deliver(m \neq \perp)$  at line 31. This is because of uniform agreement of consensus and Lemma 9.10, which causes the conditions of lines 29 and 30 be interpreted the same way for  $q$  and  $p$ . So the lemma is proved.

Case 2-2) If  $p$  never receives a SLOW message from a process in a higher iteration while waiting at line 45: We already supposed at least  $f + 1$  correct processes exist in the system (including  $p$ )<sup>7</sup>. These  $f + 1$  processes, in terms of iteration, can be either in (while  $p$  is waiting at line 45) :

Case 2-2-1) iteration higher than  $k_p$  : In this case all correct processes will GOSSIP their status to  $p$  which is waiting at line 45. all these gossip messages will be eventually received by  $p$  (fair-lossy property of links and correctness of both senders of GOSSIP and  $p$ ), so they all will cause  $p$  to increment cardinality of  $Finished_p$  set at line 12.

Case 2-2-2) iteration  $k_p$  and decided (their  $decided$  values become  $true$  using line 44): the same as previous case, they all (including  $p$  itself) will cause  $p$  to increment cardinality of  $Finished_p$  set at line 12.

Case 2-2-3) iteration  $k_p$  and not decided (their  $decided$  values is  $false$ ): If any of these processes are not waiting at line 39, It means that it will decide eventually according to termination property of consensus at line 43, and so it will have the same situation as those in Case 2-2-2. But if some of these processes are waiting at line 39, it means that in their  $Rcv$  vector, there is no non-bottom element which its id is equal to corresponding element in its  $NextId$  vector (Conditions of line 39). But we know that such an element exists for process  $p$  in its  $Rcv_p$  vector, because it already decided in iteration  $k_p$ , means that it was not blocked any more on line 39.  $Rcv_p$  vector will be GOSSIPed eventually to all those correct processes who are waiting at line 39 in iteration  $k_p$  by  $p$  (because of fair-lossy property of the links and correctness of the  $p$  and receivers). These waiting processes upon receiving GOSSIP message from  $p$  will execute line 14, because they all are in the same iteration as  $k_p$ . Since they are in iteration  $k_p$ , according to Lemma 9.10, they all have the same  $NextId$  vector as  $NextId_p$ . So  $Rcv_p$  will be adopted eventually by all such processes (using lines 15

<sup>7</sup>The number of correct processes is more than  $f$  meaning that the faulty processes are in minority in the system. In my opinion this condition is also necessary Continuity, not just for Validity, because if all processes who decided at line 43 crash immediately after ( $f$  processes), there should exist enough processes so that  $p$  does not remain blocked in line 45 forever, which means at least  $f + 1$  correct process is needed (including  $p$ ).

to 18). So none of these correct processes will remain blocked at line 39 in iteration  $k$  forever, since  $p$  was not blocked in that line, and now they all have the same  $Rcv$  and  $NextId$  vectors as  $Rcv_p$  and  $NextId_p$  respectively. So according to termination property of consensus they all will decide in line 43 and set their  $Decided$  value to true. From now on the case of 2-2-2 will apply to them and they all will cause  $p$  to increment cardinality of  $Finished_p$  set at line 12.

Case 2-2-4) iteration less than  $k_p$ : If these processes arrive at iteration  $k_p$  by their own means, then the case of 2-2-3 will apply to them. If these processes can not arrive at iteration  $k_p$ , they will receive a SLOW message from process  $p$  eventually (fair lossy links and correctness of  $p$  and receivers of SLOW message) and they all will execute line 25, which means they all will jump to iteration  $k_p$ , and they will *fork* the task *sequencer*, at line 27, in iteration  $k_p$ . From now on again the case of 2-2-3 will apply to them (correct processes in iteration  $k_p$  and not decided). At any of the above two cases, they all will cause  $p$  to increment cardinality of  $Finished_p$  set at line 12.

As you see  $p$  will eventually add all correct processes (which are at least  $f + 1$ ) to its  $Finished$  vector and will not remain blocked in line 45 forever and it will execute procedure  $Deliver()$  eventually. Since process  $q$  already  $x - delivers(m \neq \perp)$  it means that  $decision_p \neq \perp$  and is the same as  $decision_q$  (based on uniform agreement of consensus). On the other hand  $NextId_p$  is the same as the  $NextId_q$  (based on lemma 9.10). So conditions of lines 29 and 30 is true for process  $p$ , which means that  $p$  who is correct, will eventually  $x - deliver(m \neq \perp)$  and arrive at x-delivery position  $\kappa$ , so lemma is proved also in this case. ■

**Lemma 9.18:** For all  $\kappa \geq 1$ , if a process x-delivers  $\perp$  at x-delivery position  $\kappa$  then at least one correct process x-delivers at  $\kappa$  a normal message. (Continuity)

*Proof:* Suppose a process x-delivers  $\perp$  at x-delivery position  $\kappa$ . Similar to continuity. If a process x-delivers  $\perp$  at x-delivery position  $\kappa$  then some process x-deliver a normal message at x-delivery position  $\kappa$  (at least the first process who arrived at x-delivery position  $\kappa$ ). So based on Lemma 9.17 at least one correct process x-delivers a normal message at x-delivery position  $\kappa$ . ■

**Lemma 9.19:** If process  $p$  is in iteration  $k_p$  and process  $q$  is in iteration  $k_q$  so that  $k_p \leq k_q$ , then

a)  $\forall i \leq n : NextId_q[i] \geq NextId_p[i]$ .

b)  $p$  and  $q$  are in x-delivery positions  $\kappa_p$  and  $\kappa_q$  respectively in a way that  $\kappa_q \geq \kappa_p$ .

*Proof:*

a) by induction on  $k_p$  :

Base case : We suppose  $k_q = k_p$  which is trivial.

Induction step : Suppose process  $p$  is in iteration  $k_p$  and process  $q$  is in iteration  $k_q$  so that  $k_p < k_q$  and  $\forall i \leq n : NextId_q[i] \geq NextId_p[i]$ . We want to show that process  $q$  in iteration  $k_q + 1$  still have  $\forall i \leq n : NextId_q[i] \geq NextId_p[i]$ .

Consider two cases for process  $q$  and iteration  $k_q + 1$  :

Case 1) Process  $q$  arrive at iteration  $k_q + 1$  normally : It means that process  $q$  executed line 33 to arrive at iteration  $k_q + 1$ . In this case  $NextId$  vector of process  $q$  in iteration  $k_q + 1$  is bigger than or equal to  $NextId$  vector of process  $q$  in iteration  $k_q$  (line 32 either executed or not). So while process  $q$  in iteration  $k_q + 1$  we still have  $\forall i \leq n : NextId_p[i] \leq NextId_q[i]$ .

Case 2) Process  $q$  arrive at iteration  $k_q + 1$  abnormally : It means that process  $q$  executed line 25 to arrive at iteration  $k_q + 1$  which means it adopts directly or indirectly  $NextId$  vector of a process who arrived at iteration  $k_q + 1$  normally. So based on case 1, in this case also while process  $q$  in iteration  $k_q + 1$  we still have  $\forall i \leq n : NextId_q[i] \geq NextId_p[i]$ .

b ) Since part a is proved, based on definition 9.12 and 9.13 (co relation of  $\Sigma NextId$  and x-delivery position  $\kappa$ ), part b is also proved. ■

*Lemma 9.20:* For all  $\kappa \geq 1$ , if a correct process x-delivers  $m \neq \perp$  at x-delivery position  $\kappa$  then all correct processes x-deliver at  $\kappa$  a normal message or  $\perp$ .

*Proof:* Suppose a correct process, called  $p$ , x-delivers  $m \neq \perp$  at x-delivery position  $\kappa$  in corresponding iteration  $k_p$  and there exists at least  $f + 1$  correct processes in the system. These correct processes can be divided into two categories :

Category 1 ) Those who are in iteration more than  $k_p$  when  $p$  x-delivered  $m \neq \perp$  and arrived at x-delivery position  $\kappa_p$  : Based on lemma 9.19 all these correct processes are in x-delivery position more or equal to  $\kappa_p$ . According to lemma 9.11 they all x-delivered already more than  $\kappa_p$  messages which means that they all x-delivered a message (normal or  $\perp$ ) at x-delivery position  $\kappa_p$ .

Category 2 ) Those who are in iteration less or equal than  $k_p$  when  $p$  x-delivers  $m \neq \perp$  and arrives at x-delivery position  $\kappa_p$  : Consider following two cases :

Case 1 ) We show that these processes arrive at iteration  $k_p$  normally by their own means : These processes will x-deliver exactly the same number of messages that  $p$  already x-delivered. The reason is behind lemma 9.10 and agreement property of consensus, because in each iteration conditions of line 29 and 30 will be interpreted the same for all processes who execute  $deliver()$  function in any iteration, so they all will have be at the same x-delivery position as the process  $p$  in iteration  $k_p$  and x-delivering the same message in x-delivery position  $\kappa_p$  which is a normal message (based on lemma 9.11).

Case 2 ) In case that the first case is not possible, these processes arrive at iteration  $k_p$  eventually using a SLOW message from  $p$  or a process ahead of  $p$  in terms of  $k_p$ . Since  $p$  is a correct process, in case that there are no other correct processes in the system, they will eventually receive a SLOW message from  $p$  which will make them x-deliver  $\perp$  at x-delivery position  $\kappa_p$ .

As you see all correct processes will eventually x-deliver a message in x-delivery position  $\kappa_p$ . ■

*Lemma 9.21:* For all  $\kappa \geq 1$ , if some process x-delivers a normal message or  $\perp$  at x-delivery position  $\kappa$  then all

correct process x-deliver at  $\kappa$  a normal message or  $\perp$ . (Uniform Agreement)

*Proof:* From lemmas 9.18, 9.17 and 9.20. ■

*Lemma 9.22:*  $\forall p \forall q, Rcv_p[q]$  is a message which already xbcasted or  $\perp$ .

*Proof:* *Proof:*  $Rcv_p[q]$  is initially equal to  $\perp$  (look at line 4). The only place which is can obtain another value is line 9 which it will adopt an already xbcasted value. Note that at line 18 it will not obtain any other value other than an xbcast value or  $\perp$ . ■

*Lemma 9.23:* A process x-delivers a normal message  $m$  only if  $m$  was previously xbcast. (Uniform Integrity)

*Proof:* Suppose a process  $p$  x-delivers a normal message  $m$  at line 31. According to lines 29 and 30 it should be a non bottom decision of consensus at iteration  $k_p$ . Based on Validity of consensus and line 42,  $m$  is one of the elements of  $Rcv$  array of some process called  $q$ . The elements of this array for process  $q$  will take value either in line 9 or line 18. In line 9 it will be initialized with a non bottom value which is xbcast by  $q$ . In line 18 it will adopt a non bottom value which is already xbcasted by some process (based on lemma 9.22). At any case  $m$  was previously xbcast. ■

*Lemma 9.24:* If process  $p$  arrives at iteration  $k$  and process  $q$  arrive at iteration  $k'$  so that  $k = k'$ ,  $p$  and  $q$  both are in the same x-delivery position  $\kappa$  when they arrive at iterations  $k$  and  $k'$ .

*Proof:* According to lemma 9.10, processes who arrive at the same iteration, they all have the same  $NextId$  vector. So  $p$  and  $q$  both have the same  $NextId$  vector which means they already arrive at the same x-delivery position called  $\kappa$  (definition 9.12). It means they are both at x-delivery position  $\kappa$  when they arrive at iterations  $k$  and  $k'$  (definition 9.13). ■

*Lemma 9.25:* If process  $p$  arrives at x-delivery position  $\kappa$  normally and process  $q$  arrives at x-delivery position  $\kappa'$  normally so that  $\kappa = \kappa'$ ,  $p$  and  $q$  both are in the same iteration  $k$  when they arrive at x-delivery positions  $\kappa$  and  $\kappa'$ .

*Proof:* By contradiction suppose  $p$  and  $q$  are in iterations  $k_p$  and  $k_q$  ( $k_q > k_p$ ) respectively while x-delivering a normal message at the same x-delivery position. Now consider the following two cases for process  $q$  and iteration  $k_p + 1$  :

Case 1 ) process  $q$  already arrived at iteration  $k_p + 1$  (either normally or abnormally) : in this case according to lemma 9.10,  $q$  will have the same  $NextId$  vector when it arrives at iteration  $k_p + 1$  as  $NextId_p$  when  $p$  arrives at x-delivery position  $\kappa$  normally. In this case it is impossible that process  $q$  arrive at x-delivery position  $\kappa'$  normally so that  $\kappa = \kappa'$ , because the  $\Sigma NextId_q$  will be incremented at least by one when process  $q$  arrives at x-delivery position  $\kappa'$  normally and so  $\kappa \neq \kappa'$ .

Case 2 ) process  $q$  already jumped from iteration  $k_p + 1$  : It will happen when  $q$  receives an SLOW message from a process who is in an iteration more than  $k_p + 1$  and less than or equal to  $k_q$ . We want to show that  $q$  will be in x-delivery position  $\kappa$  upon receiving SLOW message

and before x-delivery at x-delivery positions  $\kappa'$ . Now we should show that the process who sent the SLOW message is in x-delivery position equal to or more than  $\kappa$  itself. It is proved in lemma 9.19. ■

*Lemma 9.26:* For all  $\kappa, \kappa' \geq 1$  if process  $p$  x-delivers at  $\kappa$  a normal message  $m$  and process  $q$  x-delivers at  $\kappa'$  a normal message  $m'$ , then  $\kappa = \kappa' \Leftrightarrow m = m'$ . (Uniform Total Order)

*Proof:* We divide the proof to two parts :

Part 1 )  $\kappa = \kappa' \Rightarrow m = m'$  : according to lemma 9.25 if  $\kappa = \kappa'$  then  $p$  and  $q$  both are in the same iteration  $k$  when they arrive at x-delivery positions  $\kappa$  and  $\kappa'$ . From now on according to Agreement property of consensus, all processes will decide the same value at the same iteration, so they all will x-deliver the same message.

Part 2 )  $m = m' \Rightarrow \kappa = \kappa'$  : By contadiction suppose  $\kappa \neq \kappa'$  but  $m = m'$  and both messages x-delivered. Suppose  $k$  and  $k'$  are the corresponding iterations of x-delivery positions  $\kappa$  and  $\kappa'$  respectively. By lemma 9.24 from  $\kappa \neq \kappa'$  we will have  $k \neq k'$ . Suppose  $m$  is the decision at consensus  $k$  and  $m'$  is the decision at consensus  $k'$ . Without affecting generality suppose  $k' > k$ . It means that if process  $p$  at iteration  $k$  x-delivers  $m$ , it will increase  $NextId_p[k \bmod |\Pi|]$  by one. It means that it is not possible to deliver  $m'$  at iteration  $k'$  which is a contradiction ( $k \bmod |\Pi| = k' \bmod |\Pi|$  and elements of  $NextId_p$  vector are non-decreasing so it is not possible that  $NextId_p[k \bmod |\Pi|] = m'.id$  anymore ). ■

*Lemma 9.27:* If there exists some message  $m = Rcv_p[p]$  which is xbcast by a correct process  $p$  but never be x-delivered then for any correct process  $q$ ,  $Rcv_q[p]$  will be equal to  $m$  eventually.

*Proof:* Consider two cases :

Case 1 )  $q = p$  : Trivial.

Case 2 )  $q \neq p$  : suppose  $p$  is in iteration  $k_p$  and  $q$  is in iteration  $k_q$ . Now consider three cases :

Case 2-1 )  $k_p = k_q$  : According to lemma 9.10 at this point  $p$  and  $q$  will have the same  $NextId$  vector. Since  $p$  and  $q$  are both correct, eventually  $q$  will receive GOSSIP message from  $p$ . At this point conditions of lines 16 and 17 are correct because  $Rcv_p[p] \neq \perp$  and  $Rcv_p[p].id = NextId_p[p] = NextId_q[p]$ . So  $Rcv_q[p]$  will be equal to  $Rcv_p[p] = m$ .

Case 2-2 )  $k_p > k_q$  : In this case process  $q$  will arrive at iteration  $k_p$  either by its own means or eventually will receive SLOW message from process  $p$ . From now on it can be reasoned according to case 1.

Case 2-3 )  $k_p < k_q$  : The same as previous case, just reverse  $p$  and  $q$ . ■

*Lemma 9.28:* If there exists some message  $m = Rcv_p[p]$  which is xbcast by a correct process  $p$  but never be x-delivered then no correct process will be blocked forever at lines 39, 43 and 45 in sequencer task.

*Proof:* We divide proof in 3 parts, one for each line :

Line 39 : Suppose correct process  $q$  is blocked in line

39 in iteration  $k_q$ . Now suppose process  $p$ , which will not be blocked in line 39, is in iteration  $k_p$ . Now consider three cases regarding  $k_p$  and  $k_q$  :

Case 1 )  $k_p = k_q$  : According to lemma 9.10 at this point  $p$  and  $q$  will have the same  $NextId$  vector. Since  $p$  and  $q$  are both correct, eventually  $Rcv_q[p]$  will be equal to  $Rcv_p[p]$  (based on lemma 9.27) and  $q$  will not be blocked at line 39 any more because  $Rcv_q[p] \neq \perp$  and  $Rcv_q[p].id = NextId_p[p]$ .

Case 2 )  $k_p > k_q$  : In this case process  $q$  will arrive at iteration  $k_p$  either by its own means or eventually will receive SLOW message from process  $p$ . From now on it can be reasoned according to case 1.

Case 3 )  $k_p < k_q$  : The same as previous case, just reverse  $p$  and  $q$ .

Line 43 : No process will remain blocked in this line because of Termination property of consensus.

Line 45 : Suppose correct process  $q$  is blocked in line 39 in iteration  $k_q$ . We know that there exist  $f + 1$  correct processes in the system. We will categorize these correct processes in three groups and will show that eventually all these three processes will increment cardinality of set  $Finished_q$  by one, making  $q$  unblocked at line 45. We categorize correct processes based on the iteration  $k$  in which they are currently in :

Category 1 :  $k = k_q$  : According to what we said about eventual unblocking in lines 43 and 45, which we showed before, all processes will decide eventually and then all these processes will increment cardinality of set  $Finished_q$  by one eventually because upon receiving GOSSIP message from them by  $q$ , Condition of line 11 is always correct.

Category 2 :  $k > k_q$  : These processes all will increment cardinality of set  $Finished_q$  by one eventually because upon receiving GOSSIP message from them by  $q$ , Condition of line 11 is always correct.

Category 3 :  $k < k_q$  : In this case process  $p$  will arrive at iteration  $k_q$  either by its own means or eventually will receive SLOW message from process  $q$ . From now on it can be reasoned according to category 1. ■

*Lemma 9.29:* If there exists some message  $m = Rcv_p[p]$  which is xbcast by a correct process  $p$  but never be x-delivered then  $Rcv_p[p]$  will be decided eventually.

*Proof:* According to lemma 9.27 for all processes like  $q$ , we will have  $Rcv_q[p] = Rcv_p[p]$  eventually. After this moment since according to Lemma 9.28, no process will be blocked in Sequencer task, rotating sender strategy will eventually elect process  $p$  as the rotating sender in line 41. So  $Rcv_q[p] = m$  will be proposed in line 42, regardless of who process  $q$  is, and  $m$  will be decided eventually at line 43 (because of validity and termination of consensus). ■

*Lemma 9.30:* If a correct process  $p$  xbcast message  $m$  then some correct process like  $q$  eventually x-delivers  $m$ . (Validity)

*Proof:* By induction on  $id_p$  ( $id_p$  is non decreasing, look at line 7, so we can define induction over it)

Base :  $m.id = 0$  : According to Lemma 9.29  $m$  with  $m.id = 0$  will be eventually decided by some process  $q$ . From now on process  $q$  will not be blocked in line 45 according to lemma 9.28 and  $q$  will execute deliver function. Since  $decision_q \neq \perp$  and  $decision_q.id = NextId_q[p] = NextId_p[p] = 0$ , conditions of lines 29 and 30 are correct and  $m$  will be x-delivered.

Induction step : Suppose correct process  $p$  xbcast message  $m$  with  $m.id = id$  then correct process  $q$  eventually x-delivers  $m$  and process  $p$  xbcast message  $m'$  with  $m'.id = id+1$ . We want to show that some correct process  $s$  will eventually x-deliver  $m'$ .

Proof : Suppose by contradiction  $m'$  will never be x-delivered. Consider two cases :

Case 1  $p = q$  : In this case upon x-delivery of  $m$ , we will have  $NextId_p[p] = NextId_p[p] + 1$ . So process  $p$  will not remain blocked on line 8, so at line 9 we will have  $Rcv_p[p] = m'$  with id equal to  $m.id + 1$ . According to lemma 9.29,  $m'$  will be decided eventually at line 43 by process some process  $s$ . From now on process  $s$  will not be blocked in line 45 according to lemma 9.28 and  $s$  will execute deliver function. Since  $p$  and  $s$  are in the same iteration, according to lemma 9.10 we have  $NextId_s[p] = NextId_p[p] = Rcv_p[p].id = m'.id$  and condition of line 30,  $NextId_s[p] = NextId_p[p] = Rcv_p[p].id = m'.id$ , will be true for process  $s$ , so x-delivery of  $m'$  will happen at process  $s$  eventually which is a contradiction.

Case 2  $p \neq q$  : In this case upon x-delivery of  $m$ , we will have  $NextId_q[p] = NextId_q[p] + 1$ . Now we will show that how it will make process  $p$  not to be blocked for ever in line 8 by making  $NextId_p[p] = NextId_p[p] + 1$ . Suppose process  $p$  is blocked at line 8 after xbcasting  $m'$  which  $m'.id = m.id + 1$ . Based on lemma 9.19  $k_q > k_p$  and since  $p$  and  $q$  both are correct,  $p$  will adopt the  $NextId_q$  at line 25 by receiving SLOW message from  $q$  eventually. So  $p$  will not remain blocked anymore at line 8. The rest of the proof is similar to case 1. ■