

Early-Delivery Dynamic Atomic Broadcast

Ziv Bar-Joseph Idit Keidar
Nancy Lynch

MIT Laboratory for Computer Science

`zivbj@mit.edu`, `idish@theory.lcs.mit.edu`, `lynch@theory.lcs.mit.edu`

April 17, 2002

Abstract

We consider a problem of atomic broadcast in a dynamic setting where processes may join, leave voluntarily, or fail (by stopping) during the course of computation. We provide a formal definition of the *Dynamic Atomic Broadcast* problem and present and analyze a new algorithm for its solution in a synchronous system, where processes have approximately synchronized clocks.

Our algorithm exhibits constant message delivery latency in the absence of failures, even during periods when participants join or leave. To the best of our knowledge, this is the first algorithm for totally ordered multicast in a dynamic setting to achieve constant latency bounds in the presence of joins and leaves. When failures occur, the latency bound is linear in the number of actual failures.

Our algorithm uses a solution to a variation on the standard distributed consensus problem, in which participants do not know a priori who the other participants are. We define the new problem, which we call *Consensus with Unknown Participants*, and give an early-deciding algorithm to solve it.

1 Introduction

We consider a problem of atomic broadcast in a *dynamic setting* where an unbounded number of participants may join, leave voluntarily, or fail (by stopping) during the course of computation. We formally define the *Dynamic Atomic Broadcast (DAB)* problem, which is an extension of the Atomic Broadcast problem [17] to a setting with infinitely many processes, any finite subset of which can participate at a given time. Just as Atomic Broadcast is a basic building block for state machine replication in a static setting [20, 27], DAB can serve as a building block for state machine replication among a dynamic set of processes.

We present and analyze a new algorithm, which we call *Atom*, for solving the DAB problem in a synchronous crash failure model. Specifically, we assume that the processes solving DAB have access to approximately-synchronized local clocks and to a lower-level dynamic network that guarantees timely message delivery between currently active processes. The challenge is to guarantee consistency among the sequences of messages delivered to different participants, while still achieving timely delivery, even in the presence of joins and leaves.

Atom exhibits *constant* message delivery latency in the absence of failures, *even during periods when participants join or leave*; this is in contrast to previous algorithms solving similar problems in the context of view-oriented group communication, e.g., [1, 9]. When failures occur, *Atom*'s latency bound is linear in the number of failures that actually occur; it does not depend on the number of potential failures, nor on the number of joins and leaves that occur.

A key difficulty for an algorithm solving DAB is that when a process fails, the network does not guarantee that the surviving processes all receive the same messages from the failed process. But the strong consistency requirements of DAB dictate that processes agree on which messages they deliver to their clients. The processes carry out a protocol to coordinate message delivery, which works roughly as follows: Each *Atom* process divides time into *slots*, using its local clock, and assigns each message sent by its client to a slot. Each process delivers messages to its client in order of slots, and within each slot, in order of sender identifiers. Each process determines the *membership* of each slot, and delivers messages only from senders that it considers to be members of the slot. To ensure consistency, the processes must agree on the membership of each slot.

Processes joining (or voluntarily leaving) the service coordinate their own join (or leave) by selecting a join-slot (or leave-slot) and informing the other processes of this choice, without delaying the normal delivery of messages. When a process fails, *Atom* uses a novel *distributed consensus service* to agree upon the slot in which it fails. The consensus service required by *Atom* differs from the standard stopping-failure consensus services studied in the distributed algorithms literature (see, e.g., [21]) in that the processes implementing the consensus service do not know a priori who the other participants are. *Atom* tracks process joins and leaves, and uses this information to approximate the active set of processes that should participate in consensus. However, different processes running *Atom* may have somewhat different perceptions of the active set, e.g., when a participant joins or leaves *Atom* at roughly the time consensus is initiated.

In order to address such uncertainties, we define a new consensus service, *consensus with unknown participants (CUP)*. When a process i initiates CUP, it submits to CUP a finite set W_i estimating the current world, in addition to i 's proposed initial consensus value v_i . The worlds suggested by different participants do not have to be identical, but some restrictions are imposed on their consistency. Consider, e.g., the case that process k joins *Atom* at roughly the time CUP is initiated. One initiator, i , may think that k has joined in time to participate and include k in W_i , while another, j , may exclude k from W_j . Process k cannot participate in the CUP algorithm in the usual way, because j would not take its value into account. On the other hand, if k does not

participate at all, i could block, waiting forever for a message from k . We address such situations by allowing k to explicitly *abstain* from an instance of CUP, i.e., to participate without providing an input. A service that uses CUP must ensure that for every i , (1) W_i includes all the processes that ever initiate this instance of CUP (unless they fail or leave prior to i 's initiation); and (2) if $j \in W_i$, (and neither i nor j fail or leave), then j participates in CUP either by initiating or by abstaining. Thus, W_i sets can differ only in the inclusion of processes that abstain, leave, or fail.

Note that once an instance of CUP has been started, no new processes (that are not included in W_i) can join the running instance. Nevertheless, CUP provides a good abstraction for solving DAB, because Atom can invoke multiple instances of CUP with different sets of participants.

We give an early-deciding algorithm to solve CUP in a fail-stop model [26], that is, in a time-free crash failure model where processes are equipped with *perfect failure detectors* [5]. The failure detector is external to CUP; it is implemented by Atom. CUP uses a strategy similar to previous early-deciding algorithms for consensus with a predetermined set of participants [13], but it also tolerates uncertainty about the set of participants, and moreover, it allows processes to leave voluntarily without incurring additional delays. The time required to reach consensus is linear in the number of failures that actually occur during an execution, and does not depend on an upper bound on the number of potential failures, nor on the number of processes that leave.

We also analyze the message-delivery latency of Atom under different failure assumptions. We show a constant latency bound for periods when no failures occur, even if joins and leaves occur. When failures occur, the latency is proportional to the number of actual failures. This is inevitable: atomic broadcast requires a number of rounds that is linear in the number of failures (see [2]).

We envision a service using Atom, or a variation of it, deployed in a large LAN, where latency is predictable and message loss is bounded. In such settings, a network with the properties we assume can be implemented using forward error correction (see [3]), or retransmissions (see [28]). The algorithm can be extended for use in environments with looser time guarantees, e.g., networks with differentiated services; we outline ideas for such an extension in Section 7.7.

In summary, this paper makes the following main contributions: (1) the definitions of two new problems for dynamic networks, expressed by the DAB and CUP services; (2) an early-delivery DAB algorithm, Atom, which exhibits constant latency in the absence failures; (3) a new early-deciding algorithm for solving CUP in a fail-stop model; and (4) the analysis of Atom's message-delivery latency under various failure assumptions.

The rest of this paper is organized as follows: Section 2 discusses related work. In Section 3, we specify the DAB service. In Section 4 we specify CUP and in Section 5, we present the CUP algorithm and its analysis. We then turn to the presentation of Atom: Section 6 specifies the environment and model assumptions for Atom, and Section 7 contains a detailed presentation of the Atom algorithm and its analysis. Section 8 concludes the paper. The Appendix contains rigorous correctness proofs for both CUP and Atom.

2 Related Work

A dynamic universe, where processes join and leave, was first considered in the context of view-oriented group communication work [7], pioneered by the Isis [4] system. The first analysis of time bounds of message delivery in synchronous group communication systems was performed by Cristian [9]. Our service resembles the services provided by group communication systems; although we do not export membership to the application, it is computed, and would be easy to export.

View-oriented group communication systems, including systems designed for synchronous systems and real-time applications (e.g., Cristian's [9], xAMp [25], and RTCAST [1]), generally run

a group membership protocol every time a process joins or leaves, and therefore delay message delivery to all processes when joins or leaves occur. Cristian’s system uses an atomic broadcast primitive to agree upon group membership. Since, unlike CUP, the atomic broadcast service works with a static universe, a process join has to be agreed upon before any new membership change is handled (voluntary leaves are not considered). Therefore, Cristian’s service exhibits constant latency only in periods in which no joins or failures occur. Latency during periods with multiple joins is not analyzed. xAMP is a group communication system supporting a variety of communication primitives for real-time applications. The presentation of xAMP in [25] focuses on the various communication primitives and assumes that a membership service is given. The delays due to failures and joins are incurred in the membership part, which is not described or analyzed. RTCAST is a real-time group communication system, for which a detailed analysis of membership latency was conducted [1]. The latency bound achieved by RTCAST is *linear* in the number of processes, even when no process fails, due to the use of a logical ring. Moreover, RTCAST makes stronger assumptions about its underlying network than we do – it uses an underlying reliable broadcast service that guarantees that correct processes deliver the same messages from faulty ones; the cost of this primitive is not considered in the analysis.

Some group membership services avoid running the full-scale membership for join and leaves by using light-weight group membership [15] services; they use an atomic broadcast service to disseminate join and leave messages in a consistent manner, without running the full-scale group membership algorithm. However, unlike our CUP service, the atomic broadcast service such systems use do not tolerate uncertainty about the set of participants. Therefore, a race condition between a join and a concurrent failure can cause such light-weight group services (e.g., [23, 12, 15]) to violate the semantics of the underlying heavy-weight membership services. Those light-weight group services that do preserve the underlying heavy-weight membership semantics (e.g. [24]), do incur extra delivery latencies whenever joins and leaves occur.

Other work on group membership in synchronous and real-time systems, e.g., [19, 18] has focused on membership maintenance in a static, fairly small, group of processes, where processes are subject to failures but no new processes can join the system. Likewise, work analyzing time bounds of synchronous atomic broadcast, e.g. [16, 10, 8], considered a static universe, where processes could fail but not join. Thus, this work did not consider the DAB problem.

In a previous paper [3], we considered a simpler problem of dynamic totally ordered broadcast without all or nothing semantics. For this problem, the linear lower bound does not apply, and we exhibited an algorithm that solves the problem in constant time even in the presence of failures.

Recent work [22, 6] considers different services, including (one shot) consensus, for infinitely many processes in asynchronous shared memory models. Chockler and Malkhi [6] present a consensus algorithm for infinitely many processes using a *static* set of active disks, a minority of which can fail. This differs from the model considered here, as in our model all system components may be ephemeral. Merritt and Taubenfeld [22] study consensus under different concurrency models; in their terminology, our model assumes *unbounded* congruency and $[1, \infty]$ -participation, which means that at least one process must participate and there is no bound on the number of participants. They show that with these assumptions, in an asynchronous shared memory model, infinitely many bits are required in order to solve consensus. The algorithms they give are not fault tolerant (they tolerate only initial failures). To the best of our knowledge, atomic broadcast has not been considered in a similar context. Moreover, these problems were not considered in message-passing models, and it is not clear that a canonical transformation from the shared memory model the message-passing model applies to a setting with infinitely many processes.

3 Dynamic Atomic Broadcast Service Specification

We now present the DAB service specification. Our universe consists of an infinite ordered set of endpoints, I . The specification of DAB is parameterized by a message alphabet, M . The signature of the DAB(M) service is presented in Figure 1.

Input:

$\text{join}_i, \text{leave}_i, \text{fail}_i, i \in I$
 $\text{mcast}_i(m), m \in M, i \in I$

Output:

$\text{join_OK}_i, \text{leave_OK}_i, i \in I$
 $\text{rcv}_i(m), m \in M, i \in I$

Figure 1: The signature of the DAB(M) service.

We do not consider recoveries from failure or rejoining after leaving. In other words, there cannot be multiple “incarnations” at a single endpoint. Instead of new incarnations, consider the same client joining at new endpoints.

Assumptions about the application: DAB(M) assumes that its application satisfies the following safety conditions:

- For each $i \in I$:
 - At most one join_i and at most one leave_i occur.
 - If leave_i occurs, then it is preceded by join_OK_i .
 - Any $\text{mcast}_i(m)$ has a preceding join_OK_i but no preceding leave_i or fail_i .
- At most one $\text{mcast}(m)$ occurs for each particular m .

DAB guarantees: Given an application that satisfies the above constraints, DAB(M) satisfies the properties we now specify.

We first specify some basic integrity properties, both safety and liveness. We later specify the properties related to the ordering and reliability of messages.

Basic safety properties:

- *Join/leave integrity:* For each i :
 - At most one join_OK_i and at most one leave_OK_i occur.
 - If join_OK_i occurs then it is preceded by join_i .
 - If leave_OK_i occurs then it is preceded by leave_i .
- *Message integrity:*
 - No two $\text{rcv}_j(m)$ actions occur for the same m and j .
 - If $\text{rcv}_j(m)$ occurs for some j then it is preceded by $\text{mcast}_i(m)$ for some i .

Basic liveness properties:

- *Eventual join*: If join_i occurs then either fail_i or join_OK_i occurs.
- *Eventual leave*: If leave_i occurs then either fail_i or leave_OK_i occurs.

To specify the ordering and reliability guarantees of DAB, we require that there be a total ordering \mathcal{S} on all the messages received by any of the endpoints, such that for all $i \in I$, the following properties are satisfied.

Safety properties:

- *Multicast order*: If $\text{mcast}_i(m)$ occurs before $\text{mcast}_i(m')$, then m precedes m' in \mathcal{S} .
- *Receive order*: If $\text{rcv}_i(m)$ occurs before $\text{rcv}_i(m')$ then m precedes m' in \mathcal{S} .
- *Multicast gap-freedom*: If $\text{mcast}_i(m)$, $\text{mcast}_i(m')$, and $\text{mcast}_i(m'')$ occur, in that order, and S contains m and m'' , then S also contains m' .
- *Receive gap-freedom*: If S contains m , m' , and m'' , in that order, and $\text{rcv}_i(m)$ and $\text{rcv}_i(m'')$ occur, then $\text{rcv}_i(m')$ also occurs.

Liveness property:

- *Multicast liveness*: If $\text{mcast}_i(m)$ occurs and no fail_i occurs, then S contains m .
- *Receive liveness*: If S contains m , m is sent by i and i does not leave or fail, then $\text{rcv}_i(m)$ occurs, and for every m' that follows m in S , $\text{rcv}_i(m')$ also occurs.

4 Consensus with Unknown Participants – Specification

In this section we define the problem of Consensus with Unknown Participants (CUP). CUP is an adaptation of the problem of fail-stop uniform consensus to a dynamic setting in which the set of participants is not known ahead of time, and in which participants can leave the algorithm voluntarily after initiating it. Moreover, participants are not assumed to initiate at the same time. CUP uses an underlying reliable network, and a perfect failure detector.

We begin with a description of CUP’s external signature (interface). We then specify the assumptions that CUP makes about its environment, including the application, the underlying network, and the external failure detector. We separate these into *safety* and *liveness* assumptions. Finally, we specify CUP’s safety and liveness guarantees. CUP’s safety guarantees depend on only the safety assumptions, that is, they are not allowed to be violated even if the liveness assumptions do not hold. On the other hand, CUP’s liveness guarantees depend on both the safety and liveness assumptions.

4.1 External Signature

The CUP specification uses the following data types:

- I , an infinite ordered set of *endpoints*. Each endpoint in I corresponds to a potential participant in CUP.
- V , a totally ordered set of *values*. Initial values and decision values are elements of V .

Input:

```

initi(v,W), v ∈ V, W ⊆ I, W finite, i ∈ I // i initiates with value v, world W
abstaini, i ∈ I // i abstains
net_rcvi(m), m ∈ MCUP, i ∈ I // i receives message m
leavei, i ∈ I // i leaves
leave_detecti(j), j, i ∈ I // i detects that j has left
faili, i ∈ I // i fails
fail_detecti(j), j, i ∈ I // i detects that j has failed

```

Output:

```

decidei(v), v ∈ V, i ∈ I // i decides on value v
net_mcasti(m), m ∈ MCUP, i ∈ I // i multicasts m

```

Figure 2: The signature of CUP.

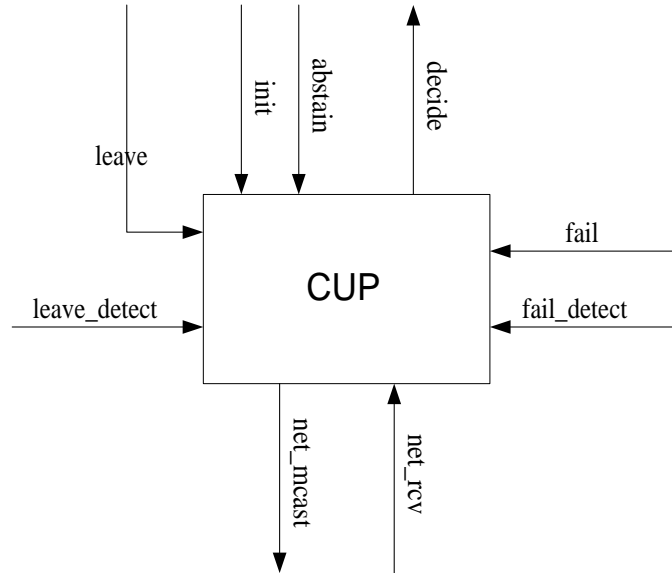


Figure 3: Interface diagram for CUP.

- M_{CUP} , a message alphabet.

The external signature of CUP is presented in Figure 2, and depicted in Figure 3.

The interface describes four kinds of interaction: “normal” interaction with clients of the CUP service, interaction with a multicast network, communication involving leaves and leave detection, and communication involving failures and failure detection.

Normal interaction with clients: A process may participate in the CUP service in two ways: it may provide an initial value, in which case we say that the process *initiates* CUP, or it may decline to provide an initial value, in which case we say that it *abstains*. Participant $i \in I$ initiates CUP using the $init_i(v,W)$ action. Here, v is i 's initial value, and W is its initial *world*, that is, the set of processes that i expects to participate in CUP. Participant i abstains using the $abstain_i$ action. Informally speaking, a participant abstains when it does not need to participate in CUP, but because of uncertainty about CUP participants, some other participant may expect it to participate.

An environment assumption ensures that, if any process expects i to participate in CUP, i will in fact participate, unless it leaves or fails. CUP reports the consensus decision value to process i using the `decidei(v)` action.

Multicast network: The network interface consists of the `net_mcast` and `net_rcv` actions.

Leaves: A participant can leave the CUP service voluntarily using the `leavei` action. We assume that the environment provides a *leave detector*: the `leave_detecti(j)` action is used to notify i that j has left the algorithm voluntarily.

Failures: The `faili` action represents the failure of endpoint i . We assume that the environment provides a failure detector, which uses the `fail_detecti(j)` action to notify i that j has failed.

4.2 Environment Assumptions

Here we list and explain the assumptions that CUP makes about its environment. We classify these as safety and liveness assumptions. Formally, each of the properties given here is a *trace property* ([21, Ch. 8]).

4.2.1 Safety assumptions

The first assumption expresses simple well-formedness conditions saying that each participant begins participating (by initiating or abstaining) at most once, leaves at most once, and fails at most once.

- *Well-formedness:* For any $i \in I$,
 1. At most one `initi` or `abstaini` event¹ occurs.
 2. At most one `leavei` event occurs.
 3. At most one `faili` event occurs.
 4. No `leavei` or `faili` precedes an `initi`.

The next assumption says that, while the worlds W suggested by different participants in their `init` events do not have to be identical, CUP's environment must guarantee that they have a certain kind of consistency. Namely, each W set submitted by an initiating participant i must include all participants that ever initiate CUP and that do not leave or fail prior to the `initi` event. This implies that every participant must be included in its own estimated world.

- *World consistency:* If `initi(*, W)` and `initj(*,*)` events occur, then either $j \in W$, or a `leavej` or `failj` event occurs before the `initi(*, W)` event.

The next property describes the correctness of the message deliveries: every message that is received was previously sent, and no message is received at the same location more than once. Moreover, the order of message receipt between particular senders and receivers is FIFO.

¹An “event” is an occurrence of an action in a sequence.

- *Message integrity*: There is a mapping from `net_rcv` events to preceding `net_mcast` events, such that the same message in M_{CUP} appears in both events, and such that no two `net_rcvi` events for the same i map to the same `net_mcast` event. Moreover, two `net_rcvi` events that map to `net_mcast` events of the same sender occur in the same order as the `net_mcast` events.

The next two properties describe assumptions about leaves and leave detection. The first says that leave detection is “accurate”, in the sense that the occurrence of a `leave_detecti(j)` implies that j has really left; it also includes a simple well-formedness condition. The second property says that leaves are handled gracefully, in the sense that the occurrence of a `leave_detecti(j)` implies that i has already received any network messages sent by j prior to leaving. Thus, a `leave_detecti(j)` is an indication that i has not lost any messages from j .

- *Accurate leave detector*: For any $i, j \in I$, at most one `leave_detecti(j)` event occurs, and if `leave_detecti(j)` occurs, then it is preceded by a `leavej`.
- *Lossless leave*: Assume `net_mcastj(m)` occurs and is followed by a `leavej`. Then if a `leave_detecti(j)` occurs, it is preceded by `net_rcvi(m)`.

The final safety assumption says that failure detection is accurate.

- *Accurate failure detector*: For any $i, j \in I$, at most one `fail_detecti(j)` event occurs, and if `fail_detecti(j)` occurs, then it is preceded by a `failj`.

Note that we do not have a failure assumption analogous to the lossless leave property; thus, failures are different from leaves in that we allow the possibility that some messages from failed processes may be lost.

4.2.2 Liveness assumptions

The first liveness assumption says that, if any process i expects another process j to participate, then j will actually do so, unless either i or j leaves or fails.

- *Init occurrence*: If an `initi(*,W)` event occurs and $j \in W$, then an `initj`, `abstainj`, `leavei`, `faili`, `leavej`, or `failj` occurs.

The next assumption describes reliability of message delivery. It says that any message that is multicast by a non-failing participant that belongs to any of the W sets submitted to CUP, is received by all the non-leaving, non-failing members of all those W sets.

- *Reliable delivery*: Define $U = \cup_{k \in I} \{ W \mid \text{init}_k(*, W) \text{ occurs} \}$. If $i, j \in U$ and `net_mcasti(m)` occurs after an `initi` or `abstaini` event, then a `net_rcvj(m)`, `leavej`, `faili`, or `failj` occurs.

The final liveness assumption says that the leaving or failure of any process that belongs to an initiator’s W set is detected by that initiator, unless it finishes by deciding, leaving, or failing.

- *Complete leave and failure detector*: If `initi(*,W)` occurs, $j \in W$, and `leavej` or `failj` occurs, then `fail_detecti(j)`, `leave_detecti(j)`, `decidei`, `leavei`, or `faili` occurs.

4.3 CUP Service Guarantees

Now we list CUP’s service guarantees. Again, we classify these as safety and liveness properties. As we noted earlier, CUP’s safety guarantees depend only on its safety assumptions, whereas CUP’s liveness guarantees depend on both its safety and liveness assumptions.

Formally, each individual property is a trace property. The complete specification consists of two general trace properties whose respective sets of traces are defined by the following predicates:

1. The conjunction of all the CUP safety assumptions implies all the CUP safety guarantees.
2. The conjunction of all the CUP safety and liveness assumptions implies all the CUP liveness guarantees.

4.3.1 Safety guarantees

The first guarantee expresses well-formedness conditions saying that only participants that have initiated can decide, and each participant decides at most once.

- *Well-formedness:* For any $i \in I$,
 1. If `decidei` occurs then it is preceded by an `initi`.
 2. At most one `decidei` occurs.

The next two guarantees are the main agreement and validity guarantees for consensus. The uniform agreement property says that everyone who decides agrees. The validity property has two parts: it says that any decision value is some participant’s initial value, and moreover, that any participant’s decision is no greater than its initial value. The latter is not a “standard” property for consensus but is needed for our use in Atom.

- *Uniform Agreement:* For any $i, j \in I$, if `decidei(v)` and `decidej(v')` both occur then $v = v'$.
- *Validity:* For any $i \in I$, if `decidei(v)` occurs then
 1. For some j , `initj(v, *)` occurs.
 2. If `initi(v', *)` occurs then $v \leq v'$.

4.3.2 Liveness guarantees

CUP provides one liveness guarantee, which says that any participant that initiates and neither leaves nor fails must eventually decide. We do not make such a guarantee for a participant that abstains, that is, participants that abstain need not be informed of the decision value.

- *Termination:* If an `initi` event occurs then a `decidei`, `leavei`, or `faili` occurs.

5 The CUP Algorithm

In this section, we present our implementation of CUP.

5.1 Modeling Assumptions and Conventions

We use the I/O automaton model of Lynch and Tuttle (see, e.g., [21, Ch. 8]), using standard precondition/effect (guarded command) pseudo-code, augmented with one new construct: effects may include statements of the form `trigger(a)`, where `a` is an output action. Formally, we assume the automaton's state contains a special FIFO buffer `trigger-buffer`. The `trigger(a)` statement adds `a` to the end of `trigger-buffer`. The action at the head of `trigger-buffer` is always enabled, and gets removed from `trigger-buffer` when it is performed. No other state changes are associated with action `a`.

The `faili` action described in the CUP interface represents the failure of endpoint i . In terms of the algorithm, we interpret this to mean that once `faili` occurs, i performs no more locally controlled actions, and input actions have no effect on the state. We treat this as a general convention, and do not include event handlers for `faili` actions in our pseudo-code.

5.2 The Algorithm

Figures 4 and 5 contain the CUP implementation for a particular endpoint $i \in I$. The algorithm includes no internal actions. Therefore, the signature consists of the actions indexed by this particular i in the external signature of CUP (see Section 4). The message alphabet M_{CUP} is specialized to the set of messages of the following forms:

- (i, r, v, W) , where $i \in I$, $r \in \mathbb{N}$, $v \in V$, and W is a finite subset of I .
- (i, OUT, r) , where $i \in I$ and $r \in \mathbb{N}$,

The algorithm proceeds in asynchronous rounds numbered $1, 2, \dots$. In each round, a process sends its current estimates of the value and the world (the set of active processes) to the other processes. Each process maintains two-dimensional arrays, `value` and `world`, in which it collects the value and world information it receives from all processes in all rounds. It records, in a variable `out[r]`, the other processes that it knows will not participate in round r because they have previously left, abstained, or decided. It also records, in a variable `failed`, the processes that it knows have failed.

```

mode  $\in \{ \perp, \text{running}, \text{done} \}$ , initially  $\perp$ 
round  $\in \mathbb{N}$ , initially 0
for each  $r \in \mathbb{N}^+$ ,  $j \in I$ :
  value[r,j]  $\in V \cup \{ \perp \}$ , initially  $\perp$ 
  world[r,j], a finite subset of  $I$  or  $\perp$ , initially  $\perp$ 
for each  $r \in \mathbb{N}^+$ 
  out[r], a finite subset of  $I$ , initially  $\{ \}$ 
  failed[r], a finite subset of  $I$ , initially  $\{ \}$ 

Derived variables:

for each  $r \in \mathbb{N}^+$ 
  out-by[r], a finite subset of  $I$ , defined as  $\cup_{r' \leq r} \text{out}[r']$ 
  failed-by[r], a finite subset of  $I$ , defined as  $\cup_{r' \leq r} \text{failed}[r']$ 

```

Figure 4: CUP _{i} state.

```

initi(v,W)
Eff: if mode = ⊥ then
    mode ← running
    round ← 1
    trigger(net_mcasti(i,1,v,W))

net_mcasti(i,r,v,W) where r ≥ 2
Pre: mode = running
    r = round + 1
    W = world[round,i] \ out[round] \ failed[round]
    // All messages for the previous round have been received.
    ∀ j ∈ W: value[round,j] ≠ ⊥
    W ≠ {} ∧ v = min{value[round,j] | j ∈ W}
    // No decision can be made.
    ¬ ∀ j ∈ world[round,i] \ out[round]:
        value[round,j] = value[round,i] ∧ world[round,j] ⊆ world[round,i]
Eff: round ← r

net_rcvi(j,r,v,W)
Eff: if mode ≠ done ∧ j ∉ failed-by[r] then
    value[r,j] ← v
    world[r,j] ← W

abstaini
Eff: if mode = ⊥ then
    mode ← done
    trigger(net_mcasti(i,OUT))

decidei(v)
Pre: mode = running
    value[round,i] ≠ ⊥
    ∀ j ∈ world[round,i] \ out[round]:
        value[round,j] = v ∧ world[round,j] ⊆ world[round,i]
Eff: mode ← done
    trigger(net_mcasti(i,OUT))

net_rcvi(j,OUT)
Eff: if mode ≠ done then
    let r = min {r' ∈ N+ | value[r',j] = ⊥}
    out[r] ← out[r] ∪ {j}

leavei
Eff: mode ← done

leave_detecti(j)
Eff: if mode ≠ done then
    let r = min {r' ∈ N+ | value[r',j] = ⊥}
    out[r] ← out[r] ∪ {j}

fail_detecti(j)
Eff: if mode ≠ done then
    let r = min {r' ∈ N+ | value[r',j] = ⊥}
    failed[r] ← failed[r] ∪ {j}

```

Figure 5: CUP_i transitions.

The code works as follows. When an $\text{init}_i(v, W)$ input occurs, process i triggers a $\text{net_mcast}(i, 1, v, W)$ to send its initial value v and estimated world W to all processes, including itself.

For each round $r \geq 2$, process i performs an explicit $\text{net_mcast}_i(i, r, v, W)$ to multicast its round r value v and world W . The world W is determined to be the set of processes that i thinks are still active, that is, the processes in i 's previous world that i does not know to be out or to have failed in round r . Process i may perform this multicast only if its round is $r-1$, it has received round $r-1$ messages from all the processes in W , and it is not currently able to decide. The value v that is sent is the minimum value that i has recorded for round $r-1$ from a process in W .

When a $\text{net_rcv}_i(j, r, v, W)$ occurs, process i puts v and W into the appropriate places in the `value` and `world` arrays.

When an abstain_i input occurs, process i sends an `OUT` message, so that other processes will know not to wait for further messages from it, and stops participating in the algorithm.

Process i can decide at a round r when it has received messages from all processes in its `world[r, i]` except those that are out at round r , such that all of these messages contain the same value and contain worlds that are subsets of `world[r, i]`. The subset requirement ensures that processes in `world[r, i]` will not consider values from processes outside of `world[r, i]` in determining their values for future rounds. When process i decides, it multicasts an `OUT` message and stops participating in the algorithm.

When a $\text{net_rcv}_i(j, \text{OUT})$ occurs, process i records that j is out of the algorithm starting from the first round for which i has not yet received a regular message from j .

When leave_i occurs, process i just stops participating in the algorithm. When $\text{leave_detect}_i(j)$ occurs, process i records that j is out; when this occurs, the lossless leave assumption ensures that i has already received all the messages j sent. The round that is recorded for the leave is the first round after the round of the last message received from j .

Process i knows that another process has failed if it learns about the failure via a `fail_detect` event.

In the next section, we prove the algorithm's correctness. In Section 5.3, we show that the algorithm is early-deciding in the sense that the number of rounds it executes is proportional to the number of actual failures that occur, and does not depend on the number of participants or on the number of processes that leave.

5.3 The Early-Deciding Property

We now show that the algorithm is early-deciding in the sense that the number of rounds it executes is proportional to the number of actual failures that occur, and does not depend on the number of participants or on the number of processes that leave.

We start with some more lemmas.

Lemma 5.1 *If $\text{init}_i(*, W)$ occurs prior to init_j , then $j \in W$.*

Proof: The environment well-formedness assumption implies that j does not leave or fail before it initiates, and hence does not leave or fail before i initiates. Therefore, by world consistency, $j \in W$. ■

Invariant 5.1 *If $(i, 1, *, W)$ and $(j, 2, *, *)$ are in the Net then $j \in W$.*

Proof: By strong induction. For the inductive step, assume that, in the final state of the execution, $(i, 1, *, W)$ and $(j, 2, *, *)$ are in the Net. Then both init_i and init_j events appear in

the execution. If init_i precedes init_j , then Lemma 5.1 implies the result, so assume that init_j precedes init_i .

Since a round 2 message from j is in the Net, a round 1 message $(j, 1, *, W')$ is also. Then Lemma 5.1 implies that $i \in W'$.

We claim that j does not leave or fail before the init_i . Suppose for the sake of contradiction that it does. Then the $\text{net_mcast}(j, 2, *, *)$ event precedes the init_i . Then environment well-formedness implies that i does not fail or leave prior to the $\text{net_mcast}(j, 2, *, *)$ event, because it initiates after this event. Also, i does not abstain, because it initiates. And i does not decide prior to the $\text{net_mcast}(j, 2, *, *)$ event, because that precedes the init_i . Therefore, $i \notin \text{failed}[1]_j \cup \text{out}[1]_j$ in the pre-state of the $\text{net_mcast}(j, 2, *, *)$ event, so $i \in \text{world}[1, j]_j \setminus \text{failed}[1]_j \cup \text{out}[1]_j$ in that state. The precondition of net_mcast implies that $\text{value}[1, i]_j \neq \perp$ in the pre-state, that is, j has received a round 1 message from i before the $\text{net_mcast}(j, 2, *, *)$. But this cannot happen, because init_i happens after the $\text{net_mcast}(j, 2, *, *)$. This contradiction implies that j does not leave or fail before the init_i . Then world consistency implies that $j \in W$, as needed. \blacksquare

In the rest of this section, we consider a situation where no failures happen from some point onward in an execution, and where the rounds of all processes are at most r at the point where failures cease. The following lemma says that all round $r+2$ messages that are ever sent have the same world component.

Lemma 5.2 *Suppose that $r > 0$. Suppose that there is a point t in an execution such that every process has round $\leq r$ at point t , and no fail events happen from t onward.*

If $\text{net_mcast}(i, r+2, v, W)$ and $\text{net_mcast}(j, r+2, v', W')$ both occur in the execution, then $W = W'$.

Proof: We show that $W \subseteq W'$. The other direction is analogous.

The two sets are determined in the precondition of net_mcast , as follows:

$W = \text{world}[r+1, i]_i \setminus \text{out}[r+1]_i \setminus \text{failed}[r+1]_i$, where the values of the last two terms are taken from the pre-state of $\text{net_mcast}(i, r+2, v, W)$, and

$W' = \text{world}[r+1, j]_j \setminus \text{out}[r+1]_j \setminus \text{failed}[r+1]_j$, where the values of the last two terms are taken from the pre-state of $\text{net_mcast}(j, r+2, v', W')$. Invariant A.7 implies that $W = \text{world}[1, i]_i \setminus \text{out-by}[r+1]_i \setminus \text{failed-by}[r+1]_i$, where the values of the last two terms are taken from the pre-state of $\text{net_mcast}(i, r+2, v, W)$, and

$W' = \text{world}[1, j]_j \setminus \text{out-by}[r+1]_j \setminus \text{failed-by}[r+1]_j$, where the values of the last two terms are taken from the pre-state of $\text{net_mcast}(j, r+2, v', W')$.

Consider some $k \in W$. The precondition of $\text{net_mcast}(i, r+2, v, W)$ implies that in the pre-state, $\text{value}[r+1, k]_i \neq \perp$, that is, i has received a round $r+1$ message from k . This means that k has previously sent a round $r+1$ message. Since (by assumption) $r > 0$, it follows that $r+1 \geq 2$, which means that k has sent a round 2 message. Invariant 5.1, applied to any state after both $\text{net_mcast}(j, 1, *, *)$ and $\text{net_mcast}(k, 2, *, *)$ have occurred, implies that k is in the world component of j 's round 1 message, and so k is put into $\text{world}[1, j]_j$ when that is defined. To prove that $k \in W'$, it suffices to show that k is never placed into either of the sets $\text{out-by}[r+1]_j$ or $\text{failed-by}[r+1]_j$.

First, we show that k is never placed into $\text{out-by}[r+1]_j$. Suppose for the sake of contradiction that k is put into $\text{out-by}[r+1]_j$ at some point in the execution. Then consider some state that occurs after this has happened, and that is not before the pre-state of $\text{net_mcast}(i, r+2, v, W)$. In

this state, we have both $\text{value}[r+1,k]_i \neq \perp$ and $k \in \text{out-by}[r+1]_j$. This contradicts Invariant A.4. Therefore, k is never placed into $\text{out-by}[r+1]_j$.

Second, we show that k is never placed into $\text{failed-by}[r+1]_j$. Suppose for the sake of contradiction that k is put into $\text{failed-by}[r+1]_j$ at some point in the execution. Then k fails in the execution, which implies that it fails before point t . But we have already noted that k sends a round $r+1$ message during the execution. It does not send this before point t , because that would mean that it would reach round $r+1$ before point t , contradiction our assumptions. So k sends the round $r+1$ message after point t , and so it cannot fail before point t , a contradiction. Therefore, k is never placed into $\text{failed-by}[r+1]_j$. ■

The next lemma says that, under the same assumptions as for the previous lemma, all the round $r+2$ messages have the same value component.

Lemma 5.3 *Suppose that $r > 0$. Suppose that there is a point t in an execution such that every process has $\text{round} \leq r$ at point t , and no fail events happen from t onward.*

If $\text{net_mcast}(i, r+2, v, W)$ and $\text{net_mcast}(j, r+2, v', W')$ both occur in the execution, then $v = v'$.

Proof: Process i determines v as the minimum of all values $\text{value}[r+1,k]_i$ for all $k \in W$, and process j determines v' as the minimum of all values $\text{value}[r+1,k]_j$ for all $k \in W'$. Lemma 5.2 implies that $W = W'$. Since values are consistent (by Invariant A.2), the sets of values over which the two minima are taken are identical. Therefore, $v = v'$. ■

Finally, we prove the main early-deciding theorem. It says that, if no failures happen from some point onward and the rounds of all processes are at most r when failures cease, then no CUP participant ever advances beyond round $r+2$. Since we have already proved termination, this implies that all active CUP participants decide by round $r+2$.

Theorem 5.4 *Suppose that $r > 0$. Suppose that there is a point t in the execution such that every process has $\text{round} \leq r$ at point t , and no fail events happen from t onward.*

Then every process always has $\text{round} \leq r+2$.

Proof: Lemmas 5.2 and 5.3 yield a common value and world for round $r+2$ messages. Fix v' and W' to be the common value and world, respectively.

We show that the precondition of $\text{net_mcast}(i, r+3, *, *)$ can never be true, which implies that such an event can never happen. This implies that every process always has $\text{round} \leq r+2$. Suppose for the sake of contradiction that the precondition of $\text{net_mcast}(i, r+3, v, W)$ is true in some reachable state s , for some fixed i .

Since the precondition holds in s , $\text{world}[r+2, i]_i \neq \perp$ in s , and so Invariant A.1 implies that some $(i, r+2, v'', W'')$ message is in the Net in s , where $v'' = \text{value}[r+2, i]_i$ and $W'' = \text{world}[r+2, i]_i$. Since v' and W' are the common value and world for round $r+2$ messages, this implies that $\text{value}[r+2, i]_i = v'$ and $\text{world}[r+2, i]_i = W'$.

We show that for all $j \in \text{world}[r+2, i]_i \setminus \text{out}[r+2]_i$, $\text{value}[r+2, j]_i = \text{value}[r+2, i]_i$ and $\text{world}[r+2]_i \subseteq \text{world}[r+2, i]_i$. This suffices to show that the final precondition fails, which yields a contradiction.

Fix $j \in \text{world}[r+2, i]_i \setminus \text{out}[r+2]_i$. Since $\text{failed}[r+2]_i = \{\}$, it follows that $j \in \text{world}[r+2, i]_i \setminus \text{out}[r+2]_i \setminus \text{failed}[r+2]_i$. The precondition of the net_mcast then implies that $\text{value}[r+2, j]_i \neq \perp$ in state s . Invariant A.1 then implies that some $(i, r+2, v''', W''')$ message is in the Net in s ,

where $v''' = \text{value}[r+2, j]_i$ and $W''' = \text{world}[r+2, i]_i$. Since v' and W' are the only value and world for round 2 messages, this implies that $\text{value}[r+2, j]_i = v'$ and $\text{world}[r+2, j]_i = W'$ in state s . Thus, $\text{value}[r+2, j]_i = \text{value}[r+2, i]_i$ and $\text{world}[r+2, j]_i \subseteq \text{world}[r+2, i]_i$, as needed. ■

Note that this proof does not work for the case where $r=0$, because of potential differences in the initial worlds of correct processes. Consider, for example, an execution in which no process ever fails, and some process, k , leaves after sending a round 1 message. Process k may be included in the initial world of process i but not in the initial world of another process j , if j initiates CUP after k leaves. In this case, i takes k 's round 1 message into account when choosing its round 2 message, while j does not (because k is not in j 's initial world). This scenario can only occur in round 1, because no process can send a round 2 message before j initiates.

For the case where $r = 0$, the best we can state is:

Corollary 5.5 *Suppose there is a point t in the execution such that every process has `round` = 0 at point t , and no fail events happen from t onward.*

Then every process always has `round` ≤ 3 .

Proof: This is immediate from Theorem 5.4, using $r = 1$. ■

5.4 Timing Assumptions

For the sake of analyzing the performance of the CUP algorithm, we use timed I/O automata [21, Ch. 23]. We can regard an ordinary I/O automaton as a special case of the timed model, in which arbitrary amounts of time can pass between events. All the safety results carry over to this model.

For this analysis, we add an extra assumption: we assume that any action that is enabled either gets performed or gets disabled by another action, before any time passes.

5.5 Latency Analysis

We now analyze the algorithm's latency in executions in which there are time bounds on certain environment actions. We assume the following bounds:

1. δ_1 is an upper bound on *message latency*. That is, if a `net_rcv(m)` event occurs, the time since the corresponding `net_mcast(m)` is at most δ_1 .
2. δ_2 is an upper bound on *failure and leave detection time*. Moreover, if a message is lost due to failure, then the failure is detected at most δ_2 after the lost message was sent. More precisely,
 - (a) Assume `initi(*, W)` occurs with $j \in W$ and `failj` or `leavej` occurs at time t . Then `fail_detecti(j)`, `leave_detecti(j)`, `decidei`, `leavei`, or `faili` occurs by time $t + \delta_2$.
 - (b) Define $U = \cup_{k \in I} \{W \mid \text{init}_k(*, W) \text{ occurs}\}$. Assume $i, j \in U$ and `net_mcastj(m)` occurs at time t but no `net_rcvi(m)` occurs. Then `fail_detecti(j)`, `leave_detecti(j)`, `decidei`, `leavei`, or `faili` occurs by time $t + \delta_2$.
3. δ_3 is an upper bound on the time difference between the initiation time of different processes. More precisely:
Assume some process initiates at time t and does not fail by time $t + \delta_1$. Assume further that `initi(*, W)` occurs. Then, every process $j \in W$ initiates, abstains, leaves, or fails by time $t + \delta_3$.

In practice, the failure detection time would be at least as large as the message latency. We therefore assume that $\delta_2 \geq \delta_1$.

We now use the above bounds on the environment to establish bounds on CUP's running times. The next lemma bounds the time it takes from when some process initiates CUP until all processes terminate round 1.

Lemma 5.6 *Assume that some process initiates CUP at time t and does not fail by time $t + \delta_1$. Then by time $t + \delta_2 + \delta_3$, every process that initiates either terminates round 1, or leaves, or fails.*

Proof: Let i be a process that initiates and does not leave or fail by time $t + \delta_2 + \delta_3$. We now show that i terminates round 1 by time $t + \delta_2 + \delta_3$. If i decides by time $t + \delta_2 + \delta_3$, then we are done. We therefore assume that i does not decide by this time.

In order to terminate round 1, i has to have a round 1 message from every process $j \in \text{world}[1, i]_i \setminus \text{out}[1]_i \setminus \text{failed}[1]_i$. That is, for every process $j \in \text{world}[1, i]_i$, i has to receive a round 1 message or an OUT message from j , or a `fail_detecti(j)` or a `leave_detecti(j)` event.

Fix a process $j \in \text{world}[1, i]_i$, i.e., j is in i 's initial world. Since some process initiates at time t , by our assumption on initiation times, j initiates, abstains, leaves, or fails by time $t + \delta_3$.

If j fails or leaves by time $t + \delta_3$, then by our assumption on failure and leave detection times, `fail_detecti(j)` or `leave_detecti(j)` occurs by time $t + \delta_2 + \delta_3$ (since we assume that i does not decide, leave, or fail by this time), and we are done.

Assume now that j does not fail or leave by time $t + \delta_3$. Since j is in i 's initial world, j either initiates or abstains by this time, at which point j sends a round 1 message or an OUT message (resp.). If i receives this message, i receives it by time $t + \delta_3 + \delta_1$. If i does not receive this message, `fail_detect(j)i` or `leave_detect(j)i` occurs by time $t + \delta_3 + \delta_2$.

Since $\delta_2 \geq \delta_1$, we get that for every $j \in \text{world}[1, i]_i$, by time $t + \delta_3 + \delta_2$, i either receives a round 1 message or an OUT message from j or a `fail_detect(j)i` or `leave_detect(j)i` event occurs. ■

The following lemma bounds the duration of subsequent rounds.

Lemma 5.7 *Assume that by time t , every process that initiates CUP either terminates round $r > 0$, or decides, or leaves, or fails. Then, by time $t + \delta_2$, every process that initiates CUP either terminates round $r+1$, or decides, or leaves, or fails.*

Proof: Consider a process i that initiates CUP and does not leave or fail or decide by time $t + \delta_2$. We now show that i terminates round $r+1$ by time $t + \delta_2$.

In order to terminate round $r+1$, i has to have a round $r+1$ message from every process $j \in \text{world}[r+1, i]_i \setminus \text{out}[r+1]_i \setminus \text{failed}[r+1]_i$. That is, for every process $j \in \text{world}[r+1, i]_i$, i has to either receive a round $r+1$ message or an OUT message from j , or a `fail_detecti(j)` or a `leave_detecti(j)` event has to occur.

Fix a process $j \in \text{world}[r+1, i]_i$. Process j must have initiated. By time t , j terminates round $r+1$, or decides, or leaves, or fails. If j leaves or fails by time t , then `fail_detect(j)i` or `leave_detect(j)i` occurs by time $t + \delta_2$. Otherwise, j sends a round $r+1$ message or an OUT message (in case it decides) by time t . If i receives this message, i receives it by time $t + \delta_1$. Otherwise, `fail_detect(j)i` or `leave_detect(j)i` occurs by time $t + \delta_2$. Since $\delta_2 \geq \delta_1$, we get that i terminates round $r+1$ by time $t + \delta_2$. ■

Using the two lemmas above, we get the following bound on the running time of an execution of CUP with r rounds.

Lemma 5.8 *Assume that some process initiates CUP at time t and does not fail by time $t + \delta_1$. If i decides at round $r > 0$, it does so by time $t + \delta_3 + r\delta_2$.*

Proof: By Lemma 5.6, by time $t + \delta_3 + \delta_2$, every process that initiates CUP either terminates round 1, or leaves, or fails. By iterative application of Lemma 5.7, we get that by time $t + \delta_3 + \delta_2 + (r - 1)\delta_2 = t + \delta_3 + r\delta_2$, every process that initiates CUP either terminates round r , or decides, or leaves, or fails. ■

As a consequence of the above lemmas and the early-deciding theorem of the previous section we get the following theorem:

Theorem 5.9 *Suppose that there is a point t in the execution such that no fail events happen from t onward. Suppose also that some process initiates CUP by time t . Then every process that decides, decides by time $t + \delta_3 + 3\delta_2$.*

Proof: Let r be the highest value of `round` of any process at time t . Since some process initiated CUP by time t , $r > 0$. By Theorem 5.4, every process that decides, decides at the end of round $r+2$ at the latest.

We consider two cases. First, if $r > 1$, then by Invariant A.12, every process that initiated CUP has either terminated round $r-1$ or left or failed by time t . By applying Lemma 5.7 three times, we get that every process that initiates CUP either terminates round $r+2$ or leaves or fails by time $t + 3\delta_2$. Therefore, in this case, every process that decides, decides by time $t + 3\delta_2$.

Next, assume that $r = 1$. Since some process initiates CUP by time t and does not fail, by Lemma 5.6, by time $t + \delta_3 + \delta_2$, every process that initiates CUP either terminates round 1, or leaves, or fails. By applying Lemma 5.7 twice, we get that every process that initiates CUP either terminates round $r+2$ or leaves or fails by time $t + \delta_3 + 3\delta_2$. Therefore, in this case, every process that decides, decides by time $t + \delta_3 + 3\delta_2$. ■

6 Environment and Model Assumptions for Atom

6.1 Timing Assumptions

We model time using a continuous global variable `now`, which holds the real time. This is a real variable, initially 0. We assume that it increases with derivative 1. Each endpoint i is equipped with a local clock, `clocki`, modeled by a continuous, bijective, monotonically increasing function from the nonnegative \mathbb{R} to the nonnegative \mathbb{R} .

We assume a bound of Γ on clock skew, where Γ is a positive real number. Specifically, for each endpoint i , we assume that in any state of the system that is reachable $|clock_i - now| \leq \Gamma/2$. That is, the difference between each local clock and the real time is at most $\Gamma/2$. It follows that the clock skew between any pair of processes is Γ , formally: in any reachable state, and for any two endpoints i and j , $|clock_i - clock_j| \leq \Gamma$.

We assume that local processing time is 0 and that actions are scheduled immediately when they are enabled. Formally, when any locally controlled action of any process that is part of our local algorithm is enabled, then before any time passes, the action is either performed or becomes disabled.

6.2 Reliable Network Assumptions

We assume that we are given a low-level reliable network service Net . Like DAB, Net is parameterized by a message alphabet, M .

The $\text{Net}(M)$ signature is defined in Figure 6. The actions are the same as those of DAB, except that they are prefixed with `net_`.

Input:

`net_joini`, `net_leavei`, `faili`, $i \in I$,
`net_mcasti(m)`, $m \in M$, $i \in I$

Output:

`net_join_OKi`, `net_leave_OKi`, $i \in I$,
`net_rcvi(m)`, $m \in M$, $i \in I$

Figure 6: The signature of the Net service.

$\text{Net}(M)$ assumes that its application satisfies the same basic safety conditions as those specified above for $\text{DAB}(M)$, except that action names are preceded with `net_`. Assuming the application satisfies these conditions, $\text{Net}(M)$ satisfies a number of safety and liveness properties.

First, Net satisfies the basic properties specified above for DAB: join/leave integrity, message integrity, eventual join, and eventual leave. All of these properties are the same as for DAB, except that action names are prefixed with `net_`.

In addition, Net guarantees FIFO delivery of messages:

- *FIFO delivery*: If `net_mcasti(m)` occurs before `net_mcasti(m')`, and `net_rcvj(m')` occurs, then `net_rcvj(m)` occurs before `net_rcvj(m')`.

$\text{Net}(M)$ also satisfies the following liveness property:

- *Eventual delivery*: Suppose `net_mcasti(m)` occurs after `net_join_OKj`, and no `faili` occurs. Then either `net_leavej` or `failj` or `net_rcvj(m)` occurs.

Additionally, the network latency is bounded by a constant nonnegative real number Δ . Formally, $\text{Net}(M)$ guarantees:

- *Message latency*: If `net_rcvj(m)` occurs, then the real time elapsed since the corresponding `net_mcasti(m)` is at most Δ .

The maximum message latency of Δ guaranteed by Net is intended to include any pre-send delay at the network module of the sending process.

Since an implementation of Net cannot predict the future, it must deliver messages within time Δ as long as no failures occur. In particular, if a message is sent more than Δ time before its sender fails, it must be delivered.

7 The Atom Algorithm

The Atom algorithm consists of a collection of processes corresponding to the different endpoints in I . It uses Net and CUP services as building blocks. It uses multiple instances of CUP.

7.1 Data Types

Atom defines the constant Θ , a positive real number. This will represent a time slot. We assume that $\Theta > \Delta$.

Recall that M represents the message alphabet of DAB. We will use M' to represent the message alphabet of Net. We define the message alphabet of Net in term of the alphabet of Atom:

- M_1 , the set of finite sequences of elements of M . These are the bulk messages processes send.
- $M_2 = M_1 \cup \{JOIN, LEAVE\} \cup \{CUP - INIT\} \times I$
- $M' = I \times M_2 \times \mathbb{N}$.

M' is the complete message alphabet of Net. Each message contains either a bulk message (sequence of client messages) for a particular slot, a request to join or leave a particular slot, or a report that process has initiated consensus on behalf of a particular endpoint. Each message is tagged with the sender and the slot.

7.2 Using the Net and CUP

The Net service alphabet is instantiated with M' . That is, Atom uses a service $Net(M')$ to implement the service $DAB(M)$. Atom uses multiple instances of CUP, at most one for each process j .

As before, a $fail_i$ action causes process i to stop. $fail_i$ actions go to all the components, i.e., Net and all instances of CUP (including dormant ones), and cause all of them to stop taking any locally controlled actions. Since $fail_i$ actions cannot be intercepted, we do not include them in the code.

$leave_i$ actions also go directly to all the local instances of CUP, including dormant ones.

7.3 Atom Algorithm Overview

The algorithm divides time, and respectively, messages, into slots. As time advances, each process advances through slot. The duration of a slot is Θ .

Each process multicasts all of its messages for a given slot in one bulk message. This is a useful abstraction that we make in order to simplify the presentation and analysis of the Atom algorithm. In practice, the bulk message does not have to be sent as one message; a standard packet assembly/disassembly layer can be used to provide all-or-nothing behavior.

Message delivery is also done in order of slots. Before delivering messages of a certain slot s , each process has to determine the *membership* of this slot, that is, the set of processes from which to deliver messages in this slot. To ensure total order, all the processes that deliver messages for a certain slot have to agree upon the membership of each slot. For each slot, messages are delivered in the order of process indices, and for each process, the messages are unpacked from its bulk message and delivered in FIFO order.

7.4 Signature

The signature of Atom at process i , $Atom_i$, is presented in Figure 7. It includes all the interaction with the client and all the interaction with the underlying network. The implementation of Atom uses CUP as a building block. Hence $Atom_i$ has additional input and output actions for interacting with CUP. Since Atom uses multiple instances of CUP, at most one for each process j , actions of

CUP automata are prefixed with $\text{CUP}(j)$. For example, process i uses the action $\text{CUP}(j).\text{init}_i$ to initiate the CUP automaton associated with process j . CUP.fail and CUP.leave are *not* output actions of Atom_i , since they are routed directly from the environment to all instances of CUP.

The signature of Atom_i also includes two internal actions, end_slot , and members . These two actions play a role in determining the membership for each slot. $\text{end_slot}(s)_i$ occurs at a time by which slot s messages from all processes should have reached process i . At this point, processes from which messages are expected but do not arrive are *suspected* to have failed. For each suspected process j , $\text{CUP}(j)$ is run to have the surviving processes agree upon j 's failure slot. This is needed because failed processes can be suspected at different slots by different surviving processes. After CUP reaches decisions about all the suspected processes that could have failed at slot s , $\text{members}(P, s)$ can occur, with P being the agreed membership for slot s . When process i performs $\text{members}(P, s)_i$, all the messages included in bulk messages that i received for slot s from processes in P are delivered (their delivery is triggered) in order of process indices.

Input:

```

joini, leavei
net_join_OKi, net_leave_OKi
mcasti(m), m ∈ M
net_rcvi(m), m ∈ M'
faili
CUP(j).decidei(v), v ∈ N

```

Output:

```

join_OKi, leave_OKi
net_joini, net_leavei
net_mcasti(i, m, s), m ∈ M2, s ∈ N
rcvi(m), m ∈ M
CUP(j).initi(v, W), v ∈ N
CUP(j).abstaini
CUP(j).leave_detecti(j), j ∈ I
CUP(j).fail_detecti(j), j ∈ I

```

Internal:

```

end_sloti(s), s ∈ N
membersi(P, s), P set of I, s ∈ N

```

Figure 7: Atom_i : Signature.

7.5 Pseudo-code

The Atom_i code is presented in Figures 8–10. The state components are presented in Figure 8.

Recall that we do not assume that processes execute the algorithm from the beginning of time. Rather, the application issues an explicit `join` event, and waits for a `join_OK`. The variable `join-slot` holds the slot at which a process starts participating in the algorithm; this will be the value of `current-slot` when `join_OK` will be issued, and the first slot for which a bulk message will be sent. If a process explicitly leaves the algorithm, its `leave-slot` holds the slot immediately following the last slot in which the process sends a bulk message. Both `join-slot` and `leave-slot` are initially ∞ , so as to be larger than any actual slot number they are compared with.

```

clock ∈ ℝ, initially ∈ [0, Γ/2]; dynamic type: continuous functions

join-slot ∈ ℕ ∪ ∞, initially ∞
leave-slot ∈ ℕ ∪ ∞, initially ∞
did-join-OK, boolean, initially false
did-leave, boolean, initially false

mcast-slots ⊆ ℕ, initially {}
ended-slots ⊆ ℕ, initially {}
reported-slots ⊆ ℕ, initially {}

for every s ∈ ℕ
  out-buf[s] ∈ M2, initially empty sequence of M
  joiners[s] ⊆ I, initially {}
  leavers[s] ⊆ I, initially {}
  suspects[s] ⊆ I, initially {}

for every s ∈ ℕ, j ∈ I
  in-buf[j,s], j ∈ I, s ∈ ℕ, finite sequence of M or ⊥, initially ⊥

for every j ∈ I \ { i }
  CUP-status[j] ∈ { idle, req, running, done }, initially idle
  CUP-req-val[j] ∈ ℕ ∪ { ⊥ }, initially ⊥
  CUP-dec-val[j] ∈ ℕ ∪ { ⊥ }, initially ⊥

derived variables:

current-slot ∈ ℕ = ⌊ clock / Θ ⌋

for every s ∈ ℕ
  alive[s] ⊆ I = { j | in-buf[s,j] ≠ ⊥ }

```

Figure 8: Atom_i: State.

The boolean flags `did-join-OK` and `did-leave` are used to ensure that `join_OK` and `net_leave` actions will not be performed more than once. The set `mcast-slots` keeps track of the slots for which the process already multicast a message (JOIN, LEAVE, or bulk). Likewise, `ended-slots` and `reported-slots` keep track of the slots for which the process already performed the `end_slot` or `members` actions, resp.

`out-buf[s]` stores the message (bulk, JOIN, or LEAVE) that is multicast for slot `s`; it initially holds an empty sequence, and in an active slot, all application messages are appended into it. A JOIN message is inserted for the slot before the `join-slot`, and a LEAVE message for the `leave-slot`. Either way, there is no overlap with a bulk message.

The variables `joiners[s]` and `leavers[s]` keep track of the processes `j` for which `join-slotj = s` (resp. `leave-slotj = s`). `suspects[s]` is the set of processes suspected in slot `s` as determined when `end_slot(s)` occurs.

The variable `in-buf[j,s]` is a finite sequence of messages received in a slot `s` bulk message from process `j`. The data type finite sequence supports assignment, extraction of the head of the queue, and testing for emptiness.

There are three variables for tracking the status and values of the different instances of CUP. `CUP-status[j]` is initially `idle`; when CUP(`j`) is initiated, it becomes `running`; if a CUP-INIT message for `j` arrives, it becomes `req`; and when there is a decision for CUP(`j`), or if the process abstains from CUP(`j`), it becomes `done`. `CUP-req-val[j]` holds the lowest slot value associated

with a CUP-INIT message for j (\perp if no such message has arrived). Finally, `CUP-dec-val[j]` holds the decision reached by `CUP(j)`, and \perp if there is none.

`alive[s]` is a derived variable, storing the set of processes from which slot s bulk messages were received.

```

joini
Eff: trigger(net_joini)

net_join_OKi
Eff: join-slot ← current-slot + 2 + [ Γ/Θ ]
    out-buf[join-slot - 1] ← JOIN

join_OKi
Pre: did-join-OK = false
    current-slot = join-slot
Eff: did-join-OK ← true

leavei
Eff: if (join-slot ∈ N) then
    leave-slot ← max(current-slot, join-slot) + 1
    out-buf[leave-slot] ← LEAVE

net_leavei
Pre: did-leave = false
    leave-slot ∈ mcast-slots
Eff: did-leave ← true

net_leave_OK
Eff: trigger(leave_OKi)

mcasti(m)
Eff: if (join-slot ≤ current-slot < leave-slot) then
    append m to out-buf[current-slot]

net_mcasti(i, m, s)
Pre: join-slot ∈ N
    join-slot - 1 ≤ s ≤ leave-slot
    current-slot = s+1
    s ∉ mcast-slots
    m = out-buf[s]
Eff: mcast-slots ← mcast-slots ∪ { s }

net_rcvi(j, JOIN, s)
Eff: joiners[s+1] ← joiners[s+1] ∪ { j }

net_rcvi(j, LEAVE, s)
Eff: leavers[s] ← leavers[s] ∪ { j }
    foreach (k such that CUP-status[k] = running) do
        trigger(CUP(k).leave_detecti(j))

net_rcvi(j, m, s), m sequence of M
Eff: in-buf[j,s] ← m

```

Figure 9: Atom_i : Transitions related to multicast, join, and leave.

In Figure 9 we present the first part of Atom 's transitions, including transitions related to joining, leaving, multicasting messages, and receiving messages from the network. Transitions related to membership and totally ordered delivery are presented in Figure 10.

When the application issues a `join`, Atom triggers `net_join`. Once the Net responds with a `net_join_OK`, Atom calculates the `join-slot` to be $2 + \lceil \Gamma/\Theta \rceil$ slots in the future. This will allow enough time for the `join` message to reach the other processes. A JOIN message is then inserted into `out-buf[join-slot-1]`. Once the `current-slot` reaches `join-slot`, `join_OK` is issued to the application.

When the application issues a `leave`, the `leave-slot` is chosen to be the ensuing slot, and a LEAVE message is inserted into `out-buf[leave-slot]`. A `net_leave` is issued after the LEAVE message has been multicast, and the `net_leave_OK` triggers a `leave_OK` to the application.

Messages multicast by the application are appended to the bulk message for the current slot in `out-buf[current-slot]`. Once a slot `s` ends, the message pertaining to this slot is multicast to the other processes using `net_mcast`. If `s = join-slot - 1`, a JOIN message is sent. If `s = leave-slot`, a LEAVE message is sent, and if `s` is between `join-slot` and `leave-slot - 1`, a bulk message is sent.

When a bulk message is received, it is stored in the appropriate `in-buf`. When a JOIN (LEAVE) message is received, the sender is added to the `joiners` (resp. `leavers`) set for the appropriate slot. Additionally, when a LEAVE message is received, `CUP.leave_detect` is triggered for all running instances of CUP.

Process i performs `end_sloti(s)` once it should have received all the slot `s` messages sent by other non-failed processes. Since slot `s` messages are sent immediately when slot `s` ends, messages are delayed at most Δ time in Net, and the clock difference is at most Γ , process i should have all the non-failed processes' slot `s` messages $\Delta + \Gamma$ time after slot `s+1` began. At this time, `clock` $> (s + 1)\Theta + \Delta + \Gamma$. Process i expects to receive slot `s` bulk messages from all the processes that are in `alive[s-1]`, except for those that are leaving in slot `s`. Any process from which a slot `s` bulk message is expected but does not arrive becomes suspected at this point, and is included in `suspects[s]`.

For every suspected process, CUP is run in order to agree upon the slot at which the process failed. The slot `s` in which the process is suspected is used as the initial value for CUP. The estimated world for CUP is `alive[s] \cup joiners[s+1]`. This way, if k joins in slot `s+1`, k is included in the estimated world. This is needed in order to satisfy the world consistency assumption of CUP, because k can detect the same failure at slot `s+1`, and therefore participate in `CUP(j)`. When i initiates `CUP(j)`, it also multicasts a (CUP-INIT, j) message. If a process k does not detect the failure and does not participate, the (CUP-INIT, j) message forces k to abstain.

Since Atom implements the failure detector for CUP, the effect of `end_sloti(s)` also triggers `CUP(k).fail_detect(j)` actions for every suspected process j , and for every currently running instance k of CUP.

Process i abstains from `CUP(j)` only if a (CUP-INIT, j) message has previously arrived, setting `CUP-status[j]i = req`, and only if `end_sloti` has already occurred for a slot value greater than `CUP-req-val[j]i`. The latter condition ensures that i abstains only from instances of CUP that it will not initiate. This is because the network guarantees that when a process fails, at most one slot bulk message from this process is lost (since we assume that $\Delta \leq \Theta$). This implies that the detection of j 's failure by two non-failed processes can occur at most one slot apart. Therefore, if `end_sloti` has already occurred for a slot value greater than `CUP-req-val[j]i`, i will never suspect j .

The `members(P, s)` action triggers the delivery of all slot `s` messages from processes in P . It can only occur once agreement has been reached about the processes to be included in P . Since the slot at which a process k is suspected by two processes i and j can differ by at most one, `membersi(P, s)` can occur after i receives decision from all instances of CUP pertaining to processes suspected in


```

end_sloti(s)
Pre: join-slot ≤ s
    leave-slot = ∞
    s ∉ ended-slots
    clock > (s+1)Θ + Δ + Γ
Eff: ended-slots ← ended-slots ∪ { s }
    suspects[s] ← (alive[s-1] ∪ joiners[s] \ leavers[s]) \ alive[s]
    foreach (j ∈ suspects[s]) do
        trigger(CUP(j).initi(s, alive[s] ∪ joiners[s+1]))
        net_mcasti(i, (CUP-INIT, j), s)
        CUP-status[j] ← running
        foreach (k such that CUP-status[k] = running) do
            trigger(CUP(k).fail_detecti(j))

net_rcvi(j, (CUP-INIT, k), s)
Eff: if (CUP-status[k] = idle ∨ CUP-req-val[k] > s) then
    CUP-status[k] ← req
    CUP-req-val[k] ← s

CUP(j).abstaini
Pre: CUP-status[j] = req
    ∃s ∈ ended-slots : s > CUP-req-val[j]
Eff: CUP-status[j] ← done

CUP(j).decidei(s)
Eff: CUP-status[j] ← done
    CUP-dec-val[j] ← s

membersi(P, s)
Pre: s = min{ ended-slots \ reported-slots }
    s + 1 ∈ ended-slots
    ∀j ∈ (suspects[s] ∪ suspects[s+1]) : CUP-status[j] = done
    P = { j ∈ alive[s] | CUP-dec-val[j] = ⊥ ∨ CUP-dec-val[j] > s }
Eff: reported-slots ← reported-slots ∪ { s }
    foreach j ∈ P, in order of indices do
        while in-buf[j,s] not empty do
            trigger(rcvi(head(in-buf[i,s])))

```

Figure 10: Atom_i: Transitions related to membership and message delivery.

slots up to $s+1$. Therefore, $\text{members}_i(P, s)$ must occur after $\text{end_slot}(s+1)$, when the suspicions for slot $s+1$ are determined. The set P includes every process j that is alive in slot s and for which there is either no CUP instance running (in which case j was not suspected), or the CUP decision value is greater than s .

7.6 Latency Analysis

In this section we analyze the latency guarantees of Atom. In Section 7.6.1 we show that in failure free executions, Atom's message latency is bounded by $\Delta + 2\Theta + 2\Gamma$. We denote this bound by Δ_{Atom} . In Section 7.6.2, we assign values to the constants that were used in the analysis of CUP in Section 5.5 (δ_1 , δ_2 , and δ_3). Then, in Section 7.6.3, we consider executions in which failures do occur but there is a long time period with no failures. We analyze the time it takes Atom to clear the backlog it has due to past failures, and reach a situation in which message latency is bounded

by the same bound as in failure free executions, namely Δ_{Atom} , barring additional failures.

The fact that once failures stop for a bounded time all messages are delivered within constant time implies that in periods with f failure, Atom's latency is at most linear in the number of failing processes.

7.6.1 Failure free executions

Lemma 7.1 *The time from when process j starts slot \mathbf{s} (i.e., `current-slotj` becomes \mathbf{s}) until process i performs `end_sloti($\mathbf{s}+1$)` is at most $\Delta + 2\Theta + 2\Gamma$.*

Proof: According to its preconditions, `end_sloti(\mathbf{s})` occurs $2\Theta + \Delta + \Gamma$ time after i starts slot \mathbf{s} . Since the difference between two processes' clocks is at most Γ , i starts slot \mathbf{s} at most Γ time after j starts this slot. ■

Lemma 7.2 *Consider an execution in which no process fails. If the application at process j performs `mcastj(\mathbf{m})` when `current-sloti = \mathbf{s}` and if process i delivers m , then i delivers m immediately after `end_sloti($\mathbf{s}+1$)` occurs.*

Proof: If i delivers m , `rcvi(\mathbf{m})` is triggered during the `Membersi(\mathbf{P}, \mathbf{s})` action. Since no process fails, `suspects[\mathbf{s}]i ∪ suspects[$\mathbf{s}+1$]i` is an empty set, and thus the only precondition that needs to be satisfied in order to perform `Membersi(\mathbf{P}, \mathbf{s})` is `$\mathbf{s}+1 \in \text{ended-slots}_i$` , which is true immediately after `end_sloti($\mathbf{s}+1$)` occurs. ■

As a direct result of these two lemmas, we get the following theorem:

Theorem 7.3 *If the application at process j performs `mcastj(\mathbf{m})` at time t , and if process i delivers m , then i delivers m by time $t + \Delta_{Atom} = t + \Delta + 2\Theta + 2\Gamma$.*

7.6.2 CUP bounds

We now assign values to the constants used in the analysis of CUP in Section 5.5. Recall, δ_1 is an upper bound on message latency; δ_2 is an upper bound on failure and leave detection time, and if a message is lost due to failure, then the failure is detected at most δ_2 after the lost message was sent; and δ_3 is an upper bound on the difference between different processes' initiation times.

Lemma 7.4 $\delta_1 = \Delta$

Proof: By definition, both Δ and δ_1 are defined to be upper bounds on the underlying network latency. ■

Lemma 7.5 $\delta_2 = \Delta + 3\Theta + 2\Gamma$

Proof: Assume that `CUP(\mathbf{k}).initi(\ast, W)` occurs with $j \in W$. Assume that one of the following happens at time t : `failj`, `leavej`, or `net_mcastj(\mathbf{m})` for a message \mathbf{m} that is lost because j subsequently fails. Let \mathbf{s} be the value of `current-slotj` at time t . Assume also that by time $t + \Delta + 3\Theta + 2\Gamma$, i does not decide, leave, or fail, so `CUP-status[\mathbf{k}]i = running` and i is active at this time. We have to show that by this time, `fail_detecti(j)` or `leave_detecti(j)` occurs.

If j fails at time t , then j 's slot \mathbf{s} message is never sent, and therefore i detects the failure and invokes `CUP(\mathbf{k}).fail_detecti(j)` during `end_sloti(\mathbf{s})` at the latest. By Lemma 7.1, this occurs

by time $t + 2\Theta + \Delta + 2\Gamma$. Likewise, if j sends a message m while $\text{current-slot}_j = \mathbf{s}$, and m is lost, then by the FIFO nature of the network, j 's slot \mathbf{s} message is also lost and i detects j 's failure during $\text{end_slot}_i(\mathbf{s})$ at the latest.

Assume next that j leaves when $\text{current-slot}_j = \mathbf{s}$, i.e., j 's leave-slot is $\mathbf{s}+1$. If i receives a LEAVE message from j , it receives it before $\text{end_slot}_i(\mathbf{s}+1)$ occurs, and immediately triggers $\text{CUP}(\mathbf{k}).\text{leave_detect}_i(j)$. Otherwise, i receives no slot $\mathbf{s}+1$ message from j and suspects j and invokes $\text{CUP}(\mathbf{k}).\text{fail_detect}_i(j)$ during $\text{end_slot}_i(\mathbf{s}+1)$. This occurs by time $t + 3\Theta + \Delta + 2\Gamma$. ■

Lemma 7.6 $\delta_3 = \Gamma + \Theta$

Proof: Assume that some process process l initiates $\text{CUP}(\mathbf{k})$ at time t and does not fail by time $t + \Delta$. Assume further that $\text{CUP}(\mathbf{k}).\text{init}_i(*, W)$ occurs with $j \in W$. We have to show that j initiates, abstains, leaves, or fails by time $t + \Gamma + \Theta$.

Process l triggers $\text{CUP}(\mathbf{k}).\text{init}_l(\mathbf{s}, *)$ during the $\text{end_slot}_l(\mathbf{s})$ action, and $k \in \text{suspects}_l[\mathbf{s}]$. If j initiates $\text{CUP}(\mathbf{k})$, there is a slot \mathbf{s}' such that j triggers $\text{CUP}(\mathbf{k}).\text{init}_j$ during the $\text{end_slot}_j(\mathbf{s}')$ action, and $k \in \text{suspects}_j[\mathbf{s}']$. By Invariant A.19, $\mathbf{s}' \leq \mathbf{s}+1$. Therefore, $\text{CUP}(\mathbf{k}).\text{init}_j$ occurs no later than time $t + \Gamma + \Theta$, and we are done.

Assume now that j does not initiate $\text{CUP}(\mathbf{k})$, and does not leave or fail by time $t + \Gamma + \Theta$. We now show that j abstains from $\text{CUP}(\mathbf{k})$ by time $t + \Gamma + \Theta$.

When $\text{CUP}(\mathbf{k}).\text{init}_l(\mathbf{s}, *)$ is triggered, l multicasts a $(\text{CUP-INIT}, \mathbf{k})$ message. By Lemma A.8, net_join_OK_j occurs before l initiates $\text{CUP}(\mathbf{k})$, that is, before l multicasts this message. Moreover, by assumption, l does not fail by time $t + \Delta$ and j does not leave or fail by time $t + \Delta$ (because $\Delta \leq \Theta$). Therefore, j receives this message by time $t + \Delta$, which is before time $t + \Gamma + \Theta$. After j receives this message, $\text{CUP-status}[\mathbf{k}]_j$ is req and $\text{CUP-req-val}[\mathbf{k}]_j$ is less than or equal to \mathbf{s} . By time $t + \Gamma + \Theta$, $\text{end_slot}_j(\mathbf{s}+1)$ occurs and the condition for $\text{CUP}(\mathbf{k}).\text{abstain}_j$ becomes true, and remains true until $\text{CUP}(\mathbf{k}).\text{abstain}_j$ occurs and changes $\text{CUP-status}[\mathbf{k}]_j$. Therefore, before any time passes, $\text{CUP}(\mathbf{k}).\text{abstain}_j$ occurs. ■

7.6.3 Failure free periods

We now consider executions in which failures do occur but there are long time periods with no failures. We analyze the time it takes Atom to clear the backlog it has due to past failures, and again reach a situation in which message latency is bounded by Δ_{Atom} , barring additional failures.

Let $t_1 = \delta_3 + 4\delta_2$, where δ_3 and δ_2 are bounds as given above for the difference between process initiation times and failure detection time, resp. From Lemmas 7.6 and 7.5 we get that $t_1 = \Gamma + \Theta + 4(\Delta + 3\Theta + 2\Gamma) = 4\Delta + 9\Gamma + 13\Theta$.

Assume that from time t to time $t' = t + t_1$ there are no failures. We now show that if a message m is sent after time t' , and there are no failures for a period of length Δ_{Atom} after m is sent, then m is delivered within Δ_{Atom} time of when it is sent. Since the delivery order preserves the FIFO order, this also implies that any message m' sent before time t' is delivered by time t' barring failures in the Δ_{Atom} time interval after m' is sent.

Theorem 7.7 *Assume no process fails between time t and $t' = t + t_1$. If $\text{mcast}(\mathbf{m})_j$ occurs at a time t'' such that $t + t_1 \leq t''$, and no failures occur from time t'' to time $t'' + \Delta_{\text{Atom}}$, and if i delivers m , then i delivers m by time $t'' + \Delta_{\text{Atom}}$.*

Proof: By Lemma 7.5, by time $t + \delta_2$ all the processes detect all the failures that occur by time t . Therefore, no process initiates an instance of CUP after time $t + \delta_2$. Since no failures occur after time $t + \delta_2$, by Theorem 5.9, all CUP instances that i initiates terminate by time $t + \delta_2 + \delta_3 + 3\delta_2 = t + t_1$.

Let \mathbf{s} be the value of `current-slotj` at time t'' (i.e., when `mcast(m)j` occurs). By Lemma 7.1, process i performs `end-sloti(s+1)` by time $t'' + \Delta + 2\Theta + 2\Gamma = t'' + \Delta_{Atom}$. At this time, there are no active CUP instances, because CUP instances pertaining to failures that occurred before time t have all terminated and no new failures occur until time $t'' + \Delta_{Atom}$. Therefore, for every slot $\mathbf{s}' \leq \mathbf{s}$, in order of slot numbers, `Members(P, s')` _{i} becomes enabled until it occurs. So `Members(P, s)` _{i} occurs before any time passes. If i delivers m , `rcvi(m)` is triggered during the `Membersi(P, s)` action, so `rcvi(m)` also occurs before any time passes. ■

7.7 Extending Atom to Cope with Late Messages

In this paper, we assumed a synchronous model with deterministic network latency guarantees. Since the network latency, Δ is expected to be of a smaller order of magnitude than Θ , it would not significantly hurt time bounds if conservative assumptions are made in the choice of Δ .

In ongoing research we are considering networks where latency bounds are more likely to be violated. For example, some networks may support differentiated services with probabilistic latency guarantees. Moreover, loss rates may exceed the bounds assumed in the implementation of the reliable network. Such networks can be represented using the timed-asynchronous [11] failure model.

Although our algorithm cannot guarantee atomic broadcast semantics while network latency and reliability guarantees are violated, it is important for the algorithm to be able to recover from such situations, and to once more provide correct semantics after network guarantees are re-established. In addition, it would be desirable to inform the application when a violation of Atom semantics occurs, and when the correct semantics are resumed (following the failure awareness approach of [14]).

There are some strategies that can be used to make Atom recover from periods in which network guarantees are violated. For example, a lost or late message can cause inaccurate failure suspicions. With Atom, if a process k is falsely suspected, it will receive a (CUP-INIT, k) message for itself. In order to recover from such a situation, we could have the process “commit suicide” in such a situation, that is inform the application of the failure and have the application re-join as a new process. The full modification of Atom for this setting is ongoing work.

8 Conclusions

We have defined two new problems, *Dynamic Atomic Broadcast* and *Consensus with Unknown Participants*. We have presented new algorithms for both problems. The latency of both of our algorithms depends linearly on the number of failures that occur during a particular execution, but does not depend on an upper bound on the potential number of failures, nor on the numbers of joins and leaves that happen during the execution.

Acknowledgments

We thank Alan Fekete and Rachid Guerraoui for comments that helped improve the presentation.

References

- [1] T. Abdelzaher, A. Shaikh, F. Jahanian, and K. Shin. RTCAST: Lightweight multicast for real-time process groups. In *IEEE Real-Time Technology and Applications Symposium (RTAS)*, June 1996.
- [2] Z. Bar-Joseph, I. Keidar, T. Anker, and N. Lynch. QoS preserving totally ordered multicast. Technical Report MIT-LCS-TR-796, MIT Laboratory for Computer Science, January 2000. Url: <http://theory.lcs.mit.edu/~idish/Abstracts/qos.html>.
- [3] Z. Bar-Joseph, I. Keidar, T. Anker, and N. Lynch. QoS preserving totally ordered multicast. In F. Butelle, editor, *5th International Conference On Principles Of Distributed Systems (OPODIS)*, pages 143–162, December 2000. Special issue of *Studia Informatica Universalis*.
- [4] K. Birman and R. van Renesse. *Reliable Distributed Computing with the Isis Toolkit*. IEEE Computer Society Press, 1994.
- [5] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *J. ACM*, 43(2):225–267, Mar. 1996.
- [6] G. Chockler and D. Malkhi. Active Disk Paxos with infinitely many processes. In *21st ACM Symposium on Principles of Distributed Computing (PODC)*, July 2002. To appear.
- [7] G. V. Chockler, I. Keidar, and R. Vitenberg. Group Communication Specifications: A Comprehensive Study. *ACM Comput. Surv.*, 33(4):1–43, December 2001. Previous version: MIT Technical Report MIT-LCS-TR-790, September 1999.
- [8] F. Cristian. Synchronous atomic broadcast for redundant broadcast channels. *The Journal of Real-Time Systems*, 2:195–212, 1990.
- [9] F. Cristian. Reaching agreement on processor group membership in synchronous distributed systems. *Distributed Computing*, 4(4):175–187, April 1991.
- [10] F. Cristian, H. Aghili, R. Strong, and D. Dolev. Atomic broadcast: From simple message diffusion to Byzantine agreement. *Inform. Comput.*, 118:158–179, April 1995. Early version in FICS15, June, 1985.
- [11] F. Cristian and C. Fetzer. The timed asynchronous distributed system model. *IEEE Transactions on Parallel and Distributed Systems*, pages 642–657, June 1999.
- [12] D. Dolev and D. Malkhi. The Transis approach to high availability cluster communication. *Commun. ACM*, 39(4):64–70, April 1996.
- [13] D. Dolev, R. Reischuk, and H. R. Strong. Early stopping in byzantine agreement. *J. ACM*, 37(4):720–741, October 1990.
- [14] C. Fetzer and F. Cristian. Fail-awareness: An approach to construct fail-safe applications. In *27th Annual International Fault-Tolerant Computing Symposium*, pages 282–291, 1997.
- [15] B. Glade, K. Birman, R. Cooper, and R. van Renesse. Lightweight process groups in the Isis system. *Distributed Systems Engineering*, 1:29–36, 1993.

- [16] A. Gopal, R. Strong, S. Toueg, and F. Cristian. Early-delivery atomic broadcast. In *9th ACM Symposium on Principles of Distributed Computing (PODC)*, pages 297–309, 1990.
- [17] V. Hadzilacos and S. Toueg. Fault-tolerant broadcasts and related problems. In S. Mullender, editor, *chapter in: Distributed Systems*. ACM Press, 1993.
- [18] S. Katz, P. Lincoln, and J. Rushby. Low-overhead time-triggered group membership. In *11th International Workshop on Distributed Algorithms (WDAG)*, pages 155–169, 1997. LNCS 1320.
- [19] H. Kopetz and G. Grunsteidl. Ttp - a protocol for fault-tolerant real-time systems. *IEEE Computer*, pages 14–23, January 1994.
- [20] L. Lamport. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM*, 21(7):558–565, July 78.
- [21] N. A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, 1996.
- [22] M. Merritt and G. Taubenfeld. Computing with infinitely many processes under assumptions on concurrency and participation. In *14th International Symposium on Distributed Computing (DISC)*, October 2000.
- [23] D. Powell. *Delta-4: A Generic Architecture for Dependable Distributed Computing*. Springer Verlag, 1991.
- [24] L. Rodrigues, K. Guo, A. Sargento, R. van Renesse, B. Glade, P. Verissimo, and K. Birman. A dynamic light-weight group service. In *15th IEEE International Symposium on Reliable Distributed Systems (SRDS)*, pages 23–25, Oct. 1996. also Cornell University Technical Report, TR96-1611, August, 1996.
- [25] L. Rodrigues and P. Verissimo. *xAMP*, A multi-primitive group communications service. In *11th IEEE International Symposium on Reliable Distributed Systems (SRDS)*, pages 112–121, Oct. 1992.
- [26] R. D. Schlichting and F. B. Schneider. Fail stop processors: an approach to designing fault-tolerant computing systems. *ACM Trans. Comput. Syst.*, 1(3):222–238, Aug. 1983.
- [27] F. B. Schneider. Implementing fault tolerant services using the state machine approach: A tutorial. *ACM Comput. Surv.*, 22(4):299–319, December 1990.
- [28] P. Verissimo, J. Rufino, and L. Rodrigues. Enforcing real-time behaviour of lan-based protocols. In *10th IFAC Workshop on Distributed Computer Control Systems*, September 1991.

A Correctness Proofs

A.1 Correctness of the CUP Algorithm

We consider the system consisting of a composition of automata CUP_i , one for each $i \in I$. We consider a restricted set of executions of this composition—those in which the environment safety assumptions are all satisfied. The invariants we state throughout this section should be interpreted as saying that the stated property is true for all states that occur in such executions.

A.1.1 General invariants

We say that a message is *in the Net* if a `net_mcast` event for that message has occurred or is in a `trigger-buffer`.

The first invariant lists an assortment of basic constraints. They can be proved using induction.

Invariant A.1 1. $\text{value}[r,i]_j = \perp$ if and only if $\text{world}[r,i]_j = \perp$.

2. If $\text{value}[r,i]_j = v \neq \perp$ and $\text{world}[r,i]_j = W \neq \perp$ then an (i,r,v,W) message is in the Net.
3. If $(i,r,*,*)$ is in the Net then $\text{round}_i \geq r$.
4. If $\text{mode}_i = \perp$ then $\text{round}_i = 0$.
5. If $\text{mode}_i = \text{running}$ then some $(i,1,*,*)$ message is in the Net.
6. If $i \in \text{failed}[r]_j$ then fail_i has occurred.
7. If $i \in \text{failed}[r]_j$ and $s \geq r$, then $\text{value}[s,i]_j = \perp$.

The next invariant expresses consistency of values and worlds of the same process at different places in the system.

Invariant A.2 1. If messages (j,r,v,W) and (j,r,v',W') are in the Net then $v = v'$ and $W = W'$.

2. If $\text{value}[r,j]_i \neq \perp$ and $\text{value}[r,j]_{i'} \neq \perp$ then $\text{value}[r,j]_i = \text{value}[r,j]_{i'}$.
3. If $\text{world}[r,j]_i \neq \perp$ and $\text{world}[r,j]_{i'} \neq \perp$ then $\text{world}[r,j]_i = \text{world}[r,j]_{i'}$.
4. If $\text{value}[r,j]_i \neq \perp$ and $\text{world}[r,j]_i \neq \perp$ and a message (j,r,v,W) is in the Net then $\text{value}[r,j]_i = v$ and $\text{world}[r,j]_i = W$.

The next two invariants describe some facts that follow from the existence of OUT messages and from the detection of leaves.

Invariant A.3 1. If an (i,OUT) message is in the Net then $\text{mode}_i = \text{done}$.

2. If $i \in \text{out}[r]_j$ then $\text{mode}_i = \text{done}$.

Proof: By induction. Part 2 uses the accurate leave detector assumption. ■

Invariant A.4 If $i \in \text{out}[r]_k$ and $s \geq r$, then no message of the form $(i,s,*,*)$ is in the Net, and for all j , $\text{value}[s,i]_j = \perp$.

Proof: By strong induction. First, we claim that a `net_mcasti` event cannot convert the invariant from true to false by falsifying the conclusion while leaving the hypothesis true. This is because, if the hypothesis is true, then $i \in \text{out}[r]_k$ in the pre-state of the `net_mcasti`, which implies, by Invariant A.3, that `modei = done`. But the precondition of `net_mcasti` requires that `modei = running`, a contradiction.

The key steps are, therefore, those that make the hypothesis true. Index i can be added to `out[r]k` by receipt of an OUT message by k or by a `leave_detectk(i)`. An OUT message may result from a previous `abstaini` that occurs when `modei = ⊥`, or a previous `decidei` event.

For `abstaini`, by Invariant A.1, we know that in the pre-state of the `abstaini`, `roundi = 0`. Then Invariant A.1 implies that in the pre-state, no message of the form $(i, *, *, *)$ is in the Net, and for all j and all s , `value[s, i]j = ⊥`. Once the `abstaini` happens, `modei` becomes `done`, which means that no later messages are sent.

For `decidei`, the FIFO assumption for message delivery implies that the `decidei` event must have occurred when `round = r - 1`. Invariant A.1 then implies that in the pre-state of the `decidei`, the conclusion of the invariant holds. Since the `decidei` event sets `modei` to `done`, i sends no further messages, so the conclusion continues to hold.

For `leave_detectk(i)`, we know by the lossless leave assumption that before the `leave_detectk(i)` occurs, k has already received every message that has ever been `net_mcast` by i . Since k explicitly checks that it has no values from i for round r , there are no such messages in the Net. ■

The following says that any value that appears anywhere in the system is some participant's initial value.

Invariant A.5 1. *If (i, r, v, W) is in the Net then there exists j and W' such that $(j, 1, v, W')$ is in the Net.*

2. *If `value[r, k]i = v ≠ ⊥` then there exists j and W' such that $(j, 1, v, W')$ is in the Net.*

Proof: We show Parts 1 and 2 together by induction on the length of a finite execution.

Base: Trivial, because no messages are initially in the Net and no values are initially non- \perp .

Inductive step: We first show part 1. The interesting steps are those in which a message (i, r, v, W) is put into the Net. If $r = 1$ then (i, r, v, W) is put into the Net by an `initi(v, W)` event, which puts the `net_mcast` into `trigger-bufferi`. But this immediately satisfies the conclusion. On the other hand, if $r \geq 2$, then (i, r, v, W) is put into the Net by an explicit `net_mcast(i, r, v, W)` step. In this case, v is obtained from a set of values already in i 's `value` array in the pre-state. The induction hypothesis, part 2, then implies that some $(j, 1, v, W')$ is already in the Net, as needed.

For part 2, the key step is a `net_rcvi(k, r, v, W)` for some W . In the pre-state of such a step, message (k, r, v, W) is in the Net. The inductive hypothesis, part 1, then implies that some $(j, 1, v, W')$ is already in the Net, as needed. ■

The following invariant asserts that processes are always in their own worlds.

Invariant A.6 1. *If a message (i, r, v, W) is in the Net then $i \in W$.*

2. *If `world[r, i]j ≠ ⊥` then $i \in \text{world}[r, i]_j$.*

Proof: We prove part 1 by induction on the length of the execution, with a trivial base case.

Inductive step: The interesting steps are those in which a message (i, r, v, W) is put into the Net. If $r = 1$, then this is done by an `initi(v, W)` step. In this case, the environment well-formedness

assumption implies that no leave_i or fail_i event precedes the init_i , and so the world consistency assumption implies that $i \in W$, as needed. On the other hand, if $r \geq 2$, then (i, r, v, W) is put into the Net by an explicit $\text{net_mcast}(i, r, v, W)$ step. In this case, the precondition says that $\text{mode}_i = \text{running}$ and $\text{world}[1, i]_i \neq \perp$ in the pre-state. In this pre-state, i is not in any $\text{failed}[r]_i$ set, because if it were, Invariant A.1 would imply that i has failed, and it would not be able to perform the net_mcast . Also, in this pre-state, i is not in any $\text{out}[r]_i$ set, by Invariant A.3. Therefore, i is included in W , because of the way W is defined.

Part 2 follows from part 1 and Invariant A.1. ■

The following invariant describes consequences of the definition of a round $r+1$ world and value:

Invariant A.7 *If $(i, r+1, v, W)$ is in the Net, for $r \geq 1$, then:*

1. For every $j \in W$, $\text{world}[r, j]_i \neq \perp$.
2. $W = \text{world}[r, i]_i \setminus \text{out}[r]_i \setminus \text{failed}[r]_i$.
3. $v = \min \{ \text{value}[r, j]_i : j \in W \}$.
4. $W = \text{world}[1, i]_i \setminus \text{out-by}[r]_i \setminus \text{failed-by}[r]_i$.

Proof: Part 1 is proved by an easy induction; the key step is $\text{net_mcast}(i, r+1, v, W)$, and the conclusion follows immediately from the precondition.

Given part 1, we prove part 2 by induction. Now the interesting steps are $\text{net_mcast}_i(i, r+1, v, W)$, $\text{net_rcv}_i(j, \text{OUT})$, $\text{leave_detect}_i(j)$, and $\text{fail_detect}_i(j)$. The fact that $\text{net_mcast}_i(i, r+1, v, W)$ yields the property follows immediately from the precondition. A $\text{net_rcv}_i(j, \text{OUT})$ event could only falsify the property if $j \in W$ and the event puts j into $\text{out}[r]_i$. However, part 1 implies that $\text{world}[r, j]_i \neq \perp$ in the post-state, and hence $\text{value}[r, j]_i \neq \perp$ in the post-state. But this would cause the post-state to violate Invariant A.4, a contradiction. A similar argument shows that $\text{leave_detect}_i(j)$ cannot falsify the property. Finally, $\text{fail_detect}_i(j)$ could only falsify the property if $j \in W$ and the event puts j into $\text{failed}[r]_i$. However, part 1 implies that $\text{world}[r, j]_i \neq \perp$ in the post-state, and hence $\text{value}[r, j]_i \neq \perp$ in the post-state. But this would cause the post-state to violate Invariant A.1, a contradiction.

We prove part 3 by induction, using part 1. This time, the interesting steps are $\text{net_mcast}_i(i, r+1, v, W)$ and $\text{net_rcv}_i(j, r, *, *)$. Again, the $\text{net_mcast}_i(i, r+1, v, W)$ step yields the property immediately. A $\text{net_rcv}_i(j, r, v', *)$ could only falsify the property if $j \in W$. But in this case we know that $\text{value}[r, j]_i \neq \perp$ in the pre-state, and then Invariant A.2 implies that $v' = \text{value}[r, j]_i$ in the pre-state. It follows that this step does not change $\text{value}[r, j]_i$, and so does not falsify the property.

Part 4 is proved by induction on r (not induction on the length of the execution), using part 2. The base case, $r = 1$, follows immediately from part 2. For the inductive step, we suppose that the claim is true for some $r \geq 1$ and show it for $r + 1$. That is, we assume that $(i, r+2, v, W)$ is in the Net. Then by part 2, $W = \text{world}[r+1, i]_i \setminus \text{out}[r+1]_i \setminus \text{failed}[r+1]_i$. Now, since $\text{world}[r+1, i]_i \neq \perp$, Invariant A.1 implies that a message of the form $(i, r+1, v', W')$ is in the Net, where $W' = \text{world}[r+1, i]_i$. By inductive hypothesis, part 4, this implies that $W' = \text{world}[1, i]_i \setminus \text{out-by}[r]_i \setminus \text{failed-by}[r]_i$. Therefore, $W = \text{world}[1, i]_i \setminus \text{out-by}[r]_i \setminus \text{failed-by}[r]_i \setminus \text{out}[r+1]_i \setminus \text{failed}[r+1]_i$. This is equal to $\text{world}[1, i]_i \setminus \text{out-by}[r+1]_i \setminus \text{failed-by}[r+1]_i$, as needed. ■

Invariant A.8 *Suppose that $\text{decide}_i(v)$ has happened at round r . Then:*

1. For all $j \in \text{world}[r, i]_i \setminus \text{out}[r]_i$, $\text{value}[r, j]_i = \text{value}[r, i]_i$ and $\text{world}[r, j]_i \subseteq \text{world}[r, i]_i$.
2. For all $j \in \text{world}[r, i]_i$, if (j, r, v', W) is in the Net then $v' = v$ and $W \subseteq \text{world}[r, i]_i$.

Proof: Part 1 follows from an easy induction: `out` can only grow, and `value` and `world` do not change once they are non- \perp . Therefore, the only interesting step is `decidei(v)`, and the result follows directly from the precondition.

For part 2, consider any state s after a `decidei(v)` has happened at round r . Suppose that $j \in \text{world}[r, i]_i$ and (j, r, v', W) is in the Net, in state s . Then Invariant A.4 implies that $j \notin \text{out}[r]_i$. Thus, $j \in \text{world}[r, i]_i \setminus \text{out}[r]_i$. Then part 1 and Invariant A.2 imply the conclusion. ■

The following invariants say that any process' values and worlds decrease as rounds increase.

Invariant A.9 For any $r \geq 1$, if a message $(i, r+1, v, W)$ is in the Net then $\text{value}[r, i]_i \neq \perp$, $v \leq \text{value}[r, i]_i$, and $W \subseteq \text{world}[r, i]_i$.

Proof: By induction. For the inductive step, the interesting case is when the last action of the execution is `net_mcast(i, r+1, v, W)`. Invariant A.6 implies that $i \in W$. Therefore, the precondition for `net_mcast(i, r+1, v, W)` implies that, in the pre-state, $\text{value}[r, i]_i \neq \perp$. Therefore, this is also true in the post-state, as needed.

Next, we show that $v \leq \text{value}[r, i]_i$. The value v is determined in the `net_mcast` event to be the minimum of the set of values of the form $\text{value}[r, j]_i$, for $j \in W$. Since $i \in W$, this minimum includes $\text{value}[r, i]_i$. Therefore, $v \leq \text{value}[r, i]_i$.

Finally, we show that $W \subseteq \text{world}[r, i]_i$. The value W is determined in the `net_mcast` event to be `world[r, i]_i \setminus \text{out-by}[r] \setminus \text{failed-by}[r]`, according to the values of the `out` and `failed` sets in the pre-state. It follows immediately that W is a subset of $\text{world}[r, i]_i$. ■

Invariant A.10 For any $r \geq 1$,

1. If $\text{value}[r+1, i]_i \neq \perp$ then $\text{value}[r, i]_i \neq \perp$, $\text{value}[r+1, i]_i \leq \text{value}[r, i]_i$, and $\text{world}[r+1, i]_i \subseteq \text{world}[r, i]_i$.
2. If $\text{value}[r, i]_i \neq \perp$ and $1 \leq s \leq r$ then $\text{value}[s, i]_i \neq \perp$, $\text{value}[r, i]_i \leq \text{value}[s, i]_i$, and $\text{world}[r, i]_i \subseteq \text{world}[s, i]_i$.

Proof: For part 1, assume that, in some reachable state, $\text{value}[r+1, i]_i \neq \perp$, and hence $\text{world}[r+1, i]_i \neq \perp$. Then Invariant A.1 implies that in the same state, a message $(i, r+1, v, W)$ must be in the Net, where $v = \text{value}[r+1, i]_i$ and $W = \text{world}[r+1, i]_i$. Invariant A.9 then yields the conclusions. Part 2 follows from part 1, using induction on $r-s$. ■

The following invariant says that, if all the messages for a particular round r are “consistent”, then so are all the messages for all later rounds.

Invariant A.11 Let W be a nonempty finite set, $v \in V$, and $r \geq 1$. Suppose that, for every $i \in W$, if a message of the form (i, r, v', W') is in the Net, then $W' \subseteq W$ and $v' = v$. Then for every $i, j \in W$ and for every $s \geq r$,

1. If a message of the form (i, s, v', W') is in the Net, then $W' \subseteq W$ and $v = v'$.
2. $\text{value}[s, i]_j$ is either v or \perp .

3. $\text{world}[s, i]_j$ is either a subset of W or \perp .

Proof: We prove part 1 by induction on the length of an execution.

Base: The conclusion of the invariant is vacuously true in the start state.

Inductive step: The interesting steps are those that put some message (i, s, v', W') into the Net, where $i \in W$ and $s \geq r$. We may restrict attention to the case where $s > r$, because if $s = r$ and the step falsifies part 1, it also falsifies the hypothesis of the invariant. Thus, the only interesting steps are of the form $\text{net_mcast}(i, s, v', W')$ where $i \in W$ and $s > r$. So consider such a step, and fix i, s, v' , and W' . Assume that the hypothesis of the invariant is true after (and hence before) the step.

We show that $W' \subseteq W$. After the net_mcast step, the message is in the Net. Invariant A.9 then implies that $W' \subseteq \text{world}[s-1, i]_i$. Invariant A.10 then implies that $\text{world}[r, i]_i \neq \perp$ and $\text{world}[s-1, i]_i \subseteq \text{world}[r, i]_i$. Therefore, $W' \subseteq \text{world}[r, i]_i$. Since $\text{world}[r, i]_i \neq \perp$, an $(i, r, *, W'')$ message is in the Net, where $W'' = \text{world}[r, i]_i$. Then since the hypothesis of the invariant is true, it follows that $\text{world}[r, i]_i \subseteq W$. Putting all the pieces together yields that $W' \subseteq W$.

Next, we argue that $v' = v$. The value v' is determined by the precondition of the net_mcast action, as the minimum of a set of values $\text{value}[s-1, j]_i$, taken over all indices j in W' . Because $W' \subseteq W$, every such index j is in W . Since $\text{value}[s-1, j]_i \neq \perp$, a $(j, s-1, v'', W'')$ message is in the Net in the pre-state of the new net_mcast , where $v'' = \text{value}[s-1, j]_i$. Our assumption that the conclusion of the invariant is true in the pre-state then implies that $\text{value}[s-1, j]_i = v$. Thus, all the values considered in the min are equal to v , which implies that $v = v'$.

This proves part 1. Parts 2 and 3 follow from part 1 and Invariant A.1. ■

Invariant A.12 *If $\text{round}_i = r > 1$ and $\text{mode}_j = \text{running}$ and j is not failed, then $\text{round}_j \geq r - 1$.*

Proof: Since j is not failed, by Invariant A.1(6), $j \notin \text{failed}[s]_i$ for any s , so $j \notin \text{failed-by}[r-1]_i$. By Invariant A.3(2), $j \notin \text{out}[s]_i$ for any s , so $j \notin \text{out-by}[r-1]_i$. Since $\text{mode}_j = \text{running}$, j initiated, and by the world consistency assumption, $j \in \text{world}[1, i]_i$. By Invariant A.7(4), j is in i 's world for round r . Therefore, i must have received a round $r-1$ message from j before moving to round r . ■

A.1.2 CUP safety guarantees

We now prove that the CUP implementation satisfies the CUP safety guarantees, assuming the environment satisfies the safety assumptions.

Theorem A.1 *The CUP algorithm satisfies well-formedness.*

Proof: This is straightforward from the code and the well-formedness assumptions on the environment. For condition 1, assume that decide_i occurs. Then in the preceding state, $\text{mode} = \text{running}$. mode is initially \perp , and the only way it becomes running is via init_i . So there must be a preceding init_i .

For condition 2, assume for the sake of contradiction that two decide_i events occur. Part 1 implies that an init_i precedes the first decide_i . The first decide_i sets mode_i to done . After this point, and before the second decide_i event occurs, mode_i must become running . This can happen only as a result of another init_i event. This means that two init_i events must occur, which contradicts the environment well-formedness assumption. Therefore, no more than one decide_i event occurs. ■

Theorem A.2 *The CUP algorithm satisfies uniform agreement.*

Proof: If at most one `decide` event occurs, the result follows immediately. So assume that there are at least two `decide` events. Consider the first `decide` event, `decidei(v)`. By the precondition, we know that in the pre-state, there exists r such that $\text{world}[r, i]_i \neq \perp$ and $\forall j \in \text{world}[r, i]_i \setminus \text{out}[r]_i, \text{value}[r, j]_i = v$ and $\text{world}[r, j]_i \subseteq \text{world}[r, i]_i$. Since i does not leave, abstain, or decide before the `decidei` event, we know that $i \notin \text{out}[r]_i$ in the pre-state; therefore, $\text{value}[r, i]_i = v$. Also consider any particular later `decide` event, `decidei'(v')`. As above, we know that in the pre-state of this event, there exists r' such that $\text{world}[r', i']_{i'} \neq \perp$ and $\forall j \in \text{world}[r', i']_{i'} \setminus \text{out}[r']_{i'}, \text{value}[r', j]_{i'} = v'$ and $\text{world}[r', j]_{i'} \subseteq \text{world}[r', i']_{i'}$. Moreover, $\text{value}[r', i']_{i'} = v'$.

We now show that $i' \in \text{world}[r, i]_i$. Since i' decides, it initiates and does not leave or fail before it decides. Since i initiates before it decides, and thus before i' decides, i' does not leave or fail before i initiates. Then the world consistency assumption implies that i' gets put into $\text{world}[1, i]_i$. If $r = 1$ then we are done, so assume that $r \geq 2$. Then the value of $\text{world}[r, i]_i$ is determined in a `net_mcasti(i, r, *, *)` step. To see that i' is included in $\text{world}[r, i]_i$, note that that set is defined in the `net_mcasti(i, r, *, *)` step to include (at least) all processes in $\text{world}[1, i]_i$ that do not leave, abstain, decide, or fail before the `net_mcasti(i, r, *, *)` event. And i' does not leave, abstain, decide, or fail by then, because this `net_mcasti(i, r, *, *)` event happens prior to the `decidei`.

We also know that $i' \notin \text{out}[r]_i$ in the pre-state of `decidei`. This is because i' has not left, abstained, or decided before the `decidei`.

Next, we show that $r' \geq r$, that is, the round at which i' decides is at least as great as the round at which i decides. Since $i' \in \text{world}[r, i]_i$ and $i' \notin \text{out}[r]_i$ in the pre-state of `decidei`, the precondition for `decidei` implies that $\text{value}[r, i']_i \neq \perp$ in the pre-state of `decidei`. This means that i' must send an `(i', r, v, *)` message. This implies that the round r' at which i' decides is at least as great as r , that is, $r' \geq r$.

Finally, we argue that $v' = v$. Invariant A.8, part 2, implies that in the pre-state of `decidei'`, if $j \in \text{world}[r, i]_i$ and if (j, r, v'', W'') is in the Net, then $v'' = \text{value}[r, i]_i$ and $W'' \subseteq \text{world}[r, i]_i$. Since $r' \geq r$ and $i' \in \text{world}[r, i]_i$, Invariant A.11, part 2, implies that in the pre-state of `decidei'`, $\text{value}[r', i']_{i'}$ is either v or \perp . Since (as noted earlier) $\text{value}[r', i']_{i'} = v'$, we have that $v = v'$. ■

Theorem A.3 *The CUP algorithm satisfies validity.*

Proof: Part 1 follows from Invariant A.5. Part 2 follows from Invariant A.10. ■

A.1.3 CUP liveness guarantees

We now show that CUP satisfies its liveness property—termination. Formally, the lemmas and theorem we state in this section should be interpreted with respect to an execution α of the composition of automata CUP_i for $i \in I$ such that:

1. All the environment safety and liveness assumptions are satisfied in α .
2. α is “weakly fair” to all actions of all CUP_i automata, in the sense that if an action is enabled from some point onward, it eventually is performed.

Lemma A.4 *Let J be the set of processes that initiate and never decide, leave, or fail, and suppose that $i \in J$. If $\text{init}_i(v, W)$ occurs and $j \in W$ then either $j \in J$ or else j abstains, leaves, decides, or fails.*

Proof: Follows from the init occurrence assumption. ■

Lemma A.5 *If process i initiates and never decides, leaves, or fails, then round_i increases without bound.*

Proof: Let J be the set of all processes that initiate and never decide, leave, or fail. Assume for the sake of contradiction that, for some process $i \in J$, round_i is bounded. Let r be the smallest round number such that for some process $i \in J$, round_i is bounded by r , and fix such $i \in J$. Process i cannot get stuck at round 0, because the init_i action immediately increments the round to 1. So we may assume that $r > 0$.

We argue that i cannot be stuck at round r , by showing that for some v, W , the $\text{net_mcast}_i(i, r+1, v, W)$ action is eventually enabled and stays enabled. Then weak fairness implies that $\text{net_mcast}_i(i, r+1, v, W)$ eventually occurs.

We claim that the last precondition of $\text{net_mcast}_i(i, r+1, *, *)$ (the negation of the decide precondition) is always true. For if not, then $\text{decide}_i(v)$ would be enabled for some v , and would stay enabled forever. This implies, by weak fairness, that decide_i occurs, a contradiction.

Next, we claim that for every $j \in \text{world}[1, i]$, either i receives a round r message from j , or else i puts j into its $\text{failed}[r']$ set or $\text{out}[r']$ set for some $r' \leq r$. Fix any such j . Lemma A.4 implies that either $j \in J$ or j eventually abstains, leaves, decides, or fails. If $j \in J$ then by choice of r , j does not get stuck at any round less than r , and so j eventually sends a round r message, which i eventually receives.

If j fails, then eventually a $\text{fail_detect}_i(j)$ occurs, which makes i put j into one of its $\text{failed}[r']$ sets. If $r' \leq r$ then we are done; on the other hand, if $r' > r$ then i receives a round r message from j .

If j abstains and does not fail, then eventually i puts j into its $\text{out}[1]$ set (which suffices because $1 \leq r$). If j leaves or decides at a round $r' \leq r$, then eventually i puts j into its $\text{out}[r']$ set. Finally, if j leaves or decides at a round $r' > r$, then eventually i receives a round r message from j .

This claim implies that eventually the precondition of $\text{net_mcast}_i(i, r+1, v, W)$ is satisfied for some v, W . Because the values and worlds can only decrease, eventually the precondition is satisfied, and remains satisfied, for the same v, W . Then weak fairness implies that the action eventually occurs, which moves j to round $r + 1$. This is a contradiction. ■

Lemma A.6 *Let J be the set of processes that initiate and never decide, leave, or fail, and suppose that $i \in J$. Then for r sufficiently large, $\text{world}[r, i]_i = J$.*

Proof: The result follows from two claims: that for all r , $J \subseteq \text{world}[r, i]_i$, and that for sufficiently large r , $\text{world}[r, i]_i$ is a subset of J .

First, we show that for all r , $J \subseteq \text{world}[r, i]_i$. World consistency implies that $J \subseteq \text{world}[1, i]_i$. Since no element of J ever abstains, leaves, fails, or decides, no element of J is ever put into any $\text{failed}[r]_i$ or $\text{out}[r]_i$. Then the definition of $\text{world}[r, i]_i$ (in $\text{net_mcast}(i, r, *, *)$) implies that for all r , $J \subseteq \text{world}[r, i]_i$.

Second, we show that for sufficiently large r , $\text{world}[r, i]_i$ is a subset of J . Let j be any element of $\text{world}[r, i]_i$. Lemma A.4 implies that if $j \notin J$, then j eventually abstains, leaves, decides, or

fails. But in any of these cases, j eventually gets put into some `failed[r]i` or `out[r]i`. This means that j is excluded from `world[r, i]i` for sufficiently large r . ■

Theorem A.7 *The CUP algorithm satisfies termination.*

Proof: We prove that every process that initiates eventually decides, leaves, or fails. Assume for the sake of contradiction that there is at least one initiator that does not decide, leave, or fail. Let J be the set of processes that initiate and never decide, leave, or fail; then J is not empty. Then Lemma A.5 implies that the rounds of all processes in J increase without bound, and Lemma A.6 implies that for sufficiently large r , `world[r, i]i = J` for all $i \in J$. Thus from some round onward, every process in J bases its new value on values heard from exactly the members of J .

Thereafter, each $i \in J$ eventually reaches some minimum value of `value[r, i]i` (by monotonicity and the fact that only finitely many values can be used). Consider a round beyond which all the minima have been attained. If these are all identical, then all processes can decide based on this value and world J , and we are done. On the other hand, if they are not all identical, then let i be a process whose minimum is larger than some other process' minimum. Then i would see a smaller value and reduce its value further, a contradiction. ■

A.2 Atom Correctness Proof: Safety Arguments

A.2.1 General Invariants

The following invariants follow immediately from the code:

Invariant A.13 *If `join-sloti ≠ ∞` then `leave-sloti > join-sloti`.*

Invariant A.14 *Suppose $s \in \text{ended-slots}_i$. Then:*

1. *If $j \in \text{joiners}[s]_i$ then `join-slotj = s`.*
2. *If $j \in \text{leavers}[s]_i$ then `leave-slotj = s`.*

Proof: Process j can be inserted into `joiners[s]i` (`leavers[s]i`) only if i receives a $(j, \text{JOIN}, s-1)$ (resp. (j, LEAVE, s)) message, which can be sent only by j and only if `join-slotj = s` (resp. `leave-slotj = s`). ■

The following invariant asserts that from the join slot onward, slot messages (bulk, join, or leave) are multicast in order.

Invariant A.15 *If `join-sloti - 1 ≤ s' ≤ s` and $s \in \text{mcast-slots}_i$ then $s' \in \text{mcast-slots}_i$.*

Proof: `join-sloti` had to have been set before `current-sloti` becomes $s'+1$ because it is always chosen to be in the future. Therefore, `net_mcasti(i, *, s')` is enabled once `current-sloti` becomes $s'+1$. This is earlier than the time at which `net_mcasti(i, *, s)` can occur, so time could not have passed beyond that point without `net_mcasti(i, *, s-1)` occurring. ■

The following invariant is central to the rest of the proof. It asserts that by the time of `end_sloti(s)`, i has all the right processes in `alive[s-1]`, `alive[s]`, `joiners[s]`, and `joiners[s+1]`.

Invariant A.16 *If $s \in \text{ended-slots}_i$ then*

1. *If $\text{join-slot}_j \leq s$ and $s \in \text{mcast-slots}_j$ then $j \in \text{alive}[s-1]_i \cup \text{joiners}[s]_i$.*
2. *If $\text{join-slot}_j \leq s+1$ and $s+1 \in \text{mcast-slots}_j$ then $j \in \text{alive}[s]_i \cup \text{joiners}[s+1]$.*

Proof: If j joined by slot s , it registered for the network before starting slot $s-1$. Moreover, if $s \in \text{mcast-slots}_j$, then by Invariant A.15, $s-1 \in \text{mcast-slots}_j$, and therefore j multicasts either a $(j, \text{JOIN}, s-1)$ or a bulk message in slot $s-1$, and by the Net's *reliable delivery* property, this message is not lost due to j 's failure because j multicasts a message in the following slot, which occurs Θ time later, and we assume that $\Theta > \Delta$, and messages sent more than Δ time before the failure are not lost. Likewise, if $\text{join-slot}_j \leq s+1$ and $s+1 \in \text{mcast-slots}_j$, then $s \in \text{mcast-slots}_j$ (by Invariant A.15), and j multicasts a bulk or join message in slot s , which is not lost due to j 's failure.

We will now show two things: first, that i joined early enough to get j 's slot $s-1$ bulk or join message; and second, that $\text{end_slot}_i(s)$ occurred late enough for i to have received j 's slot s bulk or join message.

Since i does end_slot for s , $\text{join-slot}_i \leq s$. Process i chooses its join-slot following the net_join_OK_i to be $\text{current-slot}_i + 2 + \lceil \Gamma/\Theta \rceil$, so current-slot_i becomes $s-1$ at least Γ time after the net_join_OK_i . Since the maximum clock difference between i and j is Γ , j sends its message (join or bulk) for slot $s-1$ no earlier than the time of the net_join_OK_i , so i joined early enough to get j 's message for slot s .

It is left to show that i gets j 's slot s bulk or join message for slot s before $\text{end_slot}_i(s)$. This follows from the precondition for end_slot_i which asserts that $\text{clock}_i > (s+1)\Theta + \Delta + \Gamma$. That is, that at least $\Delta + \Gamma$ time has elapsed since slot $s+1$ has begun at i . Since the clock difference between i and j is at most Γ , we get that at least Δ time has elapsed since slot $s+1$ has begun at j . Since j sends its slot s message once slot $s+1$ begins at j , and the network latency is bounded by Δ , the message reaches i before $\text{end_slot}_i(s)$. ■

The following invariants are related to the $\text{suspects}[s]$ sets.

Invariant A.17 *If $\text{suspects}[s]_i$ is not empty, then $\text{join-slot}_i \leq s$.*

Proof: $\text{suspects}[s]_i$ gets set only upon $\text{end_slot}_i(s)$, for which this is a precondition. Once join-slot_i is set to a non- ∞ value, it does not change, by the singularity of join and net_join_OK . ■

Invariant A.18 *If $j \in \text{suspects}[s]_i$ then j has failed.*

Proof: Since j gets inserted to $\text{suspects}[s]_i$ during $\text{end_slot}_i(s)$, j is in $(\text{alive}[s-1]_i \cup \text{joiners}[s]_i \setminus \text{leavers}[s]_i) \setminus \text{alive}[s]_i$. In particular, j is in $\text{alive}[s-1]_i \cup \text{joiners}[s]_i$, and therefore $\text{join-slot}_j \leq s$. Moreover, j is not in $\text{alive}[s]_i$, so by the contrapositive of Invariant A.16(2), $s+1 \notin \text{mcast-slots}_j$, which implies that j either fails or leaves before sending a slot $s+1$ message. Since j is also not in $\text{leavers}[s]_i$, j must have failed. ■

Invariant A.19 *If $j \in \text{suspects}[s]_i$ and $j \in \text{suspects}[s']_{i'}$ then $|s' - s| \leq 1$.*

Proof: Without loss of generality, assume $s' \geq s$. Since $j \in \text{suspects}[s]_i$, then $j \in \text{alive}[s-1]_i \cup \text{joiners}[s]_i$, and therefore $\text{join-slot}_j \leq s$. Moreover, j is not in $\text{alive}[s]_i$ when $s \in \text{ended-slots}_i$, so by the contrapositive of Invariant A.16(2), $s+1 \notin \text{mcast-slots}_j$. By Invariant A.15, for any slot $r > s$, $r \notin \text{mcast-slots}_j$, and therefore $j \notin \text{alive}[r]_{i'}$ for any $r > s$. Since i' suspects j in slot s' , j is in $\text{alive}[s'-1]_{i'}$, and therefore $s'-1 \leq s$. ■

The following invariant states that a process does not abstain from CUP instances pertaining to processes that it suspects.

Invariant A.20 *If $k \in \text{suspects}[s]_i$ and $\text{CUP-status}[k]_i = \text{done}$ then $\text{CUP-dec-val}[k]_i \neq \perp$.*

Proof: Assume by contradiction that the invariant is false. Since $\text{CUP-status}[k]_i = \text{done}$ while $\text{CUP-dec-val}[k]_i = \perp$, then i must have performed $\text{CUP}(k).\text{abstain}_i$. By the precondition for $\text{CUP}(k).\text{abstain}_i$, $\text{CUP-status}[k]_i$ was req when abstain_i occurred, which implies that $\text{end-slot}(s)_i$ could not have already occurred, that is, all the slots in ended-slots_i were smaller than s at the time of $\text{CUP}(k).\text{abstain}_i$. By the precondition for $\text{CUP}(k).\text{abstain}_i$, when it occurred, $\text{CUP-req-val}[k]_i$ had some non- \perp value, v , such that $v < s-1$. This, in turn, implies that a $(\text{CUP-INIT}, k)$ message with slot v had previously arrived. That means that such a message was previously sent by some j , which implies that k is added to $\text{suspects}[v]_j$, during $\text{end-slot}_j(v)$, and remains there henceforward. But $k \in \text{suspects}[s]_i$ and $v < s-1$, a contradiction to Invariant A.19. ■

Invariant A.21 *If $k \in \text{alive}[s]_i$, $k \notin \text{alive}[s]_j$, and $s \in \text{ended-slots}_j$ then $k \in \text{suspects}[s]_j$. Moreover, if $s+1 \in \text{ended-slots}_i$ then $k \in \text{suspects}[s+1]_i$.*

Proof: Since $k \in \text{alive}[s]_i$, we know that $\text{join-slot}_k \leq s < \text{leave-slot}_k$ and that $s \in \text{mcast-slots}_k$. Therefore, by Invariant A.16(1), if $s \in \text{ended-slots}_j$ then $k \in \text{alive}[s-1]_j \cup \text{joiners}[s]_j$. Additionally, k is neither in $\text{leavers}[s]_i$ nor in $\text{leavers}[s]_j$, because it does not leave at slot s . Therefore, since $k \notin \text{alive}[s]_j$, in $\text{end-slot}_j(s)$, k gets inserted into $\text{suspects}[s]_j$.

Since $k \notin \text{alive}[s]_j$, by the contrapositive of Invariant A.16(2), we get that $s+1 \notin \text{mcast-slots}_k$. That is, k does not send a bulk or leave message for slot $s+1$. Therefore, $k \notin \text{alive}[s+1]_i \cup \text{leavers}[s+1]_i$ when $\text{end-slot}_i(s+1)$ occurs, and k gets inserted into $\text{suspects}[s+1]_i$ when $s+1$ is inserted to ended-slots_i . ■

Invariant A.22 *If $k \in \text{alive}[s]_i$ and $s+1 \in \text{ended-slots}_i$ and $\text{CUP-dec-val}[k]_j \leq s$ then $k \in \text{suspects}[s+1]_i$.*

Proof: Since $k \in \text{alive}[s]_i$, $\text{join-slot}_k \leq s < \text{leave-slot}_k$. Since $\text{CUP-dec-val}[k]_j \leq s$, then by the validity property of CUP, some process l must have initiated $\text{CUP}(k)$ with an initial value $s' \leq s$. This implies that $k \in \text{suspects}[s']_l$, and therefore $k \notin \text{alive}[s']_l$ and $s' \in \text{ended-slots}_l$. By contrapositive of Invariant A.16, $s'+1 \notin \text{mcast-slots}_k$, and therefore also $s'+1 \notin \text{mcast-slots}_k$. So i does not hear a bulk or leave message from k for slot $s+1$, and $k \in \text{suspects}[s+1]_i$. ■

Lemma A.8 *Assume that for some processes j, k, l $\text{CUP}(k).\text{init}_l(v, W)$ occurs with $j \in W$, and that $\text{CUP}(k).\text{init}_i(v', W')$ also occurs. Then net-join-OK_j has occurred before $\text{CUP}(k).\text{init}_i(v', W')$.*

Proof: By the precondition for $\text{CUP}(k).\text{init}_i$, $k \in \text{suspects}[v']_i$ and $k \in \text{suspects}[v]_l$, so by Invariant A.19, $v' \geq v - 1$. When $\text{CUP}(k).\text{init}_i(v, W)$ occurs, $W = \text{alive}[v]_l \cup \text{joiners}[v+1]_l$. Since $j \in W$, this implies that $\text{join-slot}_j \leq v+1 \leq v'+2$. Assume $\text{CUP}(k).\text{init}_i(v', W')$ occurs at time t . So at time t , $v' \in \text{ended-slots}_i$. By the precondition for $\text{end-slot}_i(v')$, at time t $\text{clock}_i > (v'+1)\Theta + \Delta + \Gamma$. Since $v'+1 \geq \text{join-slot}_j - 1$, at this time $\text{clock}_i > (\text{join-slot}_j - 1)\Theta + \Delta + \Gamma$. Since the clock skew is bounded by Γ , at time t , $\text{clock}_j > (\text{join-slot}_j - 1)\Theta + \Delta$. So t is at least Δ time after j begins slot $\text{join-slot}_j - 1$. But join-slot_j is chosen to be at least 2 slots after the slot at which net-join-OK_j occurs at j , so j begins $\text{join-slot}_j - 1$ after the net-join-OK_j , i.e., before time t . ■

A.2.2 Safety environment conditions for CUP

Well-formedness $\text{CUP}(k).\text{init}_i$ only occurs when k becomes suspected at i . Once k is suspected, it is never again alive. Therefore, it is never suspected again and $\text{CUP}(k).\text{init}_i$ occurs at most once. By Invariant A.20, since k is suspected at i , i does not abstain. Thus, at most one init_i or abstain_i event occurs.

The fact that at most one leave_i event occurs and at most one fail_i event occurs is ensured by the application, since leave and fail actions are routed directly from the application to all instances of CUP.

The fact that no fail_i precedes an init_i follows from the fact that failures affect all components and processes do not take any steps after they fail.

One of the preconditions end-slot_i is that $\text{leave-slot}_i = \infty$, that is, that leave_i did not occur. Therefore, no leave_i precedes an init_i .

World consistency Assume that $\text{CUP}(k).\text{init}_i(s, W)$ occurs at time t , j does not leave or fail before time t , and $\text{CUP}(k).\text{init}_j(s', *)$ also occurs. We need to show that $j \in W$.

$\text{CUP}(k).\text{init}_i(s, W)$ is triggered during $\text{end-slot}_i(s)$. We need to show that at this time $j \in \text{alive}[s]_i \cup \text{joiners}[s+1]_i$. This is true if i receives j 's slot s bulk or join message.

By the precondition for $\text{end-slot}_i(s)$, $\text{clock}_i > (s+1)\Theta + \Delta + \Gamma$ at time t . Since the difference between a process clock and real time is at most $\Gamma/2$, the real time associated with point t is at least $(s+1)\Theta + \Delta + \Gamma/2$. By assumption, j does not fail or leave until this time.

By Invariant A.19, $s' \leq s+1$. When $\text{CUP}(k).\text{init}_j(s', *)$ occurs, $k \in \text{suspects}[s']_j$. By Invariant A.17, $\text{join-slot}_j \leq s'$. Together these two inequalities imply that $\text{join-slot}_j \leq s+1$. Therefore, if j does not fail or leave before clock_j becomes $s+1\Theta$, j multicasts its slot s bulk or join message when $\text{clock}_j = (s+1)\Theta$ (a join message is multicast if $\text{join-slot}_j = s+1$; otherwise j multicasts a slot s bulk message). When $\text{clock}_j = (s+1)\Theta$, the real time is at most $(s+1)\Theta + \Gamma/2$. If j does not fail until the real time becomes $(s+1)\Theta + \Gamma/2 + \Delta$, then j 's message is not lost, and i receives it by time $(s+1)\Theta + \Gamma/2 + \Delta$. But we assume that j does not fail or leave until this time.

Accurate failure detector $\text{CUP}(k).\text{fail_detect}_i(j)$ occurs only if for some slot s $j \in \text{suspects}[s]_i$. Therefore, by Invariant A.18, j has previously failed. Moreover, since a process is never again alive after it is suspected, it is never again suspected, and $\text{CUP}(k).\text{fail_detect}_i(j)$ does not recur.

Accurate leave detector $\text{CUP}(k).\text{leave_detect}_i(j)$ occurs only if a LEAVE message is received from j ; j sends at most one such a message and only if it actually leaves.

Lossless leave Assume a CUP process at j multicasts a message m , and subsequently, leave_j occurs. When leave_j occurs, a LEAVE message is inserted to out-buf_j to be sent in the ensuing slot. This LEAVE message is multicast after m . $\text{leave_detect}_i(j)$ occurs when this LEAVE message is received. By the FIFO property of Net, $\text{net_rcv}_j(m)$ occurs beforehand.

A.2.3 Proving the total order property

We now prove that all the process deliver messages in a consistent total order. We define the total order S as follows: Let P_s be the union of all sets P such that an action $\text{members}(P, s)_i$ occurs. The set of messages S_s is defined to be those messages included in slot s bulk messages by processes in P_s . The set of messages in S is defined to be the union of all sets S_s .

The ordering is based on slots, so that for $s < s'$, all messages in S_s precede all messages in $S_{s'}$. For messages pertaining to the same set S_s , the ordering is by process indices. For the same slot and process index, the ordering is the temporal order of sending (at the external boundary of Atom).

We have to show that every process delivers a contiguous subsequence of S . We first prove Lemma A.9, asserting that every two processes that perform a $\text{members}(P, s)$ action for a slot s do so with the same membership set P . As part of this action, processes deliver messages for slot s . Next, we prove Lemma A.10, asserting that if a process i performs $\text{members}_i(P, s)$ with $j \in P$, then i has received a bulk message for slot s from j , and therefore triggers the delivery of all the messages included in it as an effect of the $\text{members}_i(P, s)$ action. Thus, every process that performs $\text{members}(P, s)$, triggers the delivery of all the messages in S_s . These messages are delivered in order of the sender's process index, and for each process, in FIFO order. Therefore, these messages are delivered in the order defined on S_s .

Since every process performs $\text{members}(P, s)$ for a contiguous subsequence of slots, every process delivers a contiguous subsequence of the messages in S .

We now prove the lemmas:

Lemma A.9 *If $\text{members}(P, s)_i$ and $\text{members}(P', s)_j$ occur, then $P = P'$.*

Proof: Let k be a process in P . At the time $\text{members}(P, s)_i$ occurs, $k \in \text{alive}[s]_i$, $s+1 \in \text{ended-slots}_i$, and $\text{CUP-dec-val}[k]_i$ is either \perp or larger than s . Assume by way of contradiction that $k \notin P'$, then either $k \notin \text{alive}[s]_j$ or $\text{CUP-dec-val}[k]_j \leq s$ when $\text{members}(P', s)_j$ occurs.

Assume first that $k \notin \text{alive}[s]_j$. Note that $s \in \text{ended-slots}_j$ when $\text{members}(P', s)_j$ occurs, so by Invariant A.21, $k \in \text{suspects}[s]_j$ at the time $\text{members}(P', s)_j$ occurs. By the precondition for $\text{members}(P', s)_j$, $\text{CUP-status}[k]_j = \text{done}$ when it occurs, and by Invariant A.20, $\text{CUP-dec-val}[k]_j \neq \perp$, that is, $\text{CUP}(k). \text{decide}_j(v)$ occurred for some v and set $\text{CUP-dec-val}[k]_j = v$. By the well-formedness property of CUP, j initiated $\text{CUP}(k)$. Since $k \in \text{suspects}[s]_j$, k cannot be included in $\text{suspects}[s']_j$ for any $s' \neq s$, and so j initiated $\text{CUP}(k)$ with s . By the validity condition of CUP, $v \leq s$.

Since $s+1 \in \text{ended-slots}_i$ when $\text{members}(P, s)_i$ occurs, by Invariant A.21, $k \in \text{suspects}[s+1]_i$ at this time. Therefore, by the precondition for $\text{members}(P, s)_i$, $\text{CUP-status}[k]_i = \text{done}$. By A.20, $\text{CUP-dec-val}[k]_i \neq \perp$, that is, $\text{CUP}(k). \text{decide}_i$ occurred, and by the uniform agreement property, $\text{CUP-dec-val}[k]_i = \text{CUP-dec-val}[k]_j \leq s$. A contradiction.

Now, assume that $\text{CUP-dec-val}[k]_j \leq s$ when $\text{members}(P', s)_j$ occurs. Since $s+1 \in \text{ended-slots}_i$ when $\text{members}(P, s)_i$ occurs, by Invariant A.22, $k \in \text{suspects}[s+1]_i$ at this time. Therefore, by the precondition for $\text{members}(P, s)_i$, $\text{CUP-status}[k]_i = \text{done}$. By Invariant A.20, $\text{CUP-dec-val}[k]_i \neq$

\perp , that is, $\text{CUP}(\mathbf{k}).\text{decide}_i$ occurred, and by the uniform agreement property, $\text{CUP-dec-val}[\mathbf{k}]_i = \text{CUP-dec-val}[\mathbf{k}]_j \leq s$. A contradiction. ■

Lemma A.10 *If $\text{members}_i(\mathbf{P}, \mathbf{s})$ occurs, then for every $j \in \mathbf{P}$, i received a bulk message for slot \mathbf{s} from j prior to the $\text{members}_i(\mathbf{P}, \mathbf{s})$ action.*

Proof: Assume $\text{members}_i(\mathbf{P}, \mathbf{s})$ occurs. Since $j \in \mathbf{P}$, by the precondition for $\text{members}_i(\mathbf{P}, \mathbf{s})$, $j \in \text{alive}[\mathbf{s}]_i$. By definition of $\text{alive}[\mathbf{s}]$, $\text{in-buf}[\mathbf{s}, j] \neq \perp$, that is, i received a bulk message from j for slot \mathbf{s} . ■

A.3 Atom Correctness Proof: Liveness Arguments

In the liveness proof, we can use the safety guarantees of CUP, since they depend only on the safety assumptions about the environment.

A.3.1 General liveness lemmas

Lemma A.11 *Time passes. current-slot_i increases through all slot values from zero onward, as long as i does not fail.*

Lemma A.12 *If i does not leave or fail, then $\text{end-slot}_i(\mathbf{s})$ occurs for every slot $\mathbf{s} \geq \text{join-slot}_i$.*

Lemma A.13 *If i leaves and does not fail, then eventually i multicasts a $(i, \text{LEAVE}, \mathbf{s})$ message.*

A.3.2 Liveness environment conditions for CUP

Init occurrence Assume that $\text{init}_i(\mathbf{s}, \mathbf{W})$ event occurs and $j \in \mathbf{W}$, and neither i nor j leaves or fails.

Since $j \in \mathbf{W}$, $j \in \text{alive}[\mathbf{s}] \cup \text{joiners}[\mathbf{s}+1]$ at the time $\text{init}_i(\mathbf{s}, \mathbf{W})$ is triggered, which means that net_join_OK_j had already occurred prior to the $\text{init}_i(\mathbf{s}, \mathbf{W})$ event. When $\text{init}_i(\mathbf{s}, \mathbf{W})$ is triggered, i multicasts an $(\text{CUP-INIT}, \mathbf{k})$ message. Since neither i nor j leaves or fails, j receives this message.

Consider the pre-state value of $\text{CUP-status}[\mathbf{k}]$ when the $(\text{CUP-INIT}, \mathbf{k})$ message from i arrives at j . If $\text{CUP-status}[\mathbf{k}]$ is `running` or `done`, then either $\text{CUP}(\mathbf{k}).\text{init}_j$ or $\text{CUP}(\mathbf{k}).\text{abstain}_j$ had to have already occurred and we are done. Otherwise, after this step $\text{CUP-status}[\mathbf{k}] = \text{req}$, and $\text{CUP-req-val}[\mathbf{k}] = \mathbf{v}$. Since j does not leave or fail, by Lemma A.12, it eventually has slots larger than \mathbf{v} in ended-slots_j , so either $\text{CUP}(\mathbf{k}).\text{init}_j(\ast)$ or abstain_j becomes enabled, depending on whether \mathbf{k} is suspected in some slot or in none.

Reliable delivery Assume that for some processes j, k, l $\text{CUP}(\mathbf{k}).\text{init}_l(\mathbf{v}, \mathbf{W})$ occurs with $j \in \mathbf{W}$, and that either $\text{CUP}(\mathbf{k}).\text{init}_i(\mathbf{v}', \mathbf{W}')$ or $\text{CUP}(\mathbf{k}).\text{abstain}_i$ occurs. We will show that by the time that either $\text{CUP}(\mathbf{k}).\text{init}_i(\mathbf{v}', \mathbf{W}')$ or $\text{CUP}(\mathbf{k}).\text{abstain}_i$ occurs, net_join_OK_j had already occurred. This will imply that for any $\text{net_mcast}_i(\mathbf{m})$ that occurs after this event, a $\text{net_rcv}_j(\mathbf{m})$ will occur unless either i will fail, or j will fail or leave.

If $\text{CUP}(\mathbf{k}).\text{init}_i(\mathbf{v}', \mathbf{W}')$ occurs, by Lemma A.8, net_join_OK_j occurs first. Now, consider the case that $\text{CUP}(\mathbf{k}).\text{abstain}_i$ occurs. Process i can only abstain after it receives an $(\text{CUP-INIT}, \mathbf{k})$ message which could have only been sent if some other process i' has already triggered $\text{CUP}(\mathbf{k}).\text{init}_{i'}$. By Lemma A.8, net_join_OK_j must have occurred before the $\text{CUP}(\mathbf{k}).\text{init}_{i'}$ event.

Complete leave and failure detector If $\text{CUP}(k).init_i(v, W)$ occurs with $j \in W$, then $j \in \text{alive}[v]_i \cup \text{joiners}[v+1]$. Assume that i does not decide or leave or fail. Then $\text{CUP-status}[k]_i$ remains running from the time of the $\text{CUP}(k).init_i(v, W)$ event onward. If leave_j occurs, j sends a LEAVE message which i receives. When i receives j 's LEAVE message, i triggers $\text{leave_detect}_i(j)$ and we are done. Otherwise, assume j does not leave and fail_j occurs, then eventually there is a slot for which i does not receive j 's messages. Let s be the first such slot, so $j \notin \text{alive}[s]_i$ while $j \in \text{alive}[s-1]_i \cup \text{joiners}[s]_i$, so since j does not leave in s , $j \in \text{suspects}[s]$. Since $j \in \text{alive}[v]_i \cup \text{joiners}[v+1]$, $s > v$, and i triggers $\text{fail_detect}_i(j)$ while performing $\text{end_slot}_i(s)$.

A.3.3 Liveness of Atom

Eventual join Assume no fail_i occurs. When join_i occurs, net_join_i is triggered, and by fairness, eventually occurs. By the eventual join property of Net, net_join_OK_i eventually occurs. At that point, join_slot_i is set to be bigger than current_slot_i . join_slot_i does not change from that point onward, since by the join integrity property of Net, no more net_join_OK_i events occur. By Lemma A.11, current_slot_i eventually becomes equal to join_slot_i . When that happens, join_OK_i becomes enabled, and remains enabled, as long as no time passes, until it occurs. By our assumption on time passage, no time passes until join_OK_i occurs. Therefore, by fairness, it eventually occurs.

Eventual leave Assume no fail_i occurs. When leave_i occurs, leave_slot_i is set to be bigger than current_slot_i . leave_slot_i does not change from that point onward, since by our assumption on the application, no more leave_i events occur. By Lemma A.13, i eventually multicasts a $(i, \text{LEAVE}, \text{leave_slot}_i)$ message, at which point leave_slot_i is added to mcast_slots_i . When that happens, net_leave_i becomes enabled and remains enabled until it occurs. Then, by the eventual leave property of Net, net_leave_OK_i eventually occurs and triggers leave_OK_i .

Message delivery The following lemma asserts that a process that participates in the algorithm and does not leave or fail continues to perform $\text{members}(P, s)$ forever.

Lemma A.14 *If $\text{mcast}_i(m)$ occurs for some m when $s = \text{current_slot}_i$, and no fail_i or leave_i occurs, then for every $s' \geq s$, $\text{members}_i(P, s')$ occurs.*

Proof: Since $\text{mcast}_i(m)$ occurs, by our assumption about the application, it is preceded by a join_OK_i . Therefore, m is appended to $\text{out_buf}[s]_i$. By Lemma A.11, current_slot_i becomes $s+1$, and so i eventually sends its bulk message for slot s with m included in it. By liveness of Net, $\text{net_rcv}_i(i, m', s)$ occurs, where m' is i 's slot s bulk message.

By Lemma A.12, $\text{end_slot}_i(s)$ occurs for every slot $s \geq \text{join_slot}_i$. The sets $\text{suspects}[s]_i$ and $\text{suspects}[s+1]_i$ are set when $\text{end_slot}_i(s)$ (resp. $\text{end_slot}_i(s+1)$) occurs, at which point a CUP instance for each process in these sets is initiated, and the set do not change afterwards. By the termination property of CUP, these instances of CUP eventually terminate, setting the corresponding CUP-status to done. Therefore, $\text{members}(P, s)_i$ eventually becomes enabled for some P , and by fairness, occurs. ■

We now prove that the message delivery liveness property holds.

Assume $\text{mcast}_i(m)$ occurs, and no fail_i or leave_i occurs. We first show that S contains m and $\text{rcv}_i(m)$ occurs. Let $s = \text{current_slot}_i$ when $\text{mcast}_i(m)$ occurs. By Lemma A.14, $\text{members}(P,$

$s)_i$ occurs. We now show that $i \in P$. This will imply that $m \in S$ (by definition of S), and that $\text{rcv}_i(m)$ occurs (since it is triggered by $\text{members}(P, s)_i$).

To show that $i \in P$, we have to show that $j \in \text{alive}[s]_i$ and $\text{CUP-dec-val}[i]_i = \perp$ at the time $\text{members}(P, s)_i$ occurs. Since at this time $s+1 \in \text{ended-slots}_i$, by Invariant A.16, $j \in \text{alive}[s]_i$. By Invariant A.18, since i does not fail it never becomes a suspect, and therefore, no instance of CUP is run for i , and $\text{CUP-dec-val}[i]_i = \perp$.

It remains to show that for every m' that follows m in S , $\text{rcv}_i(m')$ also occurs. By definition of S , m' is included in a bulk message for some slot $s' \geq s$ from some process i' , such that $i' \in P'$ and $\text{members}_j(P', s')$ occurs for some j . By Lemma A.14, $\text{members}(P'', s')_i$ also occurs, and by Lemma A.9, $P' = P''$. Therefore, $\text{rcv}_i(m')$ is triggered by the $\text{members}(P'', s')_i$ action.