

Self-Adjusting Atomic Memory for Dynamic Systems based on Quorums On-The-Fly. Correctness Study

Emmanuelle Anceaume, Maria Gradinariu, Vincent Gramoli, Antonino
Virgillito

► **To cite this version:**

Emmanuelle Anceaume, Maria Gradinariu, Vincent Gramoli, Antonino Virgillito. Self-Adjusting Atomic Memory for Dynamic Systems based on Quorums On-The-Fly. Correctness Study. [Research Report] PI 1795, 2006, pp.15. inria-00001192

HAL Id: inria-00001192

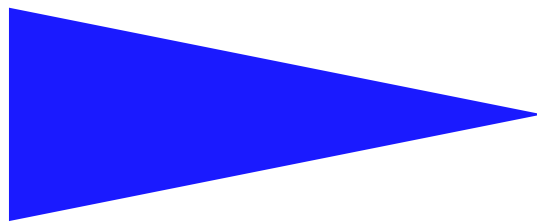
<https://hal.inria.fr/inria-00001192>

Submitted on 3 Apr 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

PUBLICATION
INTERNE
N° 1795



SELF-ADJUSTING ATOMIC MEMORY FOR DYNAMIC SYSTEMS
BASED ON QUORUMS ON-THE-FLY.
CORRECTNESS STUDY

EMMANUELLE ANCEAUME MARIA GRADINARIU
VINCENT GRAMOLI ANTONINO VIRGILLITO

Self-Adjusting Atomic Memory for Dynamic Systems based on Quorums On-The-Fly. Correctness Study

Emmanuelle Anceaume Maria Gradinariu
Vincent Gramoli Antonino Virgillito^{*}

Systèmes communicants
Projet Adept

Publication interne n1795 — Avril 2006 — 15 pages

Abstract: Atomic memory is a fundamental building block for classical distributed applications. The new emergent large scale applications such as e-auctions or e-commerce need similar fundamental abstractions able to cope with the highly dynamicity of p2p environments. In this paper we prove the correctness of a Self-Adjusting Atomic Memory that implements multi-writer/multi-reader atomic operations in dynamic systems. The architecture of this system was introduced in [5, 6]. The self-healing property guarantees the capability to cope with replica volatility (unavoidable in a p2p environment), while the self-adjusting property helps it to adapt on-the-fly to the dynamicity of client access.

Key-words: Distributed algorithm, atomicity, self-adjusting, self-healing, replication, dynamism, scalability, quorums on the fly.

(Résumé : tsvp)

^{*} Università di Roma, "La Sapienza", Italy

Correction de SAM : Une mémoire atomique auto-ajustable pour les systèmes dynamiques basée sur des quorums à la volée

Résumé : Ce papier propose la preuve de correction de SAM - une mémoire atomique auto-ajustable basée sur l'utilisation de systèmes de quorums. L'architecture de SAM a été décrite dans [5, 6]. Nous détaillons également un des modules du système qui permet de réaliser la consultations "à la volée" des quorums en lecture ou écriture.

Mots clés : atomicié, auto-ajustement, réplication, scalabilité, dynamicité, quorums à la volée

1 Introduction

The notoriety of peer-to-peer (p2p) file sharing guarantees the subsequent success of this technology in commercial applications. Some of these applications require strong computational guarantees. Internet-scale applications such as *e-auction* or *e-booking* for example, require linearizable read/write operations. That is, an auctioneer must be able to read other bids as well as writing its own bid. Alternatively, booking transaction needs persistent record. More generally, these applications need an atomic memory service able to cope with the dynamic nature of the system.

Designing an atomic memory service in p2p systems faces several problems. P2p systems are by their nature ad-hoc distributed systems without any organization or centralized control. Unlike classical distributed systems, p2p systems encompass processes (peers) that experience highly dynamic behaviors including spontaneous join and leave or change in their local connections. The high dynamicity of the network has an important impact on data availability. The use of classical distributed computing solutions, like replication for example, introduces an extra cost related to: (1) maintaining a sufficient number of replicas despite frequent disconnections, and (2) maintaining the consistency among replicas. The former problem can be solved using *self-healing* techniques while the latter finds solutions in the use of *dynamic quorums* (intersecting sets).

Another issue posed by the dynamicity of the network is the replica stress (load). An inadequate number of replicas may have tremendous impact on the replica access latency since the access latency increases with the access rate. Moreover, due to limitations of the local buffers' size a non negligible fraction of replica accesses might be lost. Consequently, the number of replicas should spontaneously adjust to the access rate.

Related Work Starting with Gifford's weighted votes [11], quorum systems [3, 16, 8, 28] have been widely used to provide consistency. Several quorum-based approaches provide mutual exclusion [22] or shared memory emulation [7]. Recently, quorum-based implementations of atomic memory for dynamic systems have been proposed in [17, 9, 12]. All these papers have a common design seed—they use reconfigurable quorum systems, work pioneered by Herlihy [14], for static distributed systems. In such systems, clients have to know the set of replicas, and thus have to participate to the reconfiguration process. In [20, 10] the authors showed that using two quorums systems concurrently preserves atomicity. This result has been later exploited in the implementation of the reconfigurable quorum systems for highly dynamic systems. That is, periodically the system proceeds to modifications of the current quorums set (referred as configuration).

The authors in [5] follow an alternative approach for implementing atomic memory in dynamic systems. This approach is based on recent achievements in the context of dynamic quorums and logical overlays. Dynamic quorums have been mainly investigated in [26, 2, 23]. Naor and Wieder [26] sought solutions for deterministic quorums using dynamic paths in a planar overlay [24]. Simultaneously, probabilistic quorums were proposed by Abraham and Malkhi [2] based on an overlay designed as a dynamic approximation of De Bruijn graphs [1]. Recently, in [28], the authors discuss the impact of dynamism on the multi-dimensional quorum systems for read-few/write-many replica control protocols. They briefly describe strategies for the design of multi-dimensional quorum systems that combine local information in order to deal with frequent replica joins and leaves and quorum sets caching in order to reduce the access latency. AndOr strategies [25] are studied in [23] in order to implement fault-tolerant storage in dynamic environment. In [21], the authors propose a solution that ensures atomic access to non replicated data in a CAN network using restrictive hypothesis: nodes nicely leave the system by notifying their neighbours. It is also assumed that nodes do not fail and the network is reliable.

Our Contributions In this paper we propose and prove the correctness of SAM, an atomic memory for dynamic systems with self-adjusting and self-healing capabilities. SAM construction brings together several new and old research areas exploiting the best of these worlds: logical overlays, dynamic quorums and replica control. SAM architecture, which is detailed in [5], is composed of three interconnected modules, each being designed to serve a specific task (i.e., data availability, linearizability and load balancing of replicated data). Object availability is guaranteed through replication among several system nodes. Replicas of the same object define a torus overlay, structure that has been proved efficient in the design of quorum systems ([15, 25]), that is, a system made of a set of subsets of nodes, such that every two of which intersects. The behavior of the atomic memory is emulated through dynamic quorum sets sampled from a deterministic overlay traversal. Unlike quorums based on dynamic paths [26] which avoid the holes in the overlay by following alternative longer paths, our traversal fully exploits the self-healing capabilities of the overlay in order to probe a minimal number of nodes.

We propose an algorithm that implements such quorum-based phases, and prove that the read/write operations are atomic—i.e., linearizable—despite asynchronism, node arrivals, node failures, and message losses. A particularity of our approach is the use of a single round-trip communication phase for read operations. This improves the efficiency of the atomic memory when read operations are frequent compared with write operations. Furthermore, our algorithm guarantees self-adjustment by expanding or shrinking the overlay with respect to the dynamicity of client accesses. We show that atomicity is not broken even when new replicas have to be added or removed from the overlay. Our work is hybrid between the reconfiguration-based systems and strategic adaptive systems. Indeed, clients or replicas do not have to be aware of the reconfiguration process. Only lightweight reconfiguration is used to achieve self-healing and self-adjusting properties, which is further exploited to design adaptive strategies for sampling read/write quorums, and balancing the replicas stress. All these properties—atomicity, self-adjustment and self-healing—are guaranteed using only constant size local information and, as far as we know, no atomic memory has been proved correct in such circumstances.

Road Map The paper is organized as follows. Section 2 introduces the system model and the problem definition. An overview of SAM is presented in Section 3. The module that handles read and write atomic operations is formally described in Section 4, and then proofs of correctness are presented in Section 5. In Section 6, we conclude and present some future research topics.

2 Model and Problem Definition

System. We consider a dynamic system \mathcal{DS} as the tuple $\mathcal{DS} = (I, X)$, where I is a set of possibly infinite nodes, and X is an unbounded universe of shared data, referred in the following as *objects*. \mathcal{DS} is subject to unpredictable changes. Nodes can leave or join the system arbitrarily often. Departure includes crash failure while recovery of a failed node is treated like a new arrival. A node is called *active* if it has joined and did not fail or leave the system since then. Otherwise the node is called *failed*. For each node, a *local history* is defined as a sequence of *events*, which consists of *connect*, *disconnect*, *sending* a message to other nodes, *receiving* a message from another node, taking computation steps, and possibly a *crash* event. If a crash event exists, it is always the last event in the local history. A collection of local histories, one for each node, is called an *execution*. Only *well formed executions* are considered, that is, if a message is received in the execution, then it was also sent during this execution.

Each object $x \in X$ has a single owner node and is replicated at some other nodes in the system. By abuse of notation, a replica of x refers to a node maintaining a copy of x . For efficiency reasons the number of replicas of an object is strictly inferior to the number of nodes in the system. Each object has a *value* that can be consulted/modified as a result of a read/write operation. It is assumed that each object has a predefined initial value. In order to read/write the value of an object a client node starts probing a set of this object's replicas. A read/write operation starts when the first replica in the set is contacted and completes when all the replicas in the set has been contacted.

Self-adjusting and self-healing atomic memory. This work aims at providing a self-adjusting and self-healing atomic memory. Atomicity is often defined in terms of an equivalence with a serial memory. In the following, the definition proposed in [19] is adopted:

Definition 1 (Atomicity) *If all the read and write operations that are invoked complete, then the read and write operations for object x can be partially ordered by an ordering \prec , so that the following conditions are satisfied:*

- *No operation has infinitely many other operations ordered before it;*
- *The partial order is consistent with the external order of the invocations and responses, that is, there does not exist read or write operations π_1 and π_2 such that π_1 completes before π_2 starts, yet $\pi_2 \prec \pi_1$;*
- *All write operations are totally ordered and every read operation is ordered with respect to all the writes;*
- *Every read operation ordered after any write returns the value of the last write preceding it in the partial order; any read operation ordered before all writes returns the initial value of the object.*

In order to be operational in a dynamic environment, two additional properties are required from an atomic memory: *self-healing* and *self-adjusting*. Self-healing aims at ensuring the availability of an object whenever failures occur while self-adjustment aims at expanding or restraining the number of replicas with respect to the dynamicity of the memory load. The load of a memory is defined regarding to the load of each of its replica in the following manner:

Let $\mathcal{L}_i(t)$ be the number of operations that some replica i has to process at time t . $\mathcal{L}_i(t)$ is referred in the following as the local load of replica i at time t . Let T_{min}^i and T_{max}^i be two application dependent parameters which define a lower and upper bound on the load of i . Replica i is *overloaded* (resp. *underloaded*) iff $\mathcal{L}_i(t) \geq T_{max}^i$ (resp. $\mathcal{L}_i(t) \leq T_{min}^i$).

Definition 2 (Self-adjusting Atomic memory) Let $DS = (I, X)$ be a dynamic system and α be an execution of DS . Let x be an object in X and let M_x be the set of x replicas. For each i in M_x let T_{min}^i and T_{max}^i the application dependent load boundaries. M_x is a self-adjusting atomic memory with respect to the load iff:

- **atomicity** α restricted to M_x verifies the atomicity definition (see Definition 1)
- **self-adjustment** $\forall i \in M_x, \exists t_{0_i}$ such that $\forall t > t_{0_i}, T_{min}^i \leq \mathcal{L}_i(t) \leq T_{max}^i$.

3 SAM Overview

SAM aims at emulating an atomic memory on top of a replicated system with self-adjusting and self-healing capabilities. This section presents an overview of SAM.

3.1 SAM Logical Overlay

Replicas of an object, referred in the following as memory, share a same logical overlay, organized in a torus topology (as for example CAN [27]). Basically, a 2-dimensional coordinate space $[0, 1) \times [0, 1)$ is shared by all the replicas of an object. A replica is responsible of a set of zones, each of those being a rectangle in the plane. The entrance and departure of a replica dynamically changes the decomposition of the zones. These zones are rectangles (union of rectangles) in the plane. Replicas owners of adjacent zones are called neighbors in the overlay and are linked by virtual links. The overlay has a torus topology in the sense that the zones over the left and right (resp. upper and lower) borders are neighbors of each other. Initially, only the owner of the object is responsible for the whole space. The bootstrapping process pushes a finite, bounded set of replicas in the network. These replicas are added to the overlay using well-known strategies [27, 26]: the owner of the object specifies randomly chosen points in the logical overlay, and the zone in which each new replica falls is split in two. Half the zone is left to the owner of the zone, and the other half is assigned to the new replica. The object initial value v_0 is replicated at the new replica. In the following we omit a more detailed description of the bootstrapping process. Note that an interesting point to explore here is the introduction of efficient incentive mechanisms to motivate nodes to host replicas (i.e., to be part of the atomic memory). Techniques from game theory or mechanism theory can be used to this end, however this topic is beyond the scope of this paper.

3.2 SAM Architecture

SAM is specified in Input/Output Automata (IOA) Language [18] and is structured as the composition of three main automata: the *Traversal* $_{x,i}$, the *LoadBalancer* $_{x,i}$, and the *Adjuster* $_{x,i}$ automata, where $i \in I$ and x is the considered object. A fourth automaton, the *CommunicationLink* $_{x,i,j}$, represents the communication medium between any two replicas i and j . In the remaining of the paper, we restrict our attention to only one object x . Thus, the x subscript is omitted. Relationship among the four automata in terms of input/output actions is depicted in Figure 1.

Briefly, the LoadBalancer automata scans the overlay to identify non overloaded replicas. It determines the strategies for reducing the memory load at nodes, and includes an operation aggregation strategy to overlap operation executions. The Traversal automaton is in charge of maintaining the consistency of the overlay applying appropriate strategies for executing linearizable operations on the replicas. Finally, the Adjuster assigns logical responsibility zones to physical replicas and maintains this correspondence consistent. It also handles the expansion or shrink of the overlay whenever requested by the Load Balancer, and the departure of non responding replicas.

Replicas are accessed by clients through read and write operations on the object. Such operations are implemented in the Traversal automaton. A read operation consists in traversing the overlay following a horizontal trajectory that wraps around all the overlay. A write operation consists in traversing the overlay following a horizontal trajectory and then a vertical one. An horizontal traversals defines a *consultation quorum*, while a vertical one is a *propagation quorum* (see Figure 2). More details are given in the following Section.

SAM starts with a read/write request from a client. A client submits a read/write request to one of the replicas. The Load Balancer of this replica does not immediately initiate a traversal but rather records this request in a local

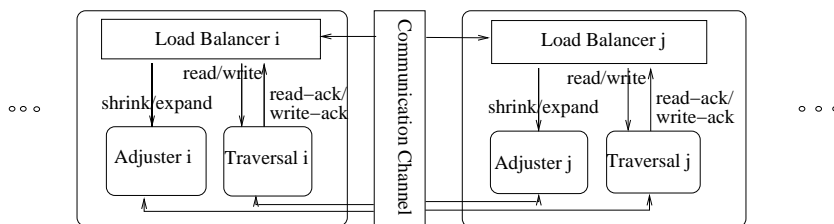


Figure 1: Overview of SAM

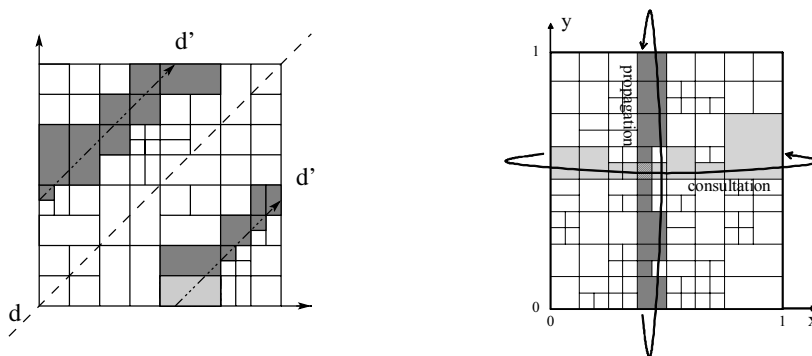


Figure 2: A Thwart Path (left grid) and a Traversal Path (right grid)

heap. Each replica periodically scans its heap to pick the requests for which a traversal has to be initiated. The strategy to pick these requests consists in choosing only one write (or read) operation among all the pending ones. More precisely, a write operation is initiated only for the most recently heaped write request (if any), while a read traversal is initiated for one of the heaped read request (if any). Old write operations can be safely discarded (no write traversals are initiated for them) as they will not influence anymore the state of the object. There is no such constraint for a read. Once a traversal is initiated, the initiating replica empties its heap, after having kept track of all the read/write requests to later return the status of the operations to the requesting clients. Clearly, this request aggregation strategy reduces the number of requests that are actually served and thus the latency of the memory.

Despite the former aggregation strategy, the number of heaped requests at an initiating replica still can grow very fast, incurring in a local overload as defined in Section 2. If a replica load catches up with a threshold T_{max} then the Load Balancer invokes a *thwart* process to find a suitable replica. The thwart process, as shown in Figure 2 shown in the appendix, checks the memory along a diagonal trajectory for a non overloaded replica. The diagonal trajectory is chosen to contact at each step a replica involved in different traversals. If such a replica is found, then it becomes the initiating replica for these requests, otherwise if the memory needs to be expanded (no overloaded replica has been found), the Load Balancer activates the Adjuster to add another replica (i.e., , it requests a node in the system to host a replica, hence to become a SAM member ¹).

Alternatively, when the load of a replica is below T_{min} , the Load Balancer invokes the Adjuster for a shrink procedure: the replica is removed from the memory. This prevents probing complexity from remaining uselessly high when a burst of requests has ended. When a replica leaves (voluntarily or not), the zone is locally healed by relying on a strategy similar to the one proposed in CAN. On the other hand, joins are triggered by SAM as follows: a replica is inserted within the quorum system only when it is required (i.e., , for expansion purpose). A complete specification of the Load Balancer and Adjuster is given in the Appendix.

¹As mentioned in the previous section, the current work does not focus on mechanisms for motivating nodes to host replicas

4 The Traversal Automaton

The Traversal is the core of SAM activity. It is in charge of maintaining the consistency of the replicas by applying appropriate strategies for executing linearizable operations via intersected sets, called *quorums*. In this Section we give a complete specification of the Traversal automaton.

4.1 Read/Write Operations

SAM guarantees the atomicity of read/write operations using a quorum system referred in the following as the *tiling quorum system*. The key notions in the definition of a tiling quorum system are the *horizontal* and *vertical* tiling sets. A set of replicas in SAM define an horizontal (resp. vertical) tiling if their responsibility zones are pairwise independent and totally cover the abscissa (resp. ordinate) of the coordinate space shared by replicas in SAM.

Definition 3 (Tiling Quorum System) A tiling quorum system \mathcal{T} is a couple $\langle \mathcal{C}, \mathcal{P} \rangle$ where \mathcal{C} is a set of horizontal tilings and \mathcal{P} a set of vertical tilings such that any horizontal tiling intersects any vertical tiling.

Let $\mathcal{T} = \langle \mathcal{C}, \mathcal{P} \rangle$ be a tiling quorum system. Elements in \mathcal{C} are referred to in the following as *consultation quorums*, while elements in \mathcal{P} are referred to as *propagation quorums*. Figure 2 in the appendix shows a consultation and a propagation quorum as well as the corresponding traversals.

Read and write operations in SAM are multi-phase operations as in [7, 17, 12, 9]. While in [7, 17] a read operation is realized in two phases, SAM, similarly to GeoQuorums [9] exploits single-phased read operations, for efficiency reasons. Then, the read operation consists in only one consultation phase, where a consultation quorum is probed. That is, the overlay is traversed horizontally (from west to east) and all the nodes encountered are requested for their object *value* and *tag* (the value's timestamp). Finally, the most updated value is returned by the read. The write strategy contains the same consultation phase to get the most recent tag, plus a propagation phase, where the value to write is written on a propagation quorum with an updated version number. To differentiate the consultation phases of read and write operations, the latter is called the *update phase*.

One-phase reads might violate linearizability in the following way: a read operation executed concurrently with a write operation may consult a fresh value not completely propagated, while a following read may consult an old value. The problem comes from the early termination of a read operation that consults a fresh value which is still being propagated. This violates the atomicity definition (see Definition 1).

GeoQuorums solves this problem by including an additional phase in the write operation, in which the initiator of a write sends a specific confirmation message when the write operation is completed. Differently, in SAM propagation terminates before termination of the read. This is achieved by using a lock mechanism to let either another read operation consult this newly written value or to order this read operation before the write. More precisely, in the propagation phase the overlay is traversed in a orthogonal direction with respect to consultation (north and south). Propagation proceeds in both senses so that each replica is visited twice in this phase: the first time the object is locked, preventing concurrent reads to return the updated value too early, while the second time the lock is released, indicating the completion of the write operation. In dynamic systems an object may be locked forever due to unforeseen leaves. We deal with this problem by assuming the use of a leasing strategy [13]. In order to guarantee the termination of the consultation phase when it crosses an infinite sequence of propagations we use the following mechanism: each consultation arriving at a replica takes a snapshot of the locally pending propagations and waits only for those propagations to end. Hence, in a finite number of steps each consultation phase terminates.

4.2 Automaton States

The state variables of the Traversal automaton are described for node i from line 7 to 37 of I/O Automaton 1. First of all, each replica maintains a tag and a value of the object (*tag* and *val* fields). Those fields are updated when the replica participates in a write operation involving the quorums it belongs to. The *op* record (cf. lines 16–29) describes an operation. It contains some identifying subfields such as its identifier *id*, its *type* which indicates if the operation associated is a read or a write, and *intr* and *intr-zone* informing about the identifier of the replica that started the operation and the starting zone. The *senses* and *phase* fields indicate respectively the sense (*E*, *N* or *S*) and the phase of the traversal in progress, while *zone* is the zone involved in the operation, among the ones managed by the replica. Finally, two last fields are required by the lock mechanism described above: *locking-prop* and *rcv-prop*. The

I/O Automaton 1 *OperationManager_i*

1: Domain:	4: $\Pi \subseteq I \times \mathbb{N}$, the set of all operation identifiers.
2: $I \subseteq \mathbb{N}$, the set of node identifiers.	5: $T \subseteq I \times \mathbb{N}$, the set of all tags.
3: V , the set of all possible values of an object.	6: M , the set of all possible messages.
7: Signature:	11: $\text{rcv}(m)_{j,i}$, $i, j \in I, m \in M$
8: Input:	12: Output:
9: $\text{read-write}(\text{type}, \text{id}, v)_i$, $i \in I, \text{type} \in \{\text{read}, \text{write}\}$,	13: $\text{snd}(m)_{i,j}$, $i, j \in I, m \in M$
10: $\text{id} \in \Pi, v \in V$	14: $\text{read-write-ack}(\text{id}, v)_i$, $i \in I, \text{id} \in \Pi, v \in V$
15: State:	27: $\text{id} \in I$
16: op a record with fields	28: $\text{val} \in V$, initially \perp
17: $\text{id} \in \Pi$	29: $\text{senses} \subseteq \{N, S, E\}$
18: $\text{intr} \in I$	30: ops , the set of all known operations
19: $\text{zone} \in \mathbb{R}^4$	31: tag , a record with fields
20: $\text{intr-zone} \in \mathbb{R}^4$	32: $ct \in \mathbb{N}$
21: $\text{type} \in \{\text{read}, \text{write}\}$	33: $\text{id} \in I$
22: $\text{phase} \in \{\text{cons}, \text{update}, \text{prop}, \text{ending}\}$	34: $\text{val} \in V$, initially v_0
23: $\text{rcv-prop} \subseteq \{N, S, E\}$	35: failed a boolean, initially false
24: $\text{locking-prop} \subseteq \Pi$, initially \perp	36: replica a boolean
25: tag , a record with fields	37: adjusting a boolean
26: $ct \in \mathbb{N}$	

former indicates, for a consultation phase, the set of pending propagations that blocks it. The latter indicates, for a

38: Transitions:	73: Input $\text{rcv}(\langle \text{id}, \text{phase}, \text{intr-zone}, z, t, v, \text{sense} \rangle)_{j,i}$
39: Input $\text{read-write}(\text{type}, \text{id}, v)_i$	74: Effect:
40: Effect:	75: if $\neg \text{failed} \wedge \text{replica} \wedge \text{phase} \neq \text{idle}$ then
41: if $\neg \text{failed} \wedge \text{replica}$ then	76: $op = op' \in ops$ st. $op'.\text{id} = \text{id}$ or \perp
42: $op.\text{id} \leftarrow \text{id}$	77: $op.\text{phase} \leftarrow \text{phase}$
43: $op.\text{type} \leftarrow \text{type}$	78: $op.\text{zone} \leftarrow z$
44: if $op.\text{type} = \text{read}$ then	79: if $op.\text{phase} = \text{cons}$ then
45: $op.\text{phase} \leftarrow \text{cons}$	80: $op.\langle \text{tag}, \text{val} \rangle \leftarrow \max(\langle \text{tag}, \text{val} \rangle, \langle t, v \rangle)$
46: for $op1$ st. $op1.\text{phase} = \text{prop}$ do	81: for $op1 \in ops$: $op1.\text{phase} = \text{prop} \wedge$
47: $op.\text{locking-prop} \cup \leftarrow \{op1.\text{id}\}$	82: $ op1.\text{rcv-prop} = 1$ do
48: else	83: $op.\text{locking-prop} \cup \leftarrow \{op1.\text{id}\}$
49: $op.\text{phase} \leftarrow \text{update}$	84: if $op.\text{intr} = i$ then
50: $op.\text{val} \leftarrow v$	85: $op.\text{phase} \leftarrow \text{ending}$
51: $op.\text{intr} \leftarrow i$	86: else if $op.\text{phase} = \text{update}$ then
52: $op.\text{zone} \in i.\text{zones}$	87: $op.\text{tag} \leftarrow \max(\text{tag}, t)$
53: $op.\text{intr-zone} \leftarrow op.\text{zone}$	88: if $op.\text{intr} = i$ then
54: $op.\text{senses} \leftarrow \{E\}$	89: $op.\text{phase} \leftarrow \text{prop}$
55: $ops \leftarrow ops \cup \{op\}$	90: $op.\text{tag} \leftarrow \langle op.\text{tag}.ct + 1, i \rangle$
	91: $op.\text{senses} \leftarrow \{N, S\}$
56: Input fail_i	92: else if $op.\text{phase} = \text{prop}$ then
57: Effect:	93: $\langle \text{tag}, \text{val} \rangle \leftarrow \langle op.\text{tag}, op.\text{val} \rangle$
58: $\text{failed} \leftarrow \text{true}$	94: $op.\text{rcv-prop} \leftarrow op.\text{rcv-prop} \cup \{\text{sense}\}$
	95: if $ op.\text{rcv-prop} = 2$ then
59: Output $\text{snd}(\langle \text{id}, \text{phase}, \text{intr-zone}, z, t, v, \text{sense} \rangle)_{i,j}$	96: for $op2$: $op2.\text{phase} = \text{cons}$ do
60: Precondition:	97: $op2.\text{locking-prop} \leftarrow \{op.\text{id}\}$
61: $\neg \text{failed} \wedge \text{replica} \wedge \text{freezing} = \emptyset$	98: if $op.\text{intr} \neq i$ then
62: $op.\text{locking-prop} = \emptyset$	99: $op.\text{senses} \leftarrow op.\text{senses} \cup \{\text{sense}\}$
63: $op \in ops$	100: $ops \leftarrow ops \cup \{op\}$
64: $\text{phase} = op.\text{phase} \notin \{\text{ending}, \text{idle}\}$	
65: $\text{id} = op.\text{id}$	101: Output $\text{read-write-ack}(\text{id}, v)_i$
66: $\langle t, v \rangle = op.\langle \text{tag}, \text{val} \rangle$	102: Precondition:
67: $\text{sense} \in op.\text{senses} \neq \emptyset$	103: $op \in ops$
68: $\text{intr-zone} \leftarrow op.\text{intr-zone}$	104: $op.\text{phase} = \text{ending}$
69: $z = \text{get-trv-zone}(\text{sense}, \text{intr-zone}, op.\text{zone})$	105: $\text{id} = op.\text{id}$
70: $j = \text{get-responsible}(z)$	106: $v = op.\text{val}$
71: Effect:	107: $op.\text{intr} = i$
72: $op.\text{senses} \leftarrow op.\text{senses} \setminus \{\text{sense}\}$	108: Effect:
	109: $op.\text{phase} \leftarrow \text{idle}$

propagation phase, the senses from where messages have been received. In overall, the number of elements in this set informs whether or not the propagation phase is pending.

Messages in a traversal are sent from neighbor to neighbor, starting at the initiator replica $op.intr$ of the operation. The messages sent during the traversal contain the operation identifier id , type $type$, initiator zone $intr-zone$, but also information about where it has to be sent (regarding to $sense$ and current zone z). Other fields that are related to the Traversal automaton are the booleans $failed$ and $replica$, indicating respectively whether the current node is crashed or not and whether it owns a copy of the object or not, and the set ops containing the operations the replica keeps track of. Finally, the $adjusting$ field is needed in case of memory adjustment, due to nodes joining or leaving the memory. If it is true, this indicates that an adjustment is pending at the current node, and no traversal participation for that node is possible before it recovers its state.

4.3 Automaton Transitions

The transitions of the *Traversal* automaton are described for node i in I/O Automaton 1, lines 38–109. A read or write operation starts at node i in the consultation phase as a result of a read-write $_i$ input event. This event is triggered by the LoadBalancer automaton and sets the op subfields. More precisely it initiates the starting phase (update or cons) sense to E , the initiator identity, the current zone ($zone$), the initiator zone ($intr-zone$), and the $locking-prop$ field if any propagation phase are pending at the current replica. Notice that it is the responsibility of the *LoadBalancer* to assign a new unique identifier to this operation, whereas the former event initializes other operation subfields.

The traversal proceeds with messages being sent by a node i to one of its east neighbors j , through a $snd_{i,j}$ event. Such an event is triggered only if i is not $adjusting$, $failed$ and is a $replica$. Note that if the $op.locking-prop$ is not empty, then the snd can never occur (cf. precondition line 62). This means that even if a replica has received a message corresponding to a traversal, while it is failed, blocked or adjusting, it does not participate by sending new messages.

When a replica i receives a message from a *Traversal* $_j$ automaton by a $recv_{j,i}$ input event, it checks (i) if it is the initiator of the operation ($op.intr = i$) and (ii) the current phase of this operation. Let ϕ be this phase.

When all the phases of an operation are done, the initiator sets the $op.phase$ subfield to ending. This occurs when the initiator of the operation receives the consultation message from the east direction, or when it has received propagation messages from both north and south senses ($op.rcv-prop = \{N, S\}$). For an update phase to complete, the initiator receives an east message and changes $op.phase$ to prop. There is a singular difference between propagation phase and update/consultation phases in the way tag and value are altered.

In all cases the replica receiving a message, updates either the operation $op.\langle tag, val \rangle$ pair or it updates its local $\langle tag, val \rangle$ pair. The most up-to-date pair is conveyed by messages as the $\langle t, v \rangle = op.\langle tag, val \rangle$ pair (l. 66). In case ϕ is a consultation or an update phase, the tag information is simply the most up-to-date one encountered: it is the maximum between the local one and the received one (l.80, 87). If ϕ is a consultation phase, the chosen value is the value related to this tag, whereas if ϕ is an update phase, the value remains the one initially set as the value to write. Finally if ϕ is a propagation then the local pair is updated with the received one (l. 93), since the value included is the value to be written.

Next, a consultation receipt at a replica i might make this consultation block. More precisely if there are pending propagations at i when the consultation receipt occurs (l.81–83), the operation $locking-prop$ field is filled which prevents further snd event from occurring. From this point on, each time one of the pending propagations completes, an element is removed from the $locking-prop$ set (l.97). When this set becomes empty again, the snd event is newly enabled, provided that i is still an active non-adjusting replica.

Moreover, the $op.sense$ subfield is set to the sense of the phase message and op field (including all its subfields) is stored in ops (l.100). Each replica sends messages with information about operations op in a sense belonging to $op.senses$ set by the mean of the snd action. Since an initiator starting a consultation or an update phase has its op field set to $\{E\}$, the first message is sent in the east sense. When a non-initiator participating replica receives the phase message, it records the op subfields to resend it simply in the same sense. If a propagation phase starts at an initiator, this one must send two messages: one in the north sense and the other in the south sense.

Note the presence of the failed input action. This action is triggered by the environment and indicates that a crash has occurred. Once this flag is set to true any other action is disabled.

Finally, we briefly describe two functions that are not specified in the *Traversal* automaton. Function $get-trv-zone$ which, given a sense (N, S, E) , the traversal initiator zone, and the operation current zone $op.zone$, returns the next zone in this sense belonging to one of the possible trajectory of the traversal (this zone is crossed by an horizontal line

also crossing the initiator zone). If many possible replicas exist it chooses one among all. Function `get-responsible` which, given a zone, returns the neighbor responsible for this zone.

5 SAM Correctness

This section proves SAM's correctness. Hence, we show that SAM is a self-adjusting atomic memory as defined in Section 2. This proof is divided in three parts. First, Theorem 5.1 shows that operations contact tiling quorum of replicas. Second, the atomicity is proved in Theorem 5.6. Finally, we prove that SAM self-adjusts (Theorem 5.7).

Assumptions. The *never-partitioning* property ensures that a set of connected nodes is never partitioned, while the *failure-detection* property ensures that a failure is detected without any mistake. We assume that those two conditions are satisfied by our overlay of replicas.

We do not show explicitly how a replica gets “a failed replica locking-prop state back”. When a replica fails and a new replica takes over its zone, or when a new replica is created the corresponding active node needs to get back *locking-prop* information. Given the *never-partitioning* property, the corresponding replica contacts its north and south neighbors to get the highest *tag*, *val*, *locking-prop*, *senses* values, and information about their phase participation before participating in any message exchange. The choice of an healing replica is done using CAN overlay mechanism and is not formally specified here.

We also assume that any sequence of external actions for replica i and object x is *well-formed*. That is (i) the first event of the sequence is either a read-write_i event or a fail_i event; (ii) a read-write_i event is immediately followed by the matching read-write-ack_i ; (iii) no fail_i event precedes any read-write_i , or read-write-ack_i event.

For the sake of stabilization of our approach, we assume the following two properties: (i) the *potential-replication* property ensures that there exists at least one active node out of the memory at any time; (ii) the *thwarting-access-rate* property means that during any thwart, the access rate (rate at which requests are received) is lower than the logical treatment rate (rate at which requests are started). If it exists a replica accepting to treat the thwarted request then the duration time of the thwart is the period between the first sent in the thwart and the acceptance time. Otherwise, it is the period between the first sent and the time the expand resulting from the last receipt (when the thwart has wrapped around the torus) occurs.

Definitions. We refer to **tag** as a mapping from TId to $\mathbb{N} \times I$ such that $op.id$ is mapped to $tag(op.id)$ if at least one of the following condition holds: (i) if $op.type = \text{read}$ then $tag(op.id) = tag_i$ when read-write-ack_i occurs, (ii) if $op.type = \text{write}$ then $tag(op.id) = tag_i$ immediately after $\text{rcv}((*, \text{prop}, \dots))_{j,i}$ occurs. Given this definition, we totally order operations by their tag, respecting the lexicographical order.

We define **inv** as a mapping from an operation to \mathbb{R}^+ such that operation π is mapped to time τ if $\text{read-write}(*, id, *)$ occurs at time τ with $\pi.id = id$.

We define **resp** as a mapping from an operation to \mathbb{R}^+ such that operation π is mapped to time τ' if $\text{read-write-ack}(id, *)$ event occurs at time τ' where $\pi.id = id$.

Next, let \prec be a total order capturing the **real-time precedence** on every event. Second, we define the history precedence, namely $<_H$, as an irreflexive partial order between operations π_1 and π_2 , such that $\pi_1 <_H \pi_2$ if and only if $\text{resp}(\pi_1) \prec \text{inv}(\pi_2)$.

Finally, we restate the **atomicity** definition of Lynch [19] using the aforementioned notations. If it exists a matching $\text{read-write-ack}_{x,i}$ event following any $\text{read-write}_{x,i}$ event then, it exists a relation $<_S$ that orders partially read and write operations such that: (i) $<_H \subseteq <_S$. (ii) For any write operation π_1 and π_2 , either $\pi_1 <_S \pi_2$ or $\pi_2 <_S \pi_1$. (iii) For any read operation π_1 and write operation π_2 , either $\pi_1 <_S \pi_2$ or $\pi_2 <_S \pi_1$. Moreover, if it exists π_2 such that $\pi_2 = \max_{<_S} \{\pi'_2 <_S \pi_1\}$ and π_2 writes value v then π_1 returns the same value v . And if no such π_2 exists, then the value returned by π_1 is v_0 . (iv) For any read or write operation π , the set $\{\pi' : \pi' <_S \pi\}$ is finite.

We define a **partial order** $<_t$ on the set of traversals such that (i) write traversals are totally ordered: a write traversal π_1 precedes another write traversal π_2 if $tag(\pi_1) < tag(\pi_2)$. (ii) read traversals are ordered between write traversals: a read traversal π_1 is ordered after all write traversals π_ℓ such that $tag(\pi_\ell) \leq tag(\pi_1)$ and before all write traversals $\pi_{\ell'}$ such that $tag(\pi_1) < tag(\pi_{\ell'})$.

Quorums properties. In this section we prove that for a given phase, our traversal mechanism involves a quorum of nodes. That is the phase set of participating replicas verifies each of the properties defining a quorum. First of all, we say operation π 's phase completes in one of the three following cases holds: (i) $rcv(\pi.id, cons, \dots)$ occurs and $op.phase$ is set to ending, (ii) $rcv(\pi.id, update, \dots)$ occurs and $op.phase$ is set to prop, or (iii) $rcv(\pi.id, prop, \dots)$ occurs and $op.phase$ is set to ending. We define the initiator of an operation π the node where π starts. Typically i is the initiator of operation π if $read-write(*, \pi.id, *)$ occurs at node i .

Theorem 5.1 *If a phase ϕ completes, all elements of a quorum Q have received messages such that:*

- *if ϕ is a propagation phase then Q is an propagation quorum.*
- *if ϕ is a consultation phase or an update then Q is a consultation quorum.*

Proof. First we show that all messages of ϕ are sent in a non-varying sense. When ϕ starts, the initiator sends messages in each sense included in the *senses* set. A non-initiator replica changes its field only by a *rcv* (or by a *snd* event to make sure messages are not resent) and it adds the receipt sense to its *senses* set. The initiator *senses* is set or updated only when a new phase starts: *phase* is set to *cons* or *update* and *senses* to $\{E\}$, or *phase* is set to *prop* while *senses* is set to $\{N, S\}$, or *phase* is set to ending and *prop* or ending.

Second, we show that replicas are contacted following a chain of adjacent zones. By examination of the *snd* action, if a $snd(*, *, *, z, *, *, sense)_{i,j}$ event occurs then z belongs to $j.zone$. Next, by $get-trv-zone_i$ function, we know that zone z is chosen such that it abuts the current zone $op.zone_i$ along the $sense_i$ axis. Lastly, the corresponding $rcv(*, *, *, z, *, *, sense)_{i,j}$ sets the current phase to z received since j is an active replica.

From now on, we show that it exists a common ordinate (resp. abscissa) to every zone, the consultation or update (resp. propagation) phase goes through. The proof follows directly from the *intr-zone* field and the definition of the $get-trv-zone$ function. First, notice that field $op.intr$ is initially set to the current node zone by a read-write event and is not modified in any further state. By definition of $get-trv-zone$ and the *snd* action, the next participant's zone is such that it exists a horizontal (resp. vertical) line that crosses it and that also crosses the $op.intr-zone$.

All these three properties leads to the conclusion. If ϕ is a consultation, the sense followed by messages is E , all zones crossed over are adjacent, and there exists a common abscissa between them, i.e., ϕ contacts a set of participants consisting of a horizontal tiling. The vertical one is shown following the same reasoning. The intersection is straightforward from tiling definition. \square

Atomicity proof. Here, we show that the partial order relation $<_t$ exists. The first lemma shows that no locked replica can participate in any read operation until it is unlocked. The second one shows that a replica receives the second propagated message after all participants have already received the first one.

Lemma 5.2 *If $|\pi_1.rcv-prop_i| = 1$ at time τ such that $rcv(\pi_2.id, cons, *, *, *, *, *)_i$ occurs at time $\tau' > \tau$ with $i \neq \pi_2.intr$, every participant or π_1 's propagation have already stored the propagated value, tag pair at time $resp(\pi_2)$.*

Proof. When $rcv(\pi_2.id, cons, *, *, *, *, *)_i$ occurs, replica i knows about operation π_1 . By assumption, it exists $\pi_1 \in ops$ such that $|\pi_1.rcv-prop_i| = 1$. By examination of the code, since $\pi_2.phase_i = cons$, the π_1 identifier is added to $locking-prop_i$ subfield of operation π_2 . Observe that no other action modifies the subfield $\pi_2.locking-prop_i$. Next, notice that the initiator of a consultation phase needs to get back the message it initially sent, before the phase can terminate. By assumption $i \neq \pi_2.intr$. Thus i must send a message before $\pi_2.intr$ receives it. No snd_i event can, however, occur until $\pi_2.locking-prop_i$ is not empty. The $\pi_2.locking-prop_i$ subfield is emptied only if i receives a second message corresponding of every pending propagations whose identifiers belong to $\pi_2.locking-prop_i$. Because of the torus topology we use, each consultation quorum as well as each propagation quorum is a ring. When one of the replicas of a quorum starts a propagation of π , it contacts its two neighbors in its propagation quorum. The two messages are sent in both senses over the ring from a single node i . Because of the uniqueness of the path, each node receives one message locking it before one of them receives a second one, unlocking it. Therefore, each of those have already set their value and tag pair to the one received when $\pi_2.locking-prop_i$ is emptied. Since the snd_i is enabled for consultation π_2 after this time, $resp(\pi_2)$ occurs when all participants of π_1 have stored the propagated tag and val pair. \square

The following lemma and corollary show that tag ordering respects real-time precedence.

Lemma 5.3 *If $resp(\pi_1) \prec inv(\pi_2)$ then $tag(\pi_1) \leq tag(\pi_2)$.*

Proof. In the first case assume that π_1 is a write operation, hence π_1 completes only if its propagation phase has already ended. If $resp(\pi_1) \prec inv(\pi_2)$, it means that π_1 has propagated its tag to a whole propagation quorum, when operation π_2 is initiated. By the quorum intersection property, it exists one element j that gets assigned this tag in every consulting quorum by time $resp(\pi_1)$. Since tag value is monotonically incremented, and the assumptions ensure that any failure is followed by a state recovery, it is clear that tag tag' , consulted during the consultation phase π_2 , is such that $tag' \geq tag(\pi_1)$. By definition of $tag(\pi_2)$, it is such that $tag(\pi_2) \geq tag'$ and putting these inequalities together yields to the result.

Now consider the second case where π_1 is considered to be a read operation. To prove that the property holds, we show that the following contraposition is true: If $tag(\pi_2) < tag(\pi_1)$ then $start(\pi_2) \prec term(\pi_1)$. Assume that π_1 contacts the consulting quorum c_1 while π_2 contacts the consulting quorum c_2 . Hence if $tag(\pi_2) < tag(\pi_1)$ then $\exists j \in c_1$ such that $j.tag > \max_{i \in c_2} \{i.tag\}$. Given that, we show that $resp(\pi_1) \prec inv(\pi_2)$ is impossible. The existence of j implies that it exists a write operation π_3 , propagating to a propagation quorum p_3 , that has propagated to j but not yet to any replica of c_2 . By the quorum intersection property, we know that $p_3 \cap c_2 \neq \emptyset$, then π_3 will eventually propagate to one element of c_2 . By lemma 5.2, j is locked until after having propagated to an element of c_2 and $resp(\pi_1) \prec inv(\pi_2)$ is impossible: if $tag(\pi_2) < tag(\pi_1)$ then $inv(\pi_2) \prec resp(\pi_1)$ and the result is also true, i.e., if $resp(\pi_1) \prec inv(\pi_2)$ then $tag(\pi_1) \leq tag(\pi_2)$. \square

Corollary 5.4 *If $resp(\pi_1) \prec inv(\pi_2)$ and π_2 is a write traversal then $tag(\pi_1) < tag(\pi_2)$.*

Proof. This follows directly from Lemma 5.3 and the definition of write operation tag. Let tag_c be the tag at the end of consultation phase of π_2 before being incremented. By examination of $rcv(\pi_2.id, phase, *, *, *, *, *)$ action that ends the consultation or update phase, we conclude that $tag(\pi_2) > tag_c$. Combining this with Lemma 5.3 leads to the conclusion. \square

The main theorem shows that the traversal ordering satisfies the atomicity definition.

Theorem 5.5 *The Traversal automaton implements atomic object.*

Proof. Two write operations get assigned different tags. This follows from the fact that two writes at the same location get assigned a different sequence number, and writes occurring at different location get assigned different tie-breaker tag. That is Part 2 is satisfied. For Part 1, assume for the sake of contradiction that $<_H \not\subseteq <_t$. That is assume that $\pi_1 <_H \pi_2$ and $\neg(\pi_1 <_t \pi_2)$. Now there are two cases: (i) If π_2 is a read operation, then $\pi_1 <_H \pi_2$ and Lemma 5.3 implies that $tag(\pi_1) \leq tag(\pi_2)$. (ii) If π_2 is a write operation, then $\pi_1 <_H \pi_2$ and Corollary 5.4 implies that $tag(\pi_1) < tag(\pi_2)$. By definition of $<_t$, both results yield a contradiction. For Part 4, since any operation π_1 in H terminates, Lemma 5.3 implies that all operations π_2 such that $\pi_1 <_H \pi_2$ is ordered after. That is, the set of operations preceding π_1 is finite. Part 3 is straightforward. \square

Corollary 5.6 *SAM implements atomic object.*

Proof. Let $inv_{LB}(\pi)$ be the time at which the *LoadBalancer* automaton starts operation π , and let $term_{LB}(\pi)$ be the time at which it ends the same operation. On the first hand, $inv_{LB}(\pi) \prec inv(\pi)$, since the traversal start is ordered by the *LoadBalancer*. On the other hand, $term(\pi) \prec term_{LB}(\pi)$ since the traversal terminates before the termination operation is acknowledged. Therefore we can map any operation of the *LoadBalancer* to the corresponding one of the *Traversal* while choosing the same serialization point for each of them (i.e., keeping the same order among operations). By assumption, tag and value are simply updated by the Adjuster, thus non-violating atomicity. \square

Self-Adjustment. We show that SAM guarantees self-adjusting property.

Theorem 5.7 *SAM self-adjusts.*

Proof. Let τ_{thw} be the duration time of thwart thw . When a replica i is overloaded, ($\mathcal{L}_i(t) > T_{max}^i(t)$) the *LoadBalancer* _{i} orders the start of thwart process thw . By *thwarting-access-rate* assumption, the access rate of i remains lower than its request treatment, thus $\mathcal{L}_i(t)$ does not increase. When the thwart terminates there are two cases: (i) if a replica j decides to treat the requests then the amount of load sent through the thwart is distributed by i to j . Since this amount is decided by i to be the set of requests over its threshold $T_{max}^i(t)$, i is no more overloaded. (ii) if no replica accept to treat the requests, then i receives back the thwart message while its $\mathcal{L}_i(t)$ has not increased since the beginning of the request. Then *potential-replication* ensures that an expand is possible and the load overhead of i is given to the new replica. Finally assume that replica i is underloaded. Then it decides to make a shrink. By leaving the memory, i is no more considered as a replica and its low load is no more taken into account. \square

6 Conclusions and Future Work

In this paper we have proved the correctness of SAM, an architecture emulating an atomic memory in highly dynamic systems exhibiting self-healing and self-adjusting properties depending on object load variations and dynamicity of the system. Properties of SAM (atomicity, self-healing and self-adjusting) are implemented using only local information whose size remains low while the system size grows. SAM provides an efficient solution to the persistent storage problem defined in [4]. It is a fundamental abstraction suitable as a building block for many distributed complex applications, especially in dynamic systems.

We are currently working on experimental results of SAM based on p2p simulator to emphasize its behavior in face of dynamism and high-scale. Future work includes latency comparison of other alternative primitives for read operations. Finally, we intend to borrow some interesting results from the mechanism theory to find incentive strategy for data replication among peers.

References

- [1] I. Abraham, B. Awerbuch, Y. Azar, Y. Bartal, D. Malkhi, and E. Pavlov. A generic scheme for building overlay networks in adversarial scenarios. In *Proc. of the 17th Int. Parallel and Distributed Processing Symposium (IPDPS)*, page 40, 2003.
- [2] I. Abraham and D. Malkhi. Probabilistic quorums for dynamic systems. In Faith Ellen Fich, editor, *Distributed algorithms*, volume 2848/2003 of *Lecture Notes in Computer Science*, pages 60–74, 2003.
- [3] D. Agrawal and A. El Abbadi. Efficient solution to the distributed mutual exclusion problem. In *Proc. of the 8th annual symposium on Principles of distributed computing (PODC)*, pages 193–200, 1989.
- [4] E. Anceaume, R. Friedman, M. Gradinariu, and M. Roy. An architecture for dynamic scalable self-managed persistent objects. In *Proc. of Int. Symposium on Distributed Objects and Applications (DOA)*, pages 1445–1462, 2004.
- [5] E. Anceaume, M. Gradinariu, V. Gramoli, and A. Virgillito. P2p architecture for self* atomic memory. In *Proc. of the 8th Int. Symposium on Parallel Architectures, Algorithms, and Networks (ISPAN)*, 2005. To appear.
- [6] E. Anceaume, M. Gradinariu, V. Gramoli, and A. Virgillito. Sam: Self* atomic memory for p2p systems. Technical Report 1717, IRISA, 2005.
- [7] H. Attiya, A. Bar-Noy, and D. Dolev. Sharing memory robustly in message-passing systems. *Journal of the ACM*, 42(1):124–142, 1995.
- [8] S. Y. Cheung, M. H. Ammar, and M. Ahamad. The grid protocol: A high performance scheme for maintaining replicated data. *IEEE Trans. Knowl. Data Eng.*, 4(6):582–592, 1992.
- [9] S. Dolev, S. Gilbert, N. Lynch, A. Shvartsman, and J. Welch. GeoQuorums: Implementing atomic memory in ad hoc networks. In *Proc. of 17th International Symposium on Distributed Computing (DISC)*, pages 306–320, 2003.
- [10] B. Englert and A. Shvartsman. Graceful quorum reconfiguration in a robust emulation of of shared memory. In *Proc. of Int. Conf. on Distributed Computing Systems (ICDCS 2000)*, pages 454–463, 2000.
- [11] D. K. Gifford. Weighted voting for replicated data. In *Proc. of the 7th ACM symposium on Operating systems principles (SOSP’79)*, pages 150–162. ACM Press, 1979.
- [12] S. Gilbert, N.A. Lynch, and A. Shvartsman. RAMBO II: Rapidly reconfigurable atomic memory for dynamic networks. In *Proc. of 17th Int. Conference on Dependable Systems and Networks (DSN)*, pages 259–269, 2003.
- [13] C. Gray and D. Cheriton. Leases: an efficient fault-tolerant mechanism for distributed file cache consistency. In *Proc. of the 12th ACM symposium on Operating systems principles (SOSP’89)*, pages 202–210. ACM Press, 1989.
- [14] M.P. Herlihy. Dynamic quorum adjustment for partitioned data. *ACM Trans. on Database Systems*, 12(2):170–194, 1987.
- [15] R. Holzman, Y. Marcus, and D. Peleg. Load balancing in quorum systems. *SIAM Journal on Discrete Mathematics*, 10(2):223–245, 1997.
- [16] A. Kumar. Hierarchical quorum consensus: A new algorithm for managing replicated data. *IEEE Trans. Computers*, 40(9):996–1004, 1991.

- [17] N. Lynch and A. Shvartsman. RAMBO: A reconfigurable atomic memory service for dynamic networks. In *Proc. of 16th International Symposium on Distributed Computing (DISC)*, pages 173–190, 2002.
- [18] N. Lynch and M. Tuttle. An introduction to input/output automata. Technical Report 3, CWI-Quaxterly, 1989.
- [19] N.A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, 1996.
- [20] N.A. Lynch and A. Shvartsman. Robust emulation of shared memory using dynamic quorums. In *Proc. of 27th Int. Symp. on Fault-Tolerant Computing (FTCS)*, pages 272–281, 1997.
- [21] Nancy A. Lynch, Dahlia Malkhi, and David Ratajczak. Atomic data access in distributed hash tables. In *IPTPS*, pages 295–305, 2002.
- [22] M. Maekawa. A \sqrt{N} algorithm for mutual exclusion in decentralized systems. *ACM Trans. on Computer Systems (TOCS)*, 3(2):145–159, 1985.
- [23] U. Nadav and M. Naor. Fault-tolerant storage in a dynamic environment. In *Proc. of the 18th Annual Conference on Distributed Computing (DISC)*, 2004.
- [24] M. Naor and U. Wieder. Novel architectures for p2p applications: the continuous-discrete approach. In *15th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 50–59, 2003.
- [25] M. Naor and A. Wool. The load, capacity, and availability of quorum systems. *SIAM Journal on Computing*, 27(2):423–447, 1998.
- [26] Moni Naor and Udi Wieder. Scalable and dynamic quorum systems. In *Proc. of the 22th annual symposium on Principles of distributed computing (PODC'03)*, pages 114–122. ACM Press, 2003.
- [27] S. Ratnasamy, P. Francis, M. Handley, R. M. Karp, and S. Shenker. A scalable content-addressable network. In *Proc. of the ACM SIGCOMM*, pages 161–172, 2001.
- [28] B. Silaghi, P. Keleher, and B. Bhattacharjee. Multi-dimensional quorum sets for read-few write-many replica control protocols. In *In Proc. of the 4th CCGRID/GP2PC*, 2004.

Appendix

A The Load Balancer and Adjuster Specification

I/O Automaton 2 *LoadBalancer_i*

1: **Signature:**

Input:

- 2: $rcv((op, type, v))_{j,i}$, $i \in I$, $type \in \{\text{read, write}\}$, $v \in V$
 3: $read\text{-}write\text{-}ack(id, v)_i$, $i \in I$, $v \in V$
 4: $rcv(m)_{j,i}$, $i, j \in I$, $m \in M$

Internal:

- 5: $thwart_i$, $i \in I$

12: **State:**

- 13: $rqst$ a record with fields
 14: $tid \in TId$
 15: $type \in \{\text{read, write}\}$
 16: $rqstr \in I$
 17: $intr \in I$
 18: $val \in V$
 19: $tag \in \mathbb{N} \times I \times \mathbb{N}$
 20: $phase \in \{\text{idle, acknowledging}\}$
 21: $batch$, an ordered set of requests $rqst$
 22:
 23: $val \in V$
 24: $failed$ a boolean
 25: $replica$ a boolean

Output:

- 6: $snd((op, type, v))_{i,j}$, $i \in I$, $type \in \{\text{read, write}\}$, $v \in V$
 7: $read\text{-}write(type, id, v)_i$, $i \in I$, $type \in \{\text{read, write}\}$,
 8: $id \in \Pi$, $v \in V$
 9: $shrink_i$, $i \in I$
 10: $expand(ra, wa)_i$, $i \in I$, $ra, wa \in \mathbb{N}$
 11: $snd(m)_{i,j}$, $i, j \in I$, $m \in M$

- 26: $heap$ a set of requests $rqst$
 27: $treating$ a set of requests $rqst$
 28: $clock \in \mathbb{R}^{>0}$
 29: $treat\text{-}time \in \mathbb{R}^{>0}$
 30: $treat\text{-}period \in \mathbb{R}^{>0}$
 31: $shrink\text{-}time \in \mathbb{R}^{>0}$
 32: $shrink\text{-}period \in \mathbb{R}^{>0}$
 33: $access\text{-}rate \in \mathbb{R}^{>0}$ the rate of local request reception.
 34: $treatment\text{-}rate \in \mathbb{R}^{>0}$ the rate of local request treatment.
 35: $fwd \in TId$
 36: $starter \in I$
 37: $order\text{-}expand$ a boolean
 38: $r\text{-}access \in \mathbb{N}$
 39: $w\text{-}access \in \mathbb{N}$
-

```

40: Transitions:
41: Input rcv( $\langle op, type, v \rangle$ )j,i
42: Effect:
43: if  $\neg failed \wedge replica$  then
44:   if  $type = read$  then
45:      $r-access \leftarrow r-access + 1$ 
46:   else
47:      $w-access \leftarrow w-access + 1$ 
48:      $ct \leftarrow w-access + r-access$ 
49:      $rqst \leftarrow \langle \langle ct, i \rangle, type, j, i, v, \perp, idle, \emptyset \rangle$ 
50:      $heap \leftarrow heap \cup \{rqst\}$ 
51:     if  $access-rate \geq treatment-rate$  then
52:        $fwd \leftarrow fwd \cup top(heap)$ 
53:        $heap \leftarrow bottom(heap)$ 
54:        $starter \leftarrow i$ 
55:        $shrink-time \leftarrow clock + shrink-period$ 

56: Output read-write( $type, id, v$ )i
57: Precondition:
58:    $\neg failed \wedge replica$ 
59:    $heap \neq \emptyset$ 
60:    $clock > treat-time$ 
61:    $W = \{w \in heap : w.type = write\}$ 
62:    $R = \{r \in heap : r.type = read\}$ 
63:   if  $W \neq \emptyset$  then
64:      $id = rqst.id = \max\{rqst'.id : rqst' \in W\}$ 
65:      $val = rqst.val$ 
66:      $type = write$ 
67:   else
68:      $id = rqst.id = \max\{rqst'.id : rqst' \in R\}$ 
69:      $val = rqst.val$ 
70:      $type = read$ 
71:      $rqst.batch \leftarrow heap \setminus \{rqst\}$ 
72:      $treat-time \leftarrow clock + treat-period$ 
73:      $treating \leftarrow treating \cup \{rqst\}$ 
74: Effect:
75:   none

76: Input read-write-ack( $id, v$ )
77: Effect:
78:   if  $\neg failed \wedge replica$  then
79:      $rqst \in treating$  st.  $rqst.id = id$ 
80:      $rqst.phase \leftarrow acknowledging$ 

81: Output snd( $\langle op-ack, v \rangle$ )i
82: Precondition:
83:    $\neg failed \wedge replica$ 
84:    $rqst \in treating$ 
85:    $rqst.phase = acknowledging$ 
86:    $rqst' \in rqst.batch \vee j = rqst.rqstr$ 
87:    $rqst.batch = \emptyset \vee j = rqst'.rqstr$ 
88:    $v = rqst.val$ 
89: Effect:
90:   if  $rqst.batch = \emptyset$  then
91:      $rqst.phase \leftarrow ending$ 
92:      $treating \leftarrow treating \setminus \{rqst\}$ 
93:   else
94:      $rqst.batch \leftarrow rqst.batch \setminus \{rqst'\}$ 

95: Output shrinki
96: Precondition:
97:    $\neg failed \wedge replica$ 
98:    $clock \geq shrink-time$ 
99:    $heap = \emptyset$ 
100: Effect:
101:   if  $r-access = w-access = 0$  then
102:      $replica \leftarrow false$ 
103:      $shrink-time \leftarrow \infty$ 
104:   else
105:      $r-access \leftarrow 0$ 
106:      $w-access \leftarrow 0$ 

107: Output expand( $ra, wa$ )i
108: Precondition:
109:    $\neg failed \wedge replica$ 
110:    $order-expand = true$ 
111:    $ra = r-access$ 
112:    $wa = w-access$ 
113: Effect:
114:    $order-expand \leftarrow false$ 

115: Internal thwarti
116: Precondition:
117:    $\neg failed \wedge replica$ 
118:    $access-rate \geq treatment-rate$ 
119: Effect:
120:    $fwd \leftarrow top(heap)$ 
121:    $heap \leftarrow bottom(heap)$ 
122:    $starter \leftarrow i$ 

123: Output snd( $\langle fwd, str \rangle$ )i,j
124: Precondition:
125:    $\neg failed \wedge replica \wedge freezing = \emptyset$ 
126:    $fwd \neq \emptyset$ 
127:    $str = starter$ 
128:    $j = get-thwart-nbr$ 
129: Effect:
130:    $fwd \leftarrow \emptyset$ 

131: Input rcv( $\langle fwd, str \rangle$ )j,i
132: Effect:
133:   if  $\neg failed \wedge replica$  then
134:     if  $str = i$  then
135:        $order-expand \leftarrow true$ 
136:     else
137:        $heap \leftarrow heap \cup \{fwd\}$ 
138:        $starter \leftarrow str$ 

```

I/O Automaton 3 *Adjuster_i* — Signature and state**1: Signature:****Input:**

- 2: shrink_i, $i \in I$
 3: expand_i, $i \in I$
 4: rcv(m)_{j,i}, $i, j \in I, m \in M$

7: State:

- 8: zones a set of records with fields
 9: $x_1 \in \mathbb{R}$
 10: $x_2 \in \mathbb{R}$
 11: $y_1 \in \mathbb{R}$
 12: $y_2 \in \mathbb{R}$
 13: nbr a record with fields
 14: zone $\in \mathbb{R}^4$
 15: $id \in I$
 16: $nid \subseteq I$
 17: $pp \in TId \times \{N, S, E\}$
 18: m a record with fields
 19: $tid \in TId$
 20: $type \in \{\text{read}, \text{write}\}$
 21: $intr \in I$
 22: $str \in I$
 23: $dir \in \{N, S, E\}$
 24: $line \in \mathbb{R} \times \mathbb{R}$
 25: tag a record with fields
 26: $ct \in \mathbb{N}$

Internal:

- 5: heal_i, $i \in I$
Output:
 6: snd(m)_{i,j}, $i, j \in I, m \in M$
 27: $id \in I$
 28: $index \in I$
 29: $val \in V$
 30: bounds a zone
 31: nbrs an array of neighbor nbr
 32: failed a boolean
 33: has-changed a boolean
 34: involving-msg a set of messages
 35: participating-msg a set of messages
 36: replica a boolean
 37: $clock \in \mathbb{R}^{>0}$
 38: $hb\text{-time} \in \mathbb{R}^{>0}$
 39: $hb\text{-period} \in \mathbb{R}^{>0}$
 40: last-split $\in I$
 41: a-last-split a mapping from I to a boolean
 42: detect-time a mapping from I to \mathbb{R}
 43: detect-period $\in \mathbb{R}^{>0}$
 44: freezing $\subseteq \{N, S, E\}$
 45: $coeff \in \mathbb{R}^{>0}$

1: Transitions:**2: Input** shrink_i**3: Effect:**

- 4: if $\neg \text{failed} \wedge \text{replica}$ then
 5: last-split = \perp
 6: replica \leftarrow false

7: Input expand(ra, wa)_i**8: Effect:**

- 9: if $\neg \text{failed} \wedge \text{replica}$ then
 10: to-add \leftarrow get-outside-node()
 11: replica \leftarrow true
 12: if $((\text{coeff} \times ra) > wa)$ then
 13: (zones.to-add.zones) \leftarrow get-zone(hor)
 14: freezing \leftarrow freezing $\cup \{N\}$
 15: else
 16: (zones.to-add.zones) \leftarrow get-zone(vert)
 17: freezing \leftarrow freezing $\cup \{E\}$
 18: $index \leftarrow |nbrs| + 1$
 19: $nbrs[index] \leftarrow$ to-add
 20: a-last-split[to-add] \leftarrow true

21: Output snd($\langle \text{heartbeat}, pp, pl, z, nid, ls, t, v \rangle$)_{i,j}**22: Precondition:**

- 23: $\neg \text{failed} \wedge \text{replica}$
 24: $j \in nbrs$
 25: $hb\text{-time} \geq \text{clock}$
 26: $pp = \text{pending-props}$
 27: $pl = \text{propagate-line}$
 28: $z \leftarrow \text{zones}$
 29: $nid \leftarrow \bigcup_{n \in nbrs} \{n.id\}$
 30: $ls \leftarrow \text{last-split}$
 31: $t \leftarrow \text{tag}$
 32: $v \leftarrow \text{val}$
 33: **Effect:**
 34: $hb\text{-time} = \text{clock} + hb\text{-period}$

35: Input rcv($\langle \text{heartbeat}, pp, pl, z, nid, ls, t, v \rangle$)_{j,i}**36: Effect:**

- 37: if $\neg \text{failed} \wedge \text{replica}$ then
 38: if $\forall index : nbrs[index].id \neq j$ then
 39: $index \leftarrow |nbrs| + 1$
 40: $nbrs[index].id \leftarrow j$

41: else

- 42: let $index$ be st. $nbrs[index].id = j$
 43: $nbrs[index].zones \leftarrow z$
 44: $nbrs[index].nbrs \leftarrow nid$
 45: $nbrs[index].pp \leftarrow pp$
 46: $nbrs[index].pl \leftarrow pl$
 47: if $ls \neq i$ then
 48: a-last-split[j] \leftarrow false
 49: if $\forall z' \in j'.zones, \bigcup_{j' \in nbrs} \{[z'.y_1, z'.y_2] : z'.x_1 > zone.x_1\} \subseteq [zone.y_1, zone.y_2]$ then
 50: freezing \leftarrow freezing $\setminus \{E\}$
 51: if $\forall z' \in j'.zones, \bigcup_{j' \in nbrs} \{[z'.x_1, z'.x_2] : z'.y_1 > zone.y_1\} \subseteq [zone.x_1, zone.x_2]$ then
 52: freezing \leftarrow freezing $\setminus \{N\}$
 53: if $\forall z' \in j'.zones, \bigcup_{j' \in nbrs} \{[z'.x_1, z'.x_2] : z'.y_1 < zone.y_1\} \subseteq [zone.x_1, zone.x_2]$ then
 54: freezing \leftarrow freezing $\setminus \{S\}$
 55: if $t > \text{tag}$ then
 56: $\langle \text{tag}, \text{val} \rangle \leftarrow \langle t, v \rangle$
 57: if freezing = \emptyset then
 58: pending-props \leftarrow update-pp($nbrs, pl$)

59: Output snd($\langle \text{expand}, v, t, neighbors \rangle$)_{i,j}**60: Precondition:**

- 61: $\neg \text{failed} \wedge \text{replica}$
 62: $j = \text{to-add}$
 63: $t = \text{tag}$
 64: $v = \text{val}$
 65: $neighbors = nbrs$
 66: **Effect:**
 67: none

68: Input rcv($\langle \text{expand}, v, t, neighbors \rangle$)_{j,i}**69: Effect:**

- 70: if $\neg \text{failed} \wedge \neg \text{replica}$ then
 71: $val \leftarrow v$
 72: $tag \leftarrow t$
 73: $nbrs \leftarrow \text{update-nbr}(\text{zone}, \text{neighbors})$
 74: last-split $\leftarrow j$
 75: replica \leftarrow true
 76: freezing $\leftarrow \{N, S, E\}$