

A Highly Available Log Service for Distributed Transaction Processing

Lásaro Camargos^{‡*}

Marcin Wieloch^{*}

Fernando Pedone^{*}

Edmundo Madeira[‡]

^{*}Faculty of Informatics
University of Lugano
6900 Lugano, Switzerland

[‡]Institute of Computing
University of Campinas
3084-971 Campinas-SP, Brazil

University of Lugano
Faculty of Informatics
Technical Report No. 2006/08
December 2006

Abstract

Processing distributed transactions hinges on enforcing atomicity and durability: resource managers participating in a transaction must agree on its outcome and the transaction updates must be permanent. We introduce the log service abstraction, which gathers resource managers' votes to commit or abort a transaction and their updates, and outputs the transaction's outcome. Updates are made durable and non-concurrent transactions are totally ordered by the service. The sequence of updates performed by a resource manager is available as a means to consistently recover resource managers without relying on their local state. As a consequence, a remote process, whose state will be recovered from the log service, can reincarnate a crashed resource manager. Moreover, the service ensures that only one process plays the role of a given resource manager at any time. We present two highly available implementations of this service and evaluate their performance running TPC-C and a microbenchmark on a distributed database.

Contents

1	Introduction	1
1.1	Atomic Commit and Paxos Commit	1
1.2	The Highly Available Log Service	1
1.3	From specs to implementations	2
1.4	Theory and practice of termination	2
1.5	Recovering resource managers	2
1.6	Our contributions	3
2	Problem statement	3
3	The Log Service	4
3.1	Terminology and notation	4
3.2	The Log Service specification	5
3.3	Termination and Recovery	5
3.4	Correctness	7
4	From abstract specifications to distributed implementations	7
4.1	Processes, communications and failures	7
4.2	Leader-election oracle	8
4.3	Consensus	8
5	Coordinated Implementation	8
5.1	RM Stubs	9
5.2	Log Service	10
6	Uncoordinated Implementation	10
6.1	Transaction Termination	12
6.2	Recovering from Failures	12
7	Evaluation	12
7.1	Micro-benchmark	14
7.2	The TPC-C benchmark	14
8	Related Work	15
9	Conclusion	16
A	Log Service	17
A.1	Log Service Constants	17
A.2	Log Service Specification	17
A.3	Correctness	23
B	Coordinated Implementation	27
B.1	Specification	27
B.2	Implementation Proof	36

C	Uncoordinated Implementation	40
C.1	Specification	40
C.2	Implementation Proof	49

1 Introduction

1.1 Atomic Commit and Paxos Commit

The problem of atomic termination of transactions is central to any distributed database system aiming at strong consistency. Terminating a distributed transaction is normally performed by an *atomic commitment protocol*, which gathers votes from all resource managers involved in the transaction and decides on the transaction's outcome based on these votes. If all votes are for commit, the transaction is committed; otherwise it is aborted. To avoid blocking when a resource manager is not available to vote, a weaker version of atomic commitment requires the outcome to be commit if none of the participants is suspected to have failed and all vote to commit the transaction. This weaker problem is called *non-blocking atomic commitment*.

In a recent paper, Gray and Lamport have introduced Paxos Commit, a non-blocking atomic commitment protocol based on the Paxos consensus algorithm [5]. In Paxos Commit each resource manager uses a consensus instance to vote. If a resource manager suspects that another one has crashed and cannot vote, the first one votes on behalf of the second resource manager, proposing in its consensus instance to abort the transaction. If more than one vote is issued for a resource manager (e.g., it is wrongly suspected), consensus ensures that participants will agree on the vote of each participant.¹ Paxos Commit is not only an elegant but also an efficient solution to non-blocking atomic commitment. It has the same latency as two-phase commit (2PC) without blocking if the transaction coordinator fails. It is cheaper than three-phase commit (3PC) and has a simpler termination protocol in case of failure of the coordinator [2]. Instead of keeping the participant votes at the coordinator, in Paxos Commit this information is stored at the *acceptors*, the processes responsible for consensus decisions.

1.2 The Highly Available Log Service

Paxos Commit suggests that information critical to transaction termination should reside neither at the coordinator, as in 2PC, nor at the resource managers, as in 3PC, but at a third party, the acceptors, whose availability can be easily parameterized. Inspired by Paxos Commit, we abstract transaction termination in terms of a *highly available log service*. This approach is justified by the following reasons:

- *Simplicity*. A non-blocking atomic commitment protocol based on a highly available log service is quite intuitive: Transaction participants submit their votes to the remote log; as soon as there is one commit vote for each participant or one abort vote stored in the log, the outcome of the transaction can be determined. Even if more than one vote is issued for a resource manager, should it be suspected to have failed, the ordering in which the votes are kept in the log ensures that the transaction outcome is deterministic.
- *Performance*. By defining transaction termination in terms of a log service, we allow other implementations to take advantage of our approach to non-blocking atomic commitment. In fact, we present in the paper an implementation of the log service that has better performance than a Paxos Commit-like implementation. We substantiate this claim with experimental results using the TPC-C benchmark and a micro-benchmark running in a cluster of nodes.
- *Fault-tolerance*. The log can be made highly available using standard replication techniques (e.g., state machine replication). Moreover, if a resource manager fails, then another instance of it can

¹Notice that in Paxos not all consensus participants are required to propose a value.

be created on an operational node using the state stored in the log. Thus, the recovery of the node hosting the resource manager is not needed for the system to resume operation, increasing the availability of resource managers.

1.3 From specs to implementations

We discuss in detail two implementations of the log service, *uncoordinated*, inspired by Paxos Commit, and *coordinated*. Both implementations rely on a series of consensus executions. The main difference between them concerns the ordering of these executions, and therefore the ordering of committed transactions in the log. While there is a separate sequence of consensus instances per resource manager in the uncoordinated implementation, leading to a *partial order of committed transactions in the log*, there is a single sequence of consensus in the coordinated version, ensuring a *total order of committed transactions in the log*. In the coordinated implementation, instead of using consensus to vote, resource managers send their votes to the coordinator, which submits the votes in a consensus instance. The nature of each approach has performance implications on both the termination of transactions (Section 3.3) and the recovery of resource managers (Section 3.3).

1.4 Theory and practice of termination

The coordinated implementation requires one more communication step than the uncoordinated implementation, corresponding to the vote sent by the participants to the coordinator. Notice that the coordinator here does not play the same role as the transaction coordinator in 2PC or 3PC. The role of the transaction coordinator in the coordinated log implementation is to gather all votes and propose them in a consensus instance. Just like in Paxos, if the coordinator fails, it can be promptly replaced by another process.

In theory, the coordinated log implementation should take longer to execute than the uncoordinated log implementation. However, performance experiments revealed that, contrary to the expectations, the coordinated approach largely outperforms the uncoordinated technique in a clustered environment. This is due to the fact that by having a single coordinator, votes from all participants in an atomic commitment, and also from concurrent atomic commitment executions, can be batched by the coordinator and proposed in a single consensus instance, saving not only in the number of messages but also in the processing done by the acceptors. These savings largely make up for the additional message latency of the coordinated approach.

1.5 Recovering resource managers

The ordering of transactions in the log has an impact on the recovery of resource managers. Upon recovery, a resource manager should find out the outcome of transactions in whose termination it was involved before the failure. The question is: How many consensus instances should be inspected until a resource manager can rebuild a consistent image of the local database? The answer to this question can be simple if the log implementation is coordinated: It suffices for the resource manager to send a “flush” request to the coordinator and wait until this request is decided in some instance. Once the request is decided in instance i , it follows from the total order of log entries that any pending transactions concerning the resource manager must have been decided before instance i . The coordinator may have them cached locally, or learn their values by re-executing the consensus instances for the missing values.

Recovery with an uncoordinated log implementation is more complex and depends on the recovery strategy and the maximum number k of transactions allowed to execute concurrently at the resource manager:

- If resource managers are statically assigned to nodes, and a node crashes, the resource manager will only recover after the node is operational again. In this case, the resource manager proposes “flush” values using its local sequence of consensus until the value is decided in some instance i . Once this happens, it should find out the decision of instances $i + 1, i + 2, \dots, i + k - 1$ since there could have been up to k concurrent transactions in execution.
- If a resource manager is dynamically assigned to a node and this one is suspected to have crashed, another instance of it is reincarnated on an available node. In the best case, the procedure described above may work here too, but if the original node did not crash, there could be two operational processes incarnating the same resource manager simultaneously, and they would propose values in the same consensus stream, an undesired situation. Eventually the “flush” request will be delivered, and then $k - 1$ more instances should be executed.

1.6 Our contributions

Section 2 reviews the atomic commit problem and specifies a property for consistent recovery. Section 3 introduces the log service specification and shows how it can be used by resource managers. Some implementation considerations are discussed in Section 4, and the coordinated and uncoordinated log implementations are presented in Sections 5 and 6, respectively. These approaches are compared in Section 7. Section 8 reviews related work and Section 9 concludes the paper.

2 Problem statement

Transactions are executed by a collection of processes called *resource managers (RMs)*. When a transaction ends, a special process called *transaction manager (TM)* starts an atomic commit (AC) protocol, in which every RM participating in the transaction votes for its commit or abort. Processes may fail by crashing and never behave maliciously (i.e., no Byzantine failures).

Atomic commit is defined by the following properties.

- **AC-Validity** If an RM decides to commit a transaction, then all RMs voted to commit the transaction.
- **AC-Agreement** It is impossible for one RM to commit a transaction and another one to abort the transaction.
- **AC-Non-Triviality** If all RMs vote to commit the transaction and no RM is suspected throughout the execution of the protocol, then the decision is commit.
- **AC-Termination** Non-faulty RMs eventually decide.

AC-Non-Triviality assumes that RMs can be “suspected” to have failed. The only assumption that we make about failure detection is that if an RM fails, then it will eventually be suspected by the other

processes.² Therefore, resource managers may be incorrectly suspected to have failed, in which case transactions would be unnecessarily aborted.

Atomic commit defines how distributed transactions terminate. A related problem is building a consistent state of a recovering resource manager. Intuitively, the problems are related because after recovering from a crash, an RM should find out which transactions were committed before the failure. Determining such transactions requires identifying previous instances of atomic commit executions.

The following property characterizes correct recovery from failures. R-Consistency can be ensured, for example, by replaying the write operations of committed transactions in the order in which these transactions committed.

- **R-Consistency** The database state of an RM after the recovery from a failure is the same as its committed state before the failure.

3 The Log Service

In the following we discuss how atomic commitment and recovery of resource managers can be implemented using an abstract log service. Before explaining the abstract algorithms, some formalism must be introduced.

3.1 Terminology and notation

The symbol “ \triangleq ” should be read as “is defined as”. For example,

$$\text{IF}(a, b, c) \triangleq \text{if } a \text{ then } b \text{ else } c$$

means that $\text{IF}(a, b, c)$ evaluates to b , if a is true, and it evaluates to c , otherwise; and

$$\text{ADDVALTOB}(val) \triangleq \mathcal{B} \leftarrow \mathcal{B} + val$$

increments \mathcal{B} by the value val . **let** and **in** specify the scope of a definition. For example,

$$\text{let SUM}(a, b) \triangleq a + b \text{ in SUM}(\text{SUM}(1, 2), \text{SUM}(3, 4))$$

defines $\text{SUM}(a, b)$ in the expression on the right of **in** only, evaluated to 10. For simplicity we assume that, with one exception (Algorithm 1 in Section 3.3), all definitions presented in the paper are performed atomically.

We use angle brackets to enclose a sequence of values as a tuple; $\langle \rangle$ is the tuple with no elements, and $\langle a, b, c \rangle$ is the tuple with elements a, b and c , in this order. The symbol “ $-$ ” in a tuple matches any element. That is, $\langle a, b \rangle = \langle -, b \rangle$ and $\langle -, b \rangle = \langle c, b \rangle$ even if $c \neq a$. We denote by $\text{Len}(s)$ the number of elements of sequence s , and use the square brackets to index its elements, e.g., $\langle a, b, c \rangle[2]$ is the element b . $a <_s b$ (resp. $>_s b$) holds if and only if a and b happen only once in s and $s[i] = a$ and $s[j] = b$ implies $i < j$ (resp. $i > j$). $s \bullet e$ is defined as the sequence s appended by element e (e.g., $\langle a, b \rangle \bullet c = \langle a, b, c \rangle$). If ss is a sequence of sets, then $ss \oplus e$ adds element e to the last set in ss , $ss[\text{Len}(ss)]$, (e.g., $\langle \{a, b\}, \{c\} \rangle \oplus d = \langle \{a, b\}, \{c, d\} \rangle$).

²This property is similar to the *completeness* property of Chandra and Toueg’s failure detectors [3].

3.2 The Log Service specification

Algorithm 1 presents the abstract behavior of the log service (lines 1–40) and of transaction termination and resource manager recovery (lines 41–54, in page 6). In Sections 4–6 we show how these behaviors can be implemented in a shared-nothing asynchronous distributed system.

The algorithm uses five data-structures:

\mathcal{V} The set of received votes, initially empty.

\mathcal{T} A partially ordered set (S, \preceq) , where S is the set of committed transactions, and \preceq is partial order relation over the elements of S . Non-concurrent transactions in S are totally ordered by \preceq , according to their commit order. For simplicity, we represent \mathcal{T} as a sequence of sets of committed transactions such that, given two sets $\mathcal{T}[i]$ and $\mathcal{T}[j]$ in \mathcal{T} , $i < j$ implies that for all transactions t and u , $t \in \mathcal{T}[j]$ and $u \in \mathcal{T}[i]$, $t \preceq u$.

\mathcal{C} and $LastC$ Helpers on building \mathcal{T} . \mathcal{C} is the set of non-terminated transactions whose votes have been issued; $LastC$ is the subset of \mathcal{C} with transactions concurrent to the last committed one. Both are initially empty.

\mathcal{R} A mapping from resource managers to the processes that are currently “incarnating” them. Initially, all resource managers are mapped to no process (i.e., \perp is an invalid process identifier).

$OUTCOME(t)$ defines the result of transaction t . It is determined by checking whether at least one vote for t is ABORT, in which case the outcome is ABORT, or if all votes are COMMIT, in which case the outcome is COMMIT. If not all votes are known to the log service, but all the votes known so far are COMMIT, the outcome is UNDEFINED.

$ISINVOLVED(t, rm)$ evaluates to TRUE if rm is a participant in t and to FALSE, otherwise, by checking if rm is in the participant list ($tset$) of any known vote for t . Because RMs may vote to ABORT a transaction at any time during its execution, including before learning about the other participants, ABORT votes may not contain the complete list of participants. Therefore, ISINVOLVED may not have conclusive information if all known votes for t are ABORT; in such a case the answer is UNKNOWN.

VOTE defines how votes are added to \mathcal{V} . A vote is added only if no previous vote for the same resource manager to the same transaction was added before. This ensures that if suspecting and suspected processes issue conflicting votes, only one vote per participant will be kept.

UPDATES(rm) evaluates to the sequence of sets of updates performed by resource manager rm , partially ordered accordingly to \mathcal{T} .

3.3 Termination and Recovery

Resource managers use the log service to terminate transactions and recover after crashes, as defined in Algorithm 1 (lines 41–54, in page 7). INCARNATE is used by process pid to incarnate resource manager rm , be it the first process to do it or a replacement incarnation. Firstly, the process sets $\mathcal{R}[rm]$ to pid and then evaluates UPDATES, described above, to get the updates executed by the previous rm incarnation. $updates$ is then scanned in order from the first to the last set, and all updates in one set are applied to the database before all the ones in the next set; no order among elements in the same set is required.

At the end of INCARNATE, if $\mathcal{R}[rm] = pid$, then pid has successfully incarnated rm and recovered the previous incarnation’s state. pid will accept and process new transactions until it crashes or another process incarnates rm (i.e., $pid \neq \mathcal{R}[rm]$). If more than one process try to incarnate rm , only the last one to execute INCARNATE will succeed.

Algorithm 1 Log service specification

```

1: Initially:
2:  $\mathcal{V} \leftarrow \emptyset$   $\triangleleft$  The history of votes.
3:  $\mathcal{T} \leftarrow \langle \rangle$   $\triangleleft$  Sequence of sets of committed transactions.
4:  $\mathcal{C} \leftarrow \emptyset$   $\triangleleft$  Set of concurrent transactions.
5:  $LastC \leftarrow \emptyset$   $\triangleleft$  Set of concurrent transactions.
6:  $\forall r \in RM, \mathcal{R}[r] \leftarrow \perp$   $\triangleleft$  Processes incarnating RMs.

7:  $OUTCOME(t) \triangleq$ 
8:   if  $\exists \langle -, t, -, ABORT, - \rangle \in \mathcal{V}$   $\triangleleft$  Any ABORTs?
9:     ABORT
10:  else if  $\exists \langle -, t, tset, -, - \rangle \in \mathcal{V} : \forall p \in tset :$ 
11:     $\langle p, t, tset, COMMIT, - \rangle \in \mathcal{V}$   $\triangleleft$  All COMMITs?
12:      COMMIT
13:  else UNDEFINED  $\triangleleft$  Neither one nor the other

14:  $ISINVOLVED(t, rm) \triangleq$ 
15:  if  $\exists \langle -, t, tset, -, - \rangle \in \mathcal{V} : rm \in tset$   $\triangleleft$  Is  $rm$  in any list?
16:    TRUE
17:  else if  $\exists \langle -, t, -, v, - \rangle \in \mathcal{V} : v = COMMIT$   $\triangleleft$  Has  $t$  committed?
18:    FALSE
19:  else
20:    UNKNOWN  $\triangleleft$  I don't know...

21:  $VOTE(\langle rm, t, tset, vote, update \rangle) \triangleq$ 
22:  if  $OUTCOME(t) = UNDEFINED$   $\triangleleft$  If  $t$  has not terminated yet
23:     $\mathcal{C} \leftarrow \mathcal{C} \cup \{t\}$   $\triangleleft$  add it to  $\mathcal{C}$ .
24:  if  $\neg \exists \langle rm, t, -, -, - \rangle \in \mathcal{V}$   $\triangleleft$  Has  $rm$  voted yet?
25:     $prevState \leftarrow OUTCOME(t)$   $\triangleleft$  ABORT or UNDEFINED.
26:     $\mathcal{V} \leftarrow \mathcal{V} \cup \{\langle rm, t, tset, vote, update \rangle\}$   $\triangleleft$  Keep this vote
27:  if  $(prevState = UNDEFINED) \wedge (OUTCOME(t) = COMMIT)$ 
28:    if  $t \in LastC$   $\triangleleft$  If  $t$  can be added to the last set...
29:       $\mathcal{T} \leftarrow \mathcal{T} \oplus t$   $\triangleleft$  ...do it...
30:    else
31:       $\mathcal{T} \leftarrow \mathcal{T} \bullet \{t\}$   $\triangleleft$  ...else add it to a new set
32:       $LastC \leftarrow \mathcal{C}$   $\triangleleft$  With a new  $LastC$ .
33:       $\mathcal{C} \leftarrow \mathcal{C} \setminus \{t\}$ 

34:  $UPDATES(rm) \triangleq$   $\triangleleft$  Set incarnating process.
35:  let  $UPD(i) \triangleq$ 
36:    if  $i = 0$ 
37:       $\langle \rangle$ 
38:    else
39:       $UPD(i - 1) \bullet \{upd : \langle rm, t, -, COMMIT, upd \rangle \in \mathcal{V} :$ 
 $t \in \mathcal{T}[i] \wedge ISINVOLVED(rm, t)\}$ 
40:  in  $UPD(Len(\mathcal{T}))$   $\triangleleft$  Return updates for  $rm$ 

```

A resource manager rm uses **TERMINATE** to vote at the termination of transaction t , of which it is a participant. If rm is willing to commit the transaction, then $vote$ equals **COMMIT** and upd contains the updates performed by t . To abort the transaction $vote$ equals **ABORT** and upd is the empty set. After casting its vote, rm waits until it learns t 's outcome. While waiting, rm monitors the other resource managers in $tset$, also involved in t . If rm suspects that some participant crashed, it votes **ABORT** on its behalf. After learning that t committed, rm will apply its updates (and possibly release the related locks, depending on the database). If rm learns that t aborted, it locally aborts the transaction. Updates are

Algorithm 1 (cont'd) Resource manager specification

```
41: INCARNATE( $rm, pid$ )  $\triangleq$ 
42:    $\mathcal{R}[rm] \leftarrow pid$ 
43:    $updates \leftarrow UPDATES(rm)$   $\triangleleft$  Get committed state.
44:   for  $i = 1$  to  $Len(updates)$   $\triangleleft$  For each set of committed transaction...
45:     apply updates in  $updates[i]$   $\triangleleft$  ...apply it to the database.

46: TERMINATE( $rm, t, tset, vote, upd$ )  $\triangleq$ 
47:   VOTE( $\langle rm, t, tset, vote, upd \rangle$ )
48:   while  $OUTCOME(t, rm) = UNDEFINED$ 
49:     wait ( $OUTCOME(t, rm) \neq UNDEFINED$ )  $\vee$  ( $suspect\ r \in tset$ )
50:     if  $suspected\ r \in tset$ 
51:       VOTE( $\langle r, t, tset, Abort, \emptyset \rangle$ )
52:     if  $OUTCOME(t, rm) = ABORT$ 
53:       abort  $t$  in the database
54:     else
55:       apply  $upd$  to database
```

made durable by the log service. We assume that the TERMINATE definition is not atomic, so multiple resource managers can vote in parallel, as well as the same resource manager can terminate distinct transactions in parallel, if its scheduling model allows it.

3.4 Correctness

Algorithm 1 solves non-blocking atomic commitment. AC-Validity and AC-Agreement are ensured by having the votes added to \mathcal{V} , from which OUTCOME deterministically evaluates to the transaction's outcome. If no resource manager is suspected and all vote COMMIT for a given transaction, then only COMMIT votes will be added to \mathcal{V} , and COMMIT is the only possible outcome for such a transaction, satisfying the AC-Non-Triviality property. Finally, since crashed resource managers are eventually suspected and have votes issued on their behalf, the protocol eventually terminates, satisfying the AC-Termination property.

The state of resource managers is changed by applying updates from write transactions. If these updates are deterministic, then given two copies of a resource manager in the same initial state and applying the same sequence of updates leads to the same final state. Assuming that transactions terminating concurrently do not interfere with each other and that update operations are deterministic, the partial order on updates provided by the log service and the deterministic method to apply them when reincarnating a resource manager, specified in Algorithm 1, correctly recovers the state before the crash and satisfies the R-Consistency property.

Formal proofs that the log service provides these properties as well as that it is implemented by protocols in Section 5 and 6 are given at Section 9.

4 From abstract specifications to distributed implementations

In this section we define the system model and some building blocks used in our implementations.

4.1 Processes, communications and failures

We assume a shared-nothing asynchronous distributed system composed of a set $\{p_1, p_2, \dots, p_n\}$ of processes. Processes may fail by crashing. Communication is by message passing using the primitives *send* and *receive*. These primitives guarantee quasi-reliable communication, that is, they ensure that if neither the sender nor the receiver of a message crashes, then the message is eventually delivered. This abstraction can be implemented on top of fair-lossy channels, for example, by periodically resubmitting messages until an acknowledgment is received by the sender. Quasi-reliable channels also ensure that messages are neither corrupted nor duplicated.

4.2 Leader-election oracle

Our coordinated implementation of the log service assumes a leader-election oracle. Participants use this oracle to determine the current coordinator. The leader-election oracle guarantees that eventually all participants will elect the same non-faulty process as the leader. Obviously, this can only be ensured if there is at least one process that eventually remains operational “forever”—in practical terms, long enough so that some useful computation can be done (e.g., deciding on a transaction’s outcome).

4.3 Consensus

In the consensus problem, a set of processes tries to agree on a common value. As in [9], we define the consensus problem over three process roles: *proposers* can propose a value (i.e., the request to be executed by the service), *acceptors* interact to choose a proposed value, and *learners* must learn the decided value. A process can play any number of these roles, and it is dubbed non-faulty if it remains up long enough to perform the consensus algorithm. Let n be the number of acceptors in the system and $f < n/2$ the maximum number of faulty processes. A correct consensus algorithm satisfies three properties:

C-Nontriviality Only a proposed value may be learned.

C-Consistency Any two learned values must be equal.

C-Progress For any proposer p and learner l , if p , l and more than $n/2$ acceptors are non-faulty, and p proposes a value, then l must learn a value.

The algorithms in Sections 5 and 6 use many instance of consensus. Therein, the primitive $propose(i, v)$ is used to propose a value v in the consensus instance i . The decision of instance i is the v in $decide(i, v)$.

5 Coordinated Implementation

The coordinated implementation is named after the coordinator, a process that serves as the log service’s interface to resource managers. The coordinator receives the votes from resource managers, makes them durable, and informs the resource managers about the transaction’s outcome. Votes are durable once they have been decided in some consensus instance; the coordinator proposes them in a sequence of instances, batching as many votes as possible to amortize the cost of termination.

For high availability, the service should replace the coordinator as soon as it is suspected to have failed. However, suspecting the coordinator too aggressively may incur unnecessary changes and different views among resource managers about which process is the coordinator. Resource managers use the leader-election oracle to elect the coordinator. The use of consensus ensures that safety is not violated

even if several processes become coordinators simultaneously. Once the leader-election oracle works properly, liveness will be ensured.

The implementation of the log service is decomposed in two parts. In Section 5.1 we present a set of stubs that run at the RMs and are used by Algorithm 1 to interact with the coordinator; in Section 5.2 we present the coordinator's protocol.

5.1 RM Stubs

The first part of Algorithm 2 initializes the data structures kept by the stubs: *outcome* is the local view of OUTCOME, and *myInc* and *rmPID* implement $\mathcal{R}[rm]$ locally.

Algorithm 2 Stubs to implement Algorithm 1.

```

1: Initially:
2:  $\forall t, outcome[t] \leftarrow \text{UNDEFINED}$   $\triangleleft$  All transactions are undecided.
3:  $myInc \leftarrow \perp$   $\triangleleft$  Have not incarnated yet.
4:  $rmPID \leftarrow \perp$   $\triangleleft$  Have not incarnated yet.

5: INCARNATE( $rm, pid$ )
6: send  $\langle \text{RECOVER}, pid, rm \rangle$  to current coordinator.
7: wait until receive  $\langle \text{RECOVERED}, rm, upd, inc \rangle$ 
8:  $myInc \leftarrow inc$ 
9:  $rmPID \leftarrow pid$   $\triangleleft$   $pid$  has incarnated  $rm$ .
10: return  $upd$ 

11:  $\mathcal{R}[rm]$ 
12: return  $rmPID$   $\triangleleft$  Either my own  $pid$  or  $\perp$ .

13: VOTE( $(rm, t, tset, vote, upd)$ )
14: if  $vote = \text{ABORT}$ 
15:    $outcome[t] \leftarrow \text{ABORT}$   $\triangleleft$  Quickly abort.
16: send  $\langle \text{VOTE}, rm, myInc, t, tset, vote, upd \rangle$  to coordinator

17: OUTCOME( $t$ )
18: return  $outcome[t]$ 

19: when receive  $\langle \text{TERMINATED}, t, out \rangle$   $\triangleleft$  Learn decision.
20:  $outcome[t] \leftarrow out$   $\triangleleft$  Learn  $t$ 's outcome.

21: when receive  $\langle \text{INCARNATE}, rm, newInc \rangle$   $\triangleleft$  Was replaced.
22:  $newInc > myInc$ 
23:  $rmPID \leftarrow \perp$   $\triangleleft$  No longer incarnates  $rm$ .

```

The INCARNATE stub sends a message RECOVER to the coordinator requesting to exchange the incarnation of resource manager *rm* to process *pid*, and waits for confirmation that the change was performed. The confirmation carries the updates executed by previous incarnations and the incarnation number of *rm*. $\mathcal{R}[rm] = pid$ evaluates to TRUE until the stub learns a bigger incarnation number, indicating that another processes took over the *rm*. The case in which more than one process believe to hold the rights on the *rm* is handled by the coordinator. To account for coordinator crashes, every time a process executes INCARNATE, it does it with a different *pid*.

VOTE forwards the resource manager's vote to the coordinator. As an optimization, if the vote is for ABORT, the resource manager aborts the transaction locally, as ABORT is the unique possible outcome of such a transaction.

The **when** clauses at the end of the algorithm execute fairly and atomically once their conditions become TRUE.

5.2 Log Service

The coordinator runs Algorithm 3. \mathcal{T}_c implements \mathcal{T} ; besides the votes of committed transactions, it also keeps the incarnation requests. \mathcal{V}_c implements \mathcal{V} as in the specification. B is the set of votes received but not yet treated by the coordinator. The coordinator executes a series of consensus instances, and i is the identifier of the instance the coordinator is waiting to terminate. Finally, $recSet$ is the set of reincarnation requests awaiting to be decided.

$\mathcal{R}[rm]$ is not explicitly kept in any data structure, but lies implicit in the subsequence of INCARNATE votes in \mathcal{T}_c : $\mathcal{R}[rm]$ equals the process in the last INCARNATE vote in \mathcal{T}_c .

When VOTE messages are received by the coordinator, they are added to B , to be proposed in the next consensus instance. Once B is no longer empty, the coordinator proposes it on instance i . For simplicity, we assume that B always fits in a consensus proposal. In some implementations of consensus, proposals may have bounded size. In such cases, B would have to be split and proposed in several consensus instances.

The coordinator waits the decision of instance i and postpones the learning of decisions of later instances. Once the coordinator learns the decision D , it first removes it from B , and then processes each of its elements. If the element is a vote then the coordinator determines its final meaning before adding it to \mathcal{V}_c . A vote's final meaning is ABORT if its issuer no longer equals $\mathcal{R}[rm]$; otherwise, it is the issuer's proposed vote. This ensures that upon recovery the new rm learns all committed transactions.

$\langle \text{RECOVER}, pid, rm \rangle$ messages asking the coordinator to change the $\mathcal{R}[rm]$ to pid are handled in the last part of the algorithm. Once such a message is received, the coordinator adds an INCARNATE vote to B . $\langle \text{INCARNATE}, pid, rm \rangle$ requests are handled after all the normal votes. They are simply appended to \mathcal{T}_c , meaning that $\mathcal{R}[rm]$ has been set to pid until another $\langle \text{INCARNATE}, pid', rm \rangle$ is appended to \mathcal{T}_c . Once the coordinator learns that $\mathcal{R}[rm]$ changed to pid , it gathers the updates performed by rm . These updates are sent to process pid , the new incarnation of rm , and a warn about the change is sent to all resource managers.

6 Uncoordinated Implementation

The main idea behind the uncoordinated log service implementation is to save time by not going through a central process. Instead of forwarding messages to the coordinator and waiting for replies as in the coordinated version, participants access local data structures and execute consensus.

Algorithm 4 defines stubs to Algorithm 1. The resource manager rm keeps \mathcal{V}_{rm} , the set of votes it has received, $commitCounter$, the counter of committed transactions involving rm , and $rmPID$, used to locally evaluate $\mathcal{R}[rm]$. The system's multiprogramming level, that is, the maximum number of transactions that can be executed concurrently by a resource manager corresponds to k . Therefore, k is also the maximum number of started transactions yet to commit, and is paramount to the uncoordinated approach. In the coordinated approach, the coordinator changes the incarnation of a resource manager by having an INCARNATE transaction added to \mathcal{T}_c . Because \mathcal{T}_c grows deterministically due to the ordering of consensus instances and their decisions, any process acting as the coordinator will see incarnation requests in the same order. In the uncoordinated approach, however, there is no total order of decisions, and a process can consider itself incarnated only after a consensus instance corroborates it,

Algorithm 3 Coordinator's protocol.

```

1: Initialization:
2:  $\mathcal{T}_c \leftarrow \emptyset$   $\triangleleft$  The sequence of terminated transactions.
3:  $\mathcal{V}_c \leftarrow \emptyset$   $\triangleleft$  The set of durable votes.
4:  $B \leftarrow \emptyset$   $\triangleleft$  Votes to be broadcast.
5:  $i \leftarrow 0$   $\triangleleft$  Current consensus instance.
6:  $recSet \leftarrow \emptyset$   $\triangleleft$  Awaiting reincarnation transactions.

7:  $\mathcal{R}[rm]$ 
8:  $rrm \leftarrow r = \langle \text{INCARNATE}, pid, rm \rangle \in \mathcal{T}_c : \forall r' = \langle \text{INCARNATE}, pid', rm \rangle \in \mathcal{T}_c, r >_{\mathcal{T}_c} r'$ 
9: return  $rrm[2]$ 

10: when receive  $\langle \text{VOTE}, rm, pid, t, tset, vote, upd \rangle$ 
11: if  $\neg \exists \langle rm, -, t, -, -, - \rangle \in B$ 
12:    $B \leftarrow B \cup \{ \langle rm, pid, t, tset, vote, upd \rangle \}$   $\triangleleft$  Only the first vote is considered.

13: when  $B \neq \emptyset$ 
14:   propose  $(i, B)$   $\triangleleft$  Propose  $B$  on instance  $i$ .

15: when decide  $(i, D)$   $\triangleleft$  Decided  $D$  on instance  $i$ .
16:    $B \leftarrow B \setminus D$   $\triangleleft$  Remove from next proposals.
17:    $C \leftarrow \emptyset$   $\triangleleft$  Temporary variable.
18:   for all  $\langle rm, pid, t, tset, vote, upd \rangle \in D$ 
19:     if  $\neg \exists \langle rm, -, t, -, -, - \rangle \in \mathcal{V}_c$ 
20:        $prevState \leftarrow \text{OUTCOME}(t)$ 
21:       if  $\mathcal{R}[rm] \neq pid$   $\triangleleft$  If  $rm$  has been reincarnated
22:          $\mathcal{V}_c \leftarrow \mathcal{V}_c \cup \langle rm, t, tset, \text{ABORT}, \emptyset \rangle$   $\triangleleft$  turn the vote into ABORT
23:       else
24:          $\mathcal{V}_c \leftarrow \mathcal{V}_c \cup \{ \langle rm, t, tset, vote, upd \rangle \}$   $\triangleleft$  otherwise add it to the set of votes.
25:         if  $(prevState = \text{UNDEFINED}) \wedge (\text{OUTCOME}(t) = \text{COMMIT})$   $\triangleleft$  If vote lead to commit of  $t$ 
26:            $C \leftarrow C \cup \{t\}$   $\triangleleft$  add it to  $C$ 
27:            $\forall p \in tset, \text{send message } \langle \text{TERMINATED}, t, \text{OUTCOME}(t) \rangle \text{ to } p$   $\triangleleft$  and warn participants.
28:           else if  $vote = \text{ABORT}$   $\triangleleft$  If lead to abortion of  $t$ 
29:              $\forall p \in tset, \text{send message } \langle \text{TERMINATED}, t, \text{ABORT} \rangle \text{ to } p$   $\triangleleft$  warn all possible.
30:            $\mathcal{T}_c \leftarrow \mathcal{T}_c \bullet C$   $\triangleleft$  Store the partial order.
31:           for all  $d = \langle \text{INCARNATE}, -, - \rangle \in D$   $\triangleleft$  Process INCARNATE votes.
32:              $\mathcal{T}_c \leftarrow \mathcal{T}_c \bullet d$ 
33:              $i \leftarrow i + 1$   $\triangleleft$  Process the next batch.

34: when receive  $\langle \text{RECOVER}, pid, rm \rangle$   $\triangleleft$  Upon request to recover
35:    $t_{inc} \leftarrow \langle \text{INCARNATE}, pid, rm \rangle$   $\triangleleft$  create a "reincarnation transaction"
36:    $recSet \leftarrow recSet \cup \{t_{inc}\}$   $\triangleleft$  put it in the awaiting set
37:    $B \leftarrow B \cup t_{inc}$   $\triangleleft$  and make it terminate.

38: when  $\exists t_{inc} = \langle \text{INCARNATE}, pid, rm \rangle \in recSet : t_{inc} \in \mathcal{T}_c$   $\triangleleft$  Once a reincarnation transaction terminates
39:    $recSet \leftarrow recSet \setminus \{t_{inc}\}$ 
40:    $\mathcal{V}^{rm} \leftarrow \{e = \langle rm, t, -, -, - \rangle \in \mathcal{V}_c : (\text{OUTCOME}(t) = \text{COMMIT}) \wedge (t <_{\mathcal{T}_c} t_{inc})\}$   $\triangleleft$  determine the updates already performed
41:    $\mathcal{U}^{rm} \leftarrow \langle u : \langle rm, -, -, u \rangle \in \mathcal{V}^{rm}, \forall \langle rm, t_1, -, u_1 \rangle, \langle rm, t_2, -, u_2 \rangle \in \mathcal{V}^{rm}, u_1 <_{\mathcal{U}^{rm}} u_2 \Rightarrow t_1 <_{\mathcal{T}_c} t_2 \rangle$   $\triangleleft$  ordered as in  $\mathcal{T}$ .
42:    $inc \leftarrow | \{r = \langle \text{INCARNATE}, -, rm \rangle \in \mathcal{T} : r <_{\mathcal{T}} \langle \text{INCARNATE}, pid, rm \rangle \} |$   $\triangleleft$  Determine the incarnation number
43:   send  $\langle \text{RECOVERED}, rm, \mathcal{U}^{rm}, inc \rangle$  to  $rm$   $\triangleleft$  and warn interested processes.
44:   send  $\langle \text{INCARNATE}, rm, inc \rangle$  to all resource managers

```

i.e., an INCARNATE transaction is decided, and it is sure that all the consensus instances that the previous incarnation had possibly voted have been treated. There can be up to $k - 1$ of such consensus instances.

6.1 Transaction Termination

The VOTE stub proposes rm 's vote for transaction t on consensus instance $inst(rm, t)$. If the rm is voting on its own behalf, the current consensus instance is stored locally. To vote for other processes, the rm must know the instances they attributed to t —notice that the rm only cares about other resource managers if it is voting COMMIT. The rm can learn the consensus instance of other rm 's from the transaction manager, when it receives the commit request. The transaction manager receives the consensus instance from the resource managers during the execution of the transactions. With the vote, the resource manager proposes also the number of transactions already committed, $commitCounter$. This information is used later during recovery, as we explain below.

The decision of a consensus instance j is learned by rm if it is one of its instances, or one associated to a transaction in which rm is participant. This decision can be of three types: (i) a vote, in which case it is added to \mathcal{V}_{rm} ; (ii) a change of incarnation of another resource manager, in which case it is added to \mathcal{V}_{rm} as an ABORT vote; and (iii) a change of rm 's incarnation to process pid' , in which case the process $pid \neq pid'$ currently incarnating learns that its incarnation lasts only until instance $j + k$. The **when** clause that processes the decisions is only activated once the resource manager has completed recovery.

6.2 Recovering from Failures

The INCARNATE stub determines the updates performed by previous incarnations and in which instance the new incarnation starts. The procedure consists of three steps: in the first step, pid (i.e., the process incarnating rm) proposes $\langle \text{INCARNATE}, pid, rm \rangle$ in its first consensus instance and until such a value is decided in an instance i , at which point pid becomes the incarnation of rm for transactions associated with consensus instances $i + k$, until displaced by another process. The second step terminates the transactions associated with instances $i + 1$ to $i + k - 1$, belonging to the incarnations being replaced. The third step determines the outcome of all transactions associated with instances smaller than $i + k$.

The protocol terminates by gathering the updates of committed transactions in the sequence \mathcal{U}_{rm} . Elements in \mathcal{U}_{rm} are ordered according to the commit counter in each vote, making it consistent with the commit order of non-concurrent transactions. Algorithm 4 does not keep the \mathcal{T} data structure nor any abridged version of it; the order of committed transactions is forgotten by the resource managers once their updates are applied. Upon recovery, \mathcal{U}_{rm} is created and kept just until the recovery is terminated.

To improve performance, implementations should use checkpoints and reduce the number of consensus instances in which resource managers must propose. Most of the remaining instances can run in parallel to reduce the recovery latency.

7 Evaluation

We performed the experiments in a cluster of nodes from the Emulab testbed [13]. Nodes were equipped with 64-bit Xeon 3GHz processors, and interconnected through a Gigabit Ethernet switch. Our communication library used datagrams of size up to 7500B.

In the graphs, points represent the average of values sampled after execution stabilization in each run of an experiment.

Algorithm 4 Algorithm 1 stubs and Algorithm 1 uncoordinated implementation, with multi-programming level equal to k (at resource manager rm).

```

1: Initialization
2:  $\mathcal{V}_{rm} \leftarrow \emptyset$ 
3:  $commitCounter \leftarrow 0$ 
4:  $rmPID \leftarrow \perp$ 
5:  $\mathcal{R}[rm]$ 
6: return  $rmPID$ 
7: OUTCOME( $t$ )
8: if  $\exists \langle -, t, -, \text{ABORT}, - \rangle \in \mathcal{V}_{rm}$ 
9:   ABORT
10: else if  $\exists \langle -, t, tset, -, - \rangle \in \mathcal{V}_{rm} : \forall s \in tset :$ 
11:    $\langle s, t, tset, \text{COMMIT}, - \rangle \in \mathcal{V}_{rm}$ 
12:   COMMIT
13:   UNDEFINED
14: VOTE( $\langle rm, t, tset, vote, update_{rm} \rangle$ )
15: propose( $inst(rm, t), \langle rm, t, tset, vote, upd, Len(\mathcal{T}_{rm}) \rangle$ )
16: when decide( $j, d$ )
17: if  $d = \langle r, t, tset, vote, upd, cn \rangle, rm \in tset$ 
18:    $\mathcal{V}_{rm} \leftarrow \mathcal{V}_{rm} \cup \{ \langle r, t, tset, vote, upd \rangle \}$ 
19: else if  $d = \langle \text{INCARNATE}, -, r \rangle, r \neq rm$ 
20:    $\mathcal{V}_{rm} \leftarrow \mathcal{V}_{rm} \cup \{ \langle r, t, \{r, rm\}, \text{ABORT}, \emptyset \rangle, inst(rm, t) = j \}$ 
21: else if  $d = \langle \text{INCARNATE}, pid', rm \rangle, pid \neq pid'$ 
22:    $rmPID \leftarrow \perp$ 
23: INCARNATE( $rm$ )
24:  $tranSet \leftarrow \emptyset$ 
25:  $i \leftarrow 0$ 
26: while TRUE
27:   propose( $i, \langle \text{INCARNATE}, pid, rm \rangle$ )
28:   wait until decide( $i, d$ )
29:   if  $d = \langle \text{INCARNATE}, pid, rm \rangle$ 
30:      $pidRM \leftarrow pid$ 
31:     break
32:   else
33:      $tranSet \leftarrow tranSet \cup \{d\}$ 
34:      $i \leftarrow i + 1$ 
35: for  $j \leftarrow 1..k - 1$ 
36:   propose( $i + j, \langle rm, t, \{rm\}, \text{ABORT}, \emptyset, 0 \rangle$ )
37:   wait until decide( $i + j, d$ )
38:   if  $d = \langle \text{INCARNATE}, -, rm \rangle$ 
39:      $rmPID \leftarrow \perp$ 
40:     return  $\langle \rangle$ 
41:    $tranSet \leftarrow tranSet \cup \{d\}$ 
42: for each  $\langle rm, t, l, v, u, cn \rangle \in tranSet$ 
43:   if  $\neg \exists \langle p', t, l', \text{ABORT}, \emptyset, cn' \rangle \in V_{rm}$ 
44:     for all  $p' \in l, \neg \exists \langle p', t, l, v', u', cn' \rangle \in V_{rm}$ 
45:       propose( $inst(p', t), \langle p', t, l, \text{ABORT}, \emptyset, 0 \rangle$ )
46:       wait until decide( $inst(p', t), d$ )
47:        $V_{rm} \leftarrow V_{rm} \cup \{d\}$ 
48:       if  $d = \langle p', t, -, \text{ABORT}, -, - \rangle$ 
49:         break
50:    $U_{rm} \leftarrow \langle u_1, \dots, u_m \rangle : e_k = \langle rm, t_k, l_k, v_k, u_k, cn_k \rangle \in V_{rm},$ 
51:    $\forall p \in l_k, \exists \langle p, t_k, l_k, \text{COMMIT}, u_p \rangle \in V_{rm},$ 
52:    $\forall i < j, cn_i < cn_j$ 

```

\triangleleft Votes I have seen.
 \triangleleft How many transactions I committed.
 \triangleleft Have not incarnated yet.
 \triangleleft Either my own pid or \perp .
 \triangleleft Any ABORTs?
 \triangleleft All COMMITs?
 \triangleleft Neither one nor the other
 \triangleleft Vote + number of committed.
 \triangleleft [Decided on instance j.]
 \triangleleft Someone else was substituted.
 \triangleleft Make d an ABORT vote.
 \triangleleft rm has been reincarnated,
 \triangleleft so give it up.
 \triangleleft Assume that 0 is the first instance identifier.
 \triangleleft Vote until incarnation changes,
 \triangleleft Reincarnated
 \triangleleft Stop the first iteration.
 \triangleleft and then, for $k - 1$ possibly open instances,
 \triangleleft vote to close them.
 \triangleleft Someone else is recovering.
 \triangleleft Give up the resource manager
 \triangleleft and stop looking for updates.
 \triangleleft Close all of those transactions:
 \triangleleft If t was not aborted,
 \triangleleft make sure to terminate it
 \triangleleft voting for others if needed.
 \triangleleft Stop on ABORT.
 \triangleleft Now order updates of known transactions
 \triangleleft that committed
 \triangleleft in commit order.

7.1 Micro-benchmark

In the first experiment we used a micro-benchmark of non-conflicting transactions, each comprising one update operation being executed in parallel at 8 RMs, and resulting, at each of them, in 100 or 7000 B of data to be logged. Figure 1 shows that for both update sizes the coordinated implementation outperforms the uncoordinated one by an order of magnitude. In case of small updates the coordinated version must perform eight times less consensus instances and in addition the coordinator can batch votes from up to nine transactions executed concurrently, while our implementation of a single RM is not capable of this grouping. For big updates the coordinator has to run as many instances as the uncoordinated version does, but the difference might be explained by the consensus executions' time: The coordinator serializes its instances so they do not interfere with each other and can be timely finished, what is supported by our observation that only a few retries happened. In the uncoordinated implementation all the RMs try to execute their instances at the same time, and as the network contention leads to many unsuccessful termination attempts, the execution becomes driven by the timeouts used.

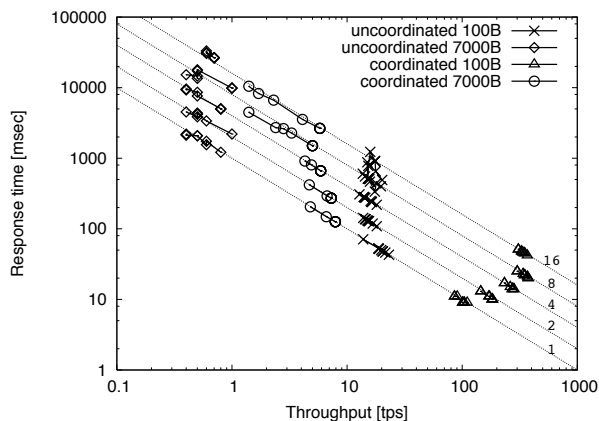


Figure 1: Throughput versus response time of the micro-benchmark. The number of clients is shown next to the lines.

7.2 The TPC-C benchmark

For TPC-C the read-only Items table was replicated on all RMs and other tables were range partitioned among RMs according to the warehouse id, so at most 15% transactions involved more than one RM. Update transactions were 92% of the workload and none of them produced data to be logged exceeding 1500B.

As Figure 2 shows, very small loads do not differ significantly in performance between configurations. By higher loads the coordinated version outperforms the uncoordinated one and the former scales much better with the number of transactions executed concurrently, which is determined by the number of client threads and RMs. Although in 85% of the cases a single transaction requires one consensus instance regardless of the termination protocol, the coordinated version can group at least five simultaneous requests, and even when a single RM could perform the batching on its own the coordinator still have higher potential for that as it can combine data from different RMs.

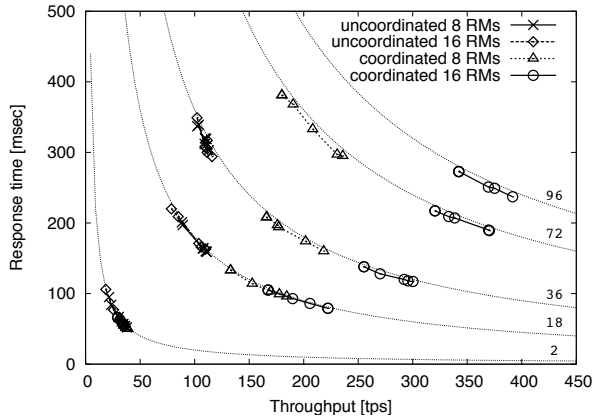


Figure 2: Throughput versus response time of TPC-C transactions. The number of clients is shown next to the curves.

8 Related Work

Stamus and Cristian [12] propose to aggregate log records generated by resource managers at the transaction managers and logged there. This approach resembles the log service we are proposing. In fact, the transaction manager acts as a log service to the transactions it manages. However, their approach uses a byzantine agreement abstraction with strong synchrony assumptions, and do not consider recovery of resource and transaction managers in nodes different from the initial ones.

The log service proposed by Daniels *et al.* [4] does consider recovery on different machines, but it considers one instance of the service per resource manager instead of a shared one, and therefore cannot provide the termination facilities of our approach.

The use of consensus or agreement protocols in atomic commitment protocols has been explored in some other works, but none of them considers the problem of interleaving commits and transaction ordering. Mohan *et al.* [11] used byzantine tolerant agreement abstractions to extend a two phase commit protocol with byzantine reliability. As in our protocols, the recovery of resource managers in their approach is based on the log provided by the agreement boxes, although not abstracted by a log service. However, they keep the “two-phases approach”, even though one agreement phase should suffice to solve the problem.

Guerraoui *et al.* [6] presented a non-blocking protocol that terminates transactions in three communication steps, from termination request to transaction termination, in good runs (i.e., without failures and suspicions). In bad runs, the protocol resorts to consensus to agree on the outcome, exchanging at least one communication step for a consensus instance.

More recently, Lamport and Gray [5] presented the Paxos Commit protocol, a non-blocking generalization of two-phase commit using the Paxos consensus protocol. The biggest difference between this protocol and the previous one is that Paxos Commit uses consensus to decide on votes, not on the transaction’s outcome. As Paxos does not require all participants to propose, a suspicious resource manager votes on the suspect’s behalf as soon as suspecting it, and may have its conservative vote or the suspect’s original vote decided in two communication steps, increasing the protocol’s latency by the time to suspect a participant. Having several consensus instances in parallel, for possibly many transactions, has a bad impact on the overall performance. We looked for a way to isolate this impact and it was by mimicking Paxos’ separation of concerns that the log service was devised.

9 Conclusion

In this paper we have introduced the specification of a log service for transaction processing systems. The log service provides atomicity and durability to transactions through a non-blocking termination protocol. The service totally orders non-concurrent transactions and, should a resource manager fail, the service can be used to recover the resource manager's state prior to the crash and start a copy of it on a different and functional node. The service safely copes with multiple copies of a resource manager.

We also presented two highly available implementations of our log service, namely coordinated and uncoordinated, and provided a comparative experimental performance evaluation. In the studied scenarios, a coordinated approach has led to a much higher transaction throughput and smaller response times.

References

- [1] M. Abadi and L. Lamport. The existence of refinement mappings. *Theoretical Computer Science*, 82(2):253–284, 1991.
- [2] P. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.
- [3] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Communications of the ACM*, 43(2):225–267, 1996.
- [4] D. S. Daniels, A. Z. Spector, and D. S. Thompson. Distributed logging for transaction processing. In U. Dayal and I. Traiger, editors, *Proceedings of the ACM SIGMOD Annual Conference*, pages 82–96, San Francisco, CA, 1987. ACM, ACM Press.
- [5] J. Gray and L. Lamport. Consensus on transaction commit. *ACM TODS*, 31(1):133–160, March 2006.
- [6] R. Guerraoui, M. Larrea, and A. Schiper. Reducing the cost for non-blocking in atomic commitment. In *ICDCS '96: Proceedings of the 16th International Conference on Distributed Computing Systems (ICDCS '96)*, page 692, Washington, DC, USA, 1996. IEEE Computer Society.
- [7] L. Lamport. Introduction to tla. Technical Report 1994-001, Palo Alto, CA, 1994.
- [8] L. Lamport. Refinement in state-based formalisms, 1996.
- [9] L. Lamport. The part-time parliament. *ACM Trans. Comput. Syst.*, 16(2):133–169, May 1998.
- [10] L. Lamport. *Specifying Systems: The TLA+ Language and Tools for Hardware and Software Engineers*. Addison-Wesley Professional, July 2002.
- [11] C. Mohan, R. Strong, and S. Finkelstein. Method for distributed transaction commit and recovery using byzantine agreement within clusters of processors. *SIGOPS Oper. Syst. Rev.*, 19(3):29–43, July 1985.
- [12] J. W. Stamos and F. Cristian. Coordinator log transaction execution protocol. *Distributed and Parallel Databases*, 1(4):383–408, 1993.
- [13] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. In *Proc. of the Fifth Symposium on Operating Systems Design and Implementation*, pages 255–270, Boston, MA, Dec. 2002. USENIX Association.

The Service's variables:

vHist the history of votes.
tHist the history of committed transactions
LastConcSet LastConcSet concurrent to the last committed.
terminatingAt transactions terminating at rm.
terminatedAt transactions terminated at rm.

VARIABLES *vHist*, *tHist*, *rm2pid*, **Log service's**
LastConcSet

The Resource Managers' variables:

terminatingAt transactions terminating at rm.
terminatedAt transactions terminated at rm.

VARIABLES *terminatingAt*, **Resource Manager's**
terminatedAt,
pid2rm

The Transaction Managers's variables:

termReq transactions requested to terminate.
part participants on each transaction.

VARIABLES *termReq*, **Transaction Managers'**
part

Defining some aliases.

svars \triangleq $\langle vHist, tHist, LastConcSet, rm2pid \rangle$
rvars \triangleq $\langle terminatingAt, terminatedAt, pid2rm \rangle$
tvars \triangleq $\langle part, termReq \rangle$
evars \triangleq $\langle badProc, suspect \rangle$
avars \triangleq $\langle svars, rvars, tvars, evars \rangle$

Types

UNKNOWN \triangleq CHOOSE $v : v \notin \{\text{TRUE}, \text{FALSE}\}$

Votes \triangleq [*rm* : *RM*, *tr* : *TID*, *tset* : SUBSET *RM*, *vt* : {"Commit", "Abort"}, *upd* : *Update*]

Invariants

Type invariant.

TypeInvariant \triangleq

- $\wedge vHist \in \text{SUBSET } Votes$
- $\wedge tHist \in \text{Seq}(\text{SUBSET } TID)$
- $\wedge LastConcSet \in \text{SUBSET } TID$
- $\wedge rm2pid \in [RM \rightarrow PID \cup \{NoPID\}]$
- $\wedge terminatingAt \in [PID \rightarrow \text{SUBSET } TID]$
- $\wedge terminatedAt \in [PID \rightarrow \text{SUBSET } TID]$
- $\wedge pid2rm \in [PID \rightarrow RM \cup \{NoRM\}]$
- $\wedge part \in [TID \rightarrow \text{UNION } \{[S \rightarrow PID] : S \in \text{SUBSET } RM\}]$
- $\wedge termReq \in \text{SUBSET } TID$

$\wedge badProc \in \text{SUBSET } PID$
 $\wedge suspect \in [PID \rightarrow [PID \rightarrow \text{BOOLEAN}]]$

Initial State

$Init \triangleq$
 $\wedge vHist = \{\}$ No vote received.
 $\wedge tHist = \langle \rangle$ No transaction committed.
 $\wedge LastConcSet = \{\}$
 $\wedge rm2pid = [r \in RM \mapsto NoPID]$ No RM is incarnated.
 $\wedge terminatingAt = [p \in PID \mapsto \{\}]$ No transactions terminating at any RM .
 $\wedge terminatedAt = [p \in PID \mapsto \{\}]$ No transactions terminating at any RM .
 $\wedge pid2rm = [r \in PID \mapsto NoRM]$ No RM is incarnated.
 $\wedge part = [t \in TID \mapsto [e \in \{\} \mapsto \{\}]]$ No participant in any transaction.
 $\wedge termReq = \{\}$ No termination request issued.
 $\wedge badProc = \{\}$ No bad process.
 $\wedge suspect = [p \in PID \mapsto [q \in PID \mapsto \text{FALSE}]]$

Operators

Operator $IsInvolved(t, rm)$ checks whether rm is involved in transaction t .

$IsInvolved(t, rm) \triangleq$ IF $\exists v \in vHist : v.tr = t \wedge rm \in v.rm$
 THEN TRUE
 ELSE IF $\exists v \in vHist : v.tr = t \wedge v.vt = \text{“Commit”}$
 THEN FALSE
 ELSE UNKNOWN

Operator $Updates(rm)$ gets the updates performed by rm in committed transactions.

$Updates(rm) \triangleq$
 LET $UpdateOn(t) \triangleq (\text{CHOOSE } v \in vHist : v.tr = t \wedge v.rm = rm).upd$
 $upd[i \in 0 .. Len(tHist)] \triangleq$
 IF $i = 0$
 THEN $\langle \rangle$
 ELSE $Append(upd[i - 1], \{UpdateOn(t) : t \in \{e \in tHist[i] : IsInvolved(rm, e) = \text{TRUE}\})$
 IN $upd[Len(tHist)]$

Operator $Outcome(h, t)$ gives the termination status of transaction t , considering the history of votes h .

$Outcome(h, t) \triangleq$ IF $\exists v \in h : v.tr = t \wedge v.vt = \text{“Abort”}$
 THEN “Abort”
 ELSE IF $\exists v \in h :$
 $\wedge v.tr = t$
 $\wedge \forall p \in v.tset :$
 $\exists vv \in h : \wedge vv.rm = p$
 $\wedge vv.tr = t$

$$\wedge vv.vt = \text{“Commit”}$$

THEN “Commit”
ELSE “Undefined”

Actions

Environment

This action crashes a good process.

$$\text{Crash} \triangleq \wedge \exists p \in (PID \setminus badProc) : badProc' = badProc \cup \{p\} \\ \wedge \text{UNCHANGED} \langle suspect \rangle$$

This action changes the suspicion status.

$$\text{ChangeSuspicion} \triangleq \wedge \exists p \in (PID \setminus badProc), q \in PID : \\ \wedge p \neq q \\ \wedge suspect' = [suspect \text{ EXCEPT } ![p][q] = \neg @] \\ \wedge \text{UNCHANGED} \langle badProc \rangle$$

EnvActions:

- process crash.
- change suspicions.

$$\text{EnvActions} \triangleq \wedge \vee \text{Crash} \\ \vee \text{ChangeSuspicion} \\ \wedge \text{UNCHANGED} \langle svars, rvars, tvars \rangle$$

Transaction Manager

Adds a resource manager to a non-terminated transaction.

The transaction manager will only succeed in the adding a resource manager r to a transaction t if r was not crashed by the time it was contacted, and replied to operation request.

$$\text{AddRM} \triangleq \wedge \exists t \in (TID \setminus termReq), \quad \begin{array}{l} t \text{ has not tried to terminate yet.} \\ r \text{ has been incarnated,} \\ \text{replied to an operation,} \\ \text{and was not a participant yet} \end{array} \\ r \in \{r \in RM : \wedge rm2pid[r] \neq NoPID \\ \wedge rm2pid[r] \notin badProc\} : \\ \wedge r \notin \text{DOMAIN } part[t] \\ \wedge part' = [part \text{ EXCEPT } ![t] = \\ [e \in \text{DOMAIN } @ \cup \{r\} \mapsto \text{IF } e \in \text{DOMAIN } @ \\ \text{THEN } @[e] \\ \text{ELSE } rm2pid[r]]] \quad \text{now it is.} \\ \wedge \text{UNCHANGED} \langle termReq, rm2pid \rangle$$

Request the termination of a transaction.

The transaction manager can, at any time, try to terminate a transaction it has not tried to terminate before, and that executed some operation.

$$\text{RequestTerm} \triangleq \wedge \exists t \in (TID \setminus termReq) : \quad \begin{array}{l} t \text{ has not tried to terminate yet.} \\ \wedge part[t] \neq \langle \rangle \\ \wedge termReq' = termReq \cup \{t\} \end{array}$$

$\wedge \text{UNCHANGED } \langle part, rm2pid \rangle$

The disjunction of Transaction Manager's actions.

- Add a resource manager as a participant.
- Try to terminate a transaction.

$TMActions \triangleq \wedge \vee AddRM$
 $\vee RequestTerm$
 $\wedge \text{UNCHANGED } \langle svars, rvars, evars \rangle$

Log Service

Action $Incarnate(pid, rm)$ is executed by process pid to incarnate resource manager rm .

$Incarnate \triangleq$
 LET $ApplyUpdates(u) \triangleq$
 LET $apply[i \in 1 .. Len(u)] \triangleq$
 IF $i = Len(u)$ THEN $ApplyUpdate(u[i])$
 ELSE $ApplyUpdate(u[i]) \wedge apply[i + 1]$
 IN IF $u = \langle \rangle$ THEN TRUE
 ELSE $apply[1]$
 IN $\exists p \in PID \setminus badProc, r \in RM :$
 $\wedge pid2rm[p] = NoRM$ but is not incarnating,
 $\wedge \vee rm2pid[r] = NoPID$ and r is not incarnated
 $\vee rm2pid[r] \neq NoPID \wedge suspect[p][rm2pid[r]]$ or its process is suspected.
 $\wedge rm2pid' = [rm2pid \text{ EXCEPT } ![r] = p]$
 $\wedge pid2rm' = [pid2rm \text{ EXCEPT } ![p] = r]$
 $\wedge ApplyUpdates(Updates(r))$
 $\wedge \text{UNCHANGED } \langle terminatingAt, terminatedAt, vHist, tHist, LastConcSet \rangle$

Action $Vote(v)$ adds v to $vHist$, if not there yet.

$Vote(v) \triangleq$
 LET $ConcSet(vh) \triangleq \{t \in \{v.tr : \wedge v \in vh\} :$ Voted transactions
 $Outcome(vh, t) = \text{"Undefined"}\}$ not terminated yet.
 $Commits \triangleq \wedge Outcome(vHist', v.tr) = \text{"Commit"}$
 $\wedge \vee \wedge v.tr \in LastConcSet$
 $\wedge tHist' = [tHist \text{ EXCEPT } ![Len(tHist)] = @ \cup \{v.tr\}]$
 $\wedge LastConcSet' = LastConcSet \setminus \{v.tr\}$
 $\vee \wedge v.tr \notin LastConcSet$
 $\wedge tHist' = Append(tHist, \{v.tr\})$
 $\wedge LastConcSet' = ConcSet(vHist') \setminus \{v.tr\}$
 $\wedge \text{UNCHANGED } rm2pid$
 $Aborts \triangleq \wedge Outcome(vHist', v.tr) = \text{"Abort"}$
 $\wedge \text{UNCHANGED } \langle tHist, LastConcSet, rm2pid \rangle$

$$\begin{aligned}
\text{IN } & \wedge \neg \exists ov \in vHist : \wedge ov.rm = v.rm \\
& \wedge ov.tr = v.tr && \text{Add } v \text{ to } vHist \\
& \wedge vHist' = vHist \cup \{v\} \\
& \wedge StayUndef \vee Commits \vee Aborts
\end{aligned}$$

Action $Terminate(rm, t)$ is executed by rm to step towards transaction t 's termination.

$$\begin{aligned}
VoteForMyself(r, t) & \triangleq \\
& \wedge t \in (termReq \\
& \quad \setminus (terminatingAt[rm2pid[r]] \cup terminatedAt[rm2pid[r]])) \quad \text{and has not tried to terminate } t. \\
& \wedge terminatingAt' = [terminatingAt \text{ EXCEPT } ![rm2pid[r]] = @ \cup \{t\}] \quad \text{Start the termination} \\
& \wedge \vee Vote([rm \mapsto r, tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
& \quad vt \mapsto \text{"Commit"}, upd \mapsto GetUpdate(r, t)]) \quad \text{voting commit.} \\
& \quad \vee Vote([rm \mapsto r, tr \mapsto t, tset \mapsto \{r\}, vt \mapsto \text{"Abort"}, upd \mapsto \{\}]) \quad \text{voting abort.} \\
& \wedge \text{UNCHANGED } \langle terminatedAt \rangle
\end{aligned}$$

$$\begin{aligned}
VoteForOthers(r, t) & \triangleq \\
& \wedge t \in terminatingAt[rm2pid[r]] \quad \text{rm tried to terminate } t \\
& \wedge Outcome(vHist, t) = \text{"Undefined"} \quad \text{t is stuck.} \\
& \wedge \exists s \in \text{DOMAIN } part[t] : \\
& \quad \wedge suspect[rm2pid[r]][rm2pid[s]] \\
& \quad \wedge Vote([rm \mapsto s, tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
& \quad \quad vt \mapsto \text{"Abort"}, upd \mapsto \{\}]) \\
& \wedge \text{UNCHANGED } \langle terminatingAt, terminatedAt \rangle
\end{aligned}$$

$$\begin{aligned}
Learn(r, t) & \triangleq \\
& \wedge Outcome(vHist, t) \neq \text{"Undefined"} \quad \text{t has terminated.} \\
& \wedge terminatingAt' = [terminatingAt \text{ EXCEPT } ![rm2pid[r]] = @ \setminus \{t\}] \quad \text{So } rm \text{ learns it.} \\
& \wedge terminatedAt' = [terminatedAt \text{ EXCEPT } ![rm2pid[r]] = @ \cup \{t\}] \\
& \wedge \text{UNCHANGED } \langle svars \rangle
\end{aligned}$$

$$\begin{aligned}
Terminate & \triangleq \\
& \exists r \in RM, t \in TID : \\
& \quad \wedge r \in \text{DOMAIN } part[t] \quad \text{r is a participant of } t. \\
& \quad \wedge part[t][r] \notin badProc \quad \text{the process is alive,} \\
& \quad \wedge part[t][r] = rm2pid[r] \quad \text{and is still the same} \\
& \quad \wedge \vee VoteForMyself(r, t) \\
& \quad \quad \vee VoteForOthers(r, t) \\
& \quad \quad \vee Learn(r, t) \\
& \quad \wedge \text{UNCHANGED } \langle pid2rm \rangle
\end{aligned}$$

$$\begin{aligned}
RMActions & \triangleq \\
& \wedge \vee Incarnate \\
& \quad \vee Terminate \\
& \wedge \text{UNCHANGED } \langle tvars, evars \rangle
\end{aligned}$$

Specification

$$\begin{aligned} Next &\triangleq \vee RMActions \\ &\vee TMActions \\ &\vee EnvActions \end{aligned}$$

$$Spec \triangleq Init \wedge \Box [Next]_{\langle svars, tvars, rvars, evars \rangle}$$

Theorems

The specification is type safe.

THEOREM $Spec \Rightarrow \Box TypeInvariant$

AC-Validity If an *RM* decides to commit a transaction, then all *RMs* voted to commit the transaction.

$$\begin{aligned} AC_Validity &\triangleq \forall t \in TID : Outcome(vHist, t) = \text{“Commit”} \\ &\Rightarrow \forall r \in \text{DOMAIN } part[t] : \\ &\quad \exists v \in vHist : \wedge v.rm = r \\ &\quad \wedge v.tr = t \\ &\quad \wedge v.vote = \text{“Commit”} \end{aligned}$$

AC-Agreement It is impossible for one *RM* to commit a transaction and another one to abort the transaction.

$$AC_Agreement \triangleq \text{TRUE} \quad \text{From the definition of } Outcome.$$

AC-Non-Triviality If all *RMs* vote to commit the transaction and no *RM* is suspected throughout the execution of the protocol, then the decision is commit.

$$\begin{aligned} AC_Non_Triviality &\triangleq \forall t \in TID : \wedge \forall r, s \in part[t] : \Box (\neg suspect[rm2pid[r]][rm2pid2[s]]) \\ &\quad \wedge \forall r \in part[t] : \exists v \in vHist : \wedge v.rm = r \\ &\quad \wedge v.t = t \\ &\quad \wedge v.vote = \text{“Commit”} \\ &\Rightarrow \\ &\quad \wedge \Diamond Outcome(vHist) = \text{“Commit”} \end{aligned}$$

AC-Termination Non-faulty *RMs* eventually decide.

$$\begin{aligned} AC_Termination &\triangleq \wedge \forall t \in TID : \Diamond \Box (Outcome(vHist, t) \in \{\text{“Commit”}, \text{“Abort”}\}) \\ &\quad \wedge \forall t \in TID : \forall r \in part[t] : \Diamond \Box (t \in terminatedAt[r]) \end{aligned}$$

THEOREM $Spec \Rightarrow \wedge AC_Validity \wedge AC_Non_Triviality \wedge AC_Termination$

A.3 Correctness

We want to prove that the log service’s specification (LSS) we gave at Section 3 satisfies the atomic commitment and the R-Consistency properties, recalled below. We first prove an invariant of the algorithm and then proceed to prove each property individually.

- **AC-Validity** If an *RM* decides to commit a transaction, then all *RMs* voted to commit the transaction.
- **AC-Agreement** It is impossible for one *RM* to commit a transaction and another one to abort the transaction.

- **AC-Non-Triviality** If all RMs vote to commit the transaction and no RM is suspected throughout the execution of the protocol, then the decision is commit.
- **AC-Termination** Non-faulty RMs eventually decide.
- **R-Consistency** The database state of an RM after the recovery from a failure is the same as its committed state before the failure.

Invariant 1 *At any point in time there is at most one process incarnating any given resource manager.*

PROOF: At the initial state, defined by *Init*, no resource manager is incarnated by any process. Resource managers are incarnated only by executing action *Incarnate*, that replaces the previous resource manager by the new one. Therefore, at most one process can be incarnating a resource manager at any point in time, and this relation is kept in *rm2pid*. \square

Invariant 2 *Given a transaction t , if $Outcome(vHist, t) \neq \text{“Undefined”}$ at some point in time, then it will equal $Outcome(vHist, t)$ at any later time.*

PROOF: The variable *vHist* is only changed in action *Vote*, where a vote v is added to *vHist* only if there was no other vote for the same transaction $v.tr$ and resource manager $v.rm$ in *vHist*. Therefore, the only change allowed to *vHist* is the addition of new votes.

At the initial state $vHist = \{\}$ and, by the definition of action *Outcome*, $Outcome(\{\}, t)$ equals “Undefined”. From this state, “Abort” and “Commit” votes for t can be added to *vHist*.

By the definition of *Outcome*, if an “Abort” is ever added, then $Outcome(vHist, t)$ will equal “Abort” at any future evaluation. If a “Commit” vote is added for each resource manager involved in t , then $Outcome(vHist, t)$ will equal “Commit” and, because no “Abort” vote for t could be added later, it will equal “Commit” at any future evaluation. \square

Proposition A.1 (AC-Validity) *LSS satisfies the AC-Validity property.*

PROOF: By the specification, the only action that issues “Commit” votes is *VoteForMyself*. As its name suggests, the action is executed for a resource manager r only to cast its own vote and, therefore, a vote v such that $v.vt = \text{“Commit”}$ and $v.rm = r$ can only be casted by r itself.

By specification, $Outcome(vHist, t)$ will equal “Commit” only if a “Commit” vote has been issued for t from all resource managers involved. By the previous paragraph, we know that each resource manager involved in t must have issued its own “Commit” vote for t , and the **AC-Validity** property is true. \square

Proposition A.2 (AC-Agreement) *LSS satisfies the AC-Validity property.*

PROOF: Suppose that $Outcome(vHist, t)$ evaluated to “Commit” at some point in time for some transaction t . By the initial state it was initially evaluated to “Undefined” and turned to “Commit” at some later state. By the Invariant A.2, it must have turned directly from “Undefined” to “Commit”, and will be “Commit” on any future evaluation. Therefore, any resource manager either sees “Commit” or “Undefined”; because a resource manager will give t for terminated only if $Outcome(vHist, t) \neq \text{“Undefined”}$, all resource managers will see the same termination outcome. Changing “Commit” for “Abort” renders the equivalent result, and therefore the **AC-Agreement** property is true. \square

Proposition A.3 (AC-Non-Triviality) *LSS satisfies the AC-Non-Triviality property.*

PROOF: By the specification, for any transaction t , a resource manager will only vote “Abort” abort on behalf of another resource manager if it suspects that it is crashed. If there are no suspicions, then all resource manager will vote for themselves. If no resource manager votes “Abort” for itself, then only “Commit” votes will be issued and added to $vHist$. Because, by assumption, all crashes are eventually suspected, the lack of suspicions implies that no resource manager crashed. Therefore, eventually all resource managers involved in t will have their votes added to $vHist$ and, at this point in time, $Outcome(vHist, t)$ will turn from “Undefined” to “Commit”, satisfying the **AC-Non-Triviality** property. \square

To prove the next property we must assume some kind of fairness on the system; we assume the following: actions that become enabled and remain in such state until executed, are eventually executed. This property is equivalent to the Weak Fairness described in [10].

Proposition A.4 (AC-Termination) *LSS satisfies the AC-Termination property.*

Assuming that any transaction t that is started will eventually be requested to terminate, resource managers will eventually crash or vote for themselves. As, by assumption, non-faulty resource managers eventually suspect any crashed ones and vote on their behalf if they have issued their own votes, eventually some “Abort” vote or all “Commit” votes for t will be gathered by the log service, t will be terminated, and resource managers that did not crash will learn the transaction’s outcome, hence, satisfying the **AC-Termination** property. \square

Proposition A.5 (R-Consistency) *LSS satisfies the R-Consistency property.*

Let s be the state of some given database. We denote by $s.u$ the state obtained by an update u to the database at state s . To prove the the **R-Consistency** property we use the following assumptions regarding this change of states. We call two transactions “non-concurrent” if their termination procedure was executed within non-overlapping periods of time.

Assumption 1 *Applying an update is a deterministic operation. That is, given two databases at states s_1 and s_2 and update u , $(s_1 = s_2) \Rightarrow (s_1.u = s_2.u)$.*

Assumption 2 *Updates to different data items are commutable. That is, given two updates u_1 and u_2 and a database in some state s , if u_1 and u_2 do not write to the same data item, then $s.u_1.u_2 = s.u_2.u_1$.*

Assumption 3 *Concurrent transactions do not access the same data items. I.e., if two transactions execute their commit procedure in parallel, then they do not read items written by each other.*

PROOF SKETCH: We divide the proof in several steps. First we show that any non-concurrent transactions that committed, who possibly accessed the same data items, are totally ordered in $tHist$, according to their termination order. Then we show that the reincarnation procedure apply the updates of committed transactions in the same order they were committed (the order in $tHist$) and, finally, show that this ensures that a recovered resource manager has the same committed stated as before crashing or being replaced.

1. ASSUME: (C, \leq_C) is the poset where C is the set of committed transactions of some run of the system, \leq_C their commit order, and
 $vHist = (H, \leq_H)$.

PROVE: $\forall t_1, t_2 \in C : t_1 \leq_C t_2 \Rightarrow t_1 \leq_H t_2$

PROOF: Let t_1 and t_2 be two non-concurrent transactions that committed and, without loss of generality, let t_1 be the one to be committed first. By the specification, all resource managers involved in t_1 executed the action *Vote* with “Commit” vote for t_1 before any has voted for t_2 .

By the definition of *Vote*, when the first vote for t_1 was issued, t_1 was added to *ConcSet*. After the last vote, t_1 was removed from *ConcSet* and added to *vHist*, and *LastConcSet* was changed to a set not containing t_2 . When t_2 receives its first vote to commit, it is added to *ConcSet*, and upon the last vote, it is removed from *ConcSet* and added to *vHist*. Because t_2 cannot belong to the *LastConcSet* defined when t_1 was committed, t_2 will belong to a set in *tHist* different from the one t_1 belongs.

Recalling the meaning of the data-structure *tHist*, if a transaction t_1 belongs to the set *tHist*[i] for some natural number i , then $t_1 \leq_H t_2$, for any transaction t_2 that belongs to *tHist*[j] for any $j > i$.

2. ASSUME: r is a resource manager,
 t is a transaction that committed, and
 r was involved in t .

PROVE: If *Updates*(r) is evaluated after t committed, then
 $\exists i \in 1..Len(Updates(r)) : Updates(r)[i] = u$,
 where u are the updates executed by r on transaction t .

PROOF: Because t committed, by the definition of *Vote*, there must be a set in *tHist* to whom t belongs. By de definition of *Updates* and *IsInvolved*, r will be identified as participant of t and, by the definition of *UpdateOn*, inside the definition of *Updates*, r 's updates on t are in *Updates*(r).

3. ASSUME: r is a resource manager,
 t_1 and t_2 are non-concurrent transactions,
 t_1 and t_2 committed,
 t_1 committed before t_2 ,
 r was involved in t_1 and t_2 ,
 r 's updates on t_1 and t_2 are u_1 and u_2 , respectively.

PROVE: $(Updates(r)[i] = u_1) \wedge (Updates(r)[j] = u_2) \Rightarrow i < j$

PROOF: The definition of *Updates* constructs *Updates*(r) accessing *tHist* backwards, but orderly. Because each update is added to before the previous one in the sequence, the final result is that updates are in the same order as their respective transaction. By the step 2, both u_1 and u_2 are in *Updates*(r) and, by step 1, t_1 is ordered before t_2 in *tHist*, and u_1 appear before u_2 in *Updates*(r).

4. Q.E.D.

PROOF: By the definition of *ApplyUpdates*, updates are applied sequentially. Consequently, by steps 1,2, and 3, all updates of non-concurrent transactions are applied in the recovering resource manager in the same order they were originally executed. By Assumption 3, concurrent transactions access only different items and are, by Assumption 2, commutable. Finally, by Assumption 3, the recovering resource manager must have the same committed as in the previous incarnation. \square

B Coordinated Implementation

B.1 Specification

MODULE *CoordLogService*

This module is specifies the log service's Coordinated Implementation.

EXTENDS *Naturals*, *FiniteSets*, *Sequences*,
LogServiceConstants,
Consensus

CONSTANTS *Coord*

The environment's variables:

badProc processes that crashed
suspect process X process suspicions.
msgs messages sent.

VARIABLES *badProc*,
suspect,
msgs

The Coordinators's variables:

vHistAt votes received per coordinator.
tHistAt transactions that committed.
recSet set of awaiting recovery transactions.
bSet sets of votes to be proposed.
instances instance to be used by the coordinator(s).

VARIABLES *vHist*,
tHist,
recSet,
bSet,
instances

The Resource Managers' variables:

terminatingAt transactions terminating at rm.
terminatedAt transactions terminated at rm.
pid2rm PID \mapsto RM.
incarns the incarnation of each rm.
outcome local view of Outcome.

VARIABLES *terminatingAt*,
terminatedAt,
pid2rm,
incarns,
outcome

The Transaction Managers's variables:

termReq transactions requested to terminate.
part participants on each transaction, and incarnating process.

VARIABLES *termReq*,
part

Defining some aliases.

$$\begin{aligned}
svars &\triangleq \langle vHist, tHist, recSet, bSet, instances \rangle \\
rvars &\triangleq \langle terminatingAt, terminatedAt, pid2rm, incarns, outcome \rangle \\
tvars &\triangleq \langle part, termReq \rangle \\
evars &\triangleq \langle badProc, suspect \rangle \\
ovars &\triangleq \langle decision \rangle \\
avars &\triangleq \langle svars, rvars, tvars, evars, ovars, msgs \rangle
\end{aligned}$$

Types

$$Incarnating \triangleq \text{CHOOSE } p : p \notin (PID \cup \{NoRM\})$$

Vote extended with *PID*

$$\begin{aligned}
EVotes &\triangleq [rm : RM, tr : TID, tset : \text{SUBSET } RM, \\
&\quad vt : \{\text{"Commit"}, \text{"Abort"}\}, upd : Update, pid : PID]
\end{aligned}$$

Votes to incarnate a resource manager.

$$IncarnT \triangleq [vt : \{\text{"Incarnate"}\}, pid : PID, rm : RM]$$

Types of messages exchanged.

$$\begin{aligned}
Msgs &\triangleq [type : \{\text{"Recover"}\}, rm : RM, pid : PID] \cup \\
&\quad [type : \{\text{"Recovered"}\}, rm : RM, pid : PID, \\
&\quad \quad upd : Seq(Update), inc : Nat] \cup \\
&\quad [type : \{\text{"Incarnated"}\}, rm : RM, inc : Nat] \cup \\
&\quad [type : \{\text{"Vote"}\}, rm : RM, pid : PID, tr : TID, \\
&\quad \quad tset : \text{SUBSET } RM, vt : \{\text{"Commit"}, \text{"Abort"}\}, upd : Update] \cup \\
&\quad [type : \{\text{"Terminated"}\}, tr : TID, out : \{\text{"Commit"}, \text{"Abort"}\}]
\end{aligned}$$

Invariants

Type invariant.

$$\begin{aligned}
TypeInvariant &\triangleq \\
&\quad \wedge vHist \in \text{SUBSET } EVotes \\
&\quad \wedge tHist \in Seq(TID \cup IncarnT) \\
&\quad \wedge bSet \in [Coord \rightarrow \text{SUBSET } (EVotes \cup IncarnT)] \\
&\quad \wedge recSet \in \text{SUBSET } IncarnT \\
&\quad \wedge instances \in Nat \\
&\quad \wedge terminatingAt \in [PID \rightarrow \text{SUBSET } TID] \\
&\quad \wedge terminatedAt \in [PID \rightarrow \text{SUBSET } TID] \\
&\quad \wedge pid2rm \in [PID \rightarrow RM \cup \{Incarnating, NoRM\}] \\
&\quad \wedge outcome \in [PID \rightarrow [TID \rightarrow \{\text{"Undefined"}, \text{"Abort"}, \text{"Commit"}\}]] \\
&\quad \wedge part \in [TID \rightarrow \text{UNION } \{[S \rightarrow PID] : S \in \text{SUBSET } RM\}] \\
&\quad \wedge termReq \in \text{SUBSET } TID \\
&\quad \wedge badProc \in \text{SUBSET } PID \\
&\quad \wedge suspect \in [PID \rightarrow [PID \rightarrow \text{BOOLEAN}]] \\
&\quad \wedge msgs \in \text{SUBSET } Msgs
\end{aligned}$$

$\wedge \text{ConsensusTypeInv}$

Initial State

$\text{Init} \triangleq$

$\wedge vHist$	$= \{\}$	No vote received.
$\wedge tHist$	$= \langle \rangle$	No transaction committed.
$\wedge bSet$	$= [c \in Coord \mapsto \{\}]$	
$\wedge recSet$	$= [c \in Coord \mapsto \{\}]$	
$\wedge instances$	$= 0$	
$\wedge terminatingAt$	$= [p \in PID \mapsto \{\}]$	No transactions terminating.
$\wedge terminatedAt$	$= [p \in PID \mapsto \{\}]$	No transactions terminating.
$\wedge pid2rm$	$= [r \in PID \mapsto NoRM]$	No RM is incarnated.
$\wedge outcome$	$= [p \in PID \mapsto [t \in TID \mapsto \text{"Undefined"}]]$	
$\wedge incarns$	$= [p \in PID \mapsto 0]$	
$\wedge part$	$= [t \in TID \mapsto [e \in \{\} \mapsto \{\}]]$	No transaction started.
$\wedge termReq$	$= \{\}$	No termination request issued.
$\wedge badProc$	$= \{\}$	No bad process.
$\wedge suspect$	$= [p \in PID \mapsto [q \in PID \mapsto FALSE]]$	
$\wedge msgs$	$= \{\}$	
$\wedge \text{ConsensusInit}$		

Operators

The pid of the currently incarnating r .

$rm2pid(r) \triangleq$

```

IF  $\exists p \in PID, i \in 1 .. Len(tHist) :$ 
   $\wedge tHist[i] = [vt \mapsto \text{"Incarnate"}, pid \mapsto p, rm \mapsto r]$ 
   $\wedge \neg \exists j \in i + 1 .. Len(tHist), op \in PID :$ 
     $tHist[j] = [vt \mapsto \text{"Incarnate"}, pid \mapsto op, rm \mapsto r]$ 
THEN CHOOSE  $p \in PID :$ 
   $\exists i \in 1 .. Len(tHist) :$ 
     $\wedge tHist[i] = [vt \mapsto \text{"Incarnate"}, pid \mapsto p, rm \mapsto r]$ 
     $\wedge \neg \exists j \in i + 1 .. Len(tHist), op \in PID :$ 
       $tHist[j] = [vt \mapsto \text{"Incarnate"}, pid \mapsto op, rm \mapsto r]$ 
ELSE  $NoPID$ 

```

Operator $Outcome(h, t)$ gives the termination status of transaction t , considering the history of votes h .

$Outcome(h, t) \triangleq$

```

IF  $\exists v \in h : v.tr = t \wedge v.vt = \text{"Abort"}$ 
THEN  $\text{"Abort"}$ 
ELSE IF  $\exists v \in h :$ 
   $\wedge v.tr = t$ 
   $\wedge \forall p \in v.tset :$ 
     $\exists vv \in h : \wedge vv.rm = p$ 

```


$$\begin{aligned} & \wedge vv.tr = t \\ & \wedge vv.vt = \text{"Commit"} \\ \text{THEN "Commit"} \\ \text{ELSE "Undefined"} \end{aligned}$$

Actions

Environment

This action crashes a good process.

$$\begin{aligned} \text{Crash} & \triangleq \\ & \wedge \exists p \in (PID \setminus badProc) : badProc' = badProc \cup \{p\} \\ & \wedge \text{UNCHANGED } \langle suspect \rangle \end{aligned}$$

This action changes the suspicion status.

$$\begin{aligned} \text{ChangeSuspicion} & \triangleq \\ & \wedge \exists p \in (PID \setminus badProc), q \in PID : \\ & \quad \wedge p \neq q \\ & \quad \wedge suspect' = [suspect \text{ EXCEPT } ![p][q] = \neg @] \\ & \wedge \text{UNCHANGED } \langle badProc \rangle \end{aligned}$$

EnvActions:

- process crash.
- change suspicions.

$$\begin{aligned} \text{EnvActions} & \triangleq \wedge \text{Crash} \vee \text{ChangeSuspicion} \\ & \wedge \text{UNCHANGED } \langle msgs, svars, tvars, rvars, ovars \rangle \end{aligned}$$

Transaction Manager

Adds a resource manager to a non-terminated transaction.

The transaction manager will only succeed in the adding a resource manager r to a transaction t if r was not crashed by the time it was contacted, and replied to operation request.

$$\begin{aligned} \text{AddRM} & \triangleq \\ & \wedge \exists t \in (TID \setminus termReq), & t \text{ has not tried to terminate yet.} \\ & r \in \{r \in RM : \wedge rm2pid(r) \neq NoPID & r \text{ has been incarnated,} \\ & \quad \wedge rm2pid(r) \notin badProc\} : & \text{replied to an operation,} \\ & \wedge r \notin \text{DOMAIN } part[t] & \text{and was not a participant yet} \\ & \wedge part' = [part \text{ EXCEPT } ![t] = \\ & \quad [e \in \text{DOMAIN } @ \cup \{r\} \mapsto \text{IF } e \in \text{DOMAIN } @ \\ & \quad \quad \text{THEN } @[e] \\ & \quad \quad \text{ELSE } rm2pid(r)]] \text{ now it is.} \\ & \wedge \text{UNCHANGED } \langle termReq \rangle \end{aligned}$$

Request the termination of a transaction.

The transaction manager can, at any time, try to terminate a transaction it has not tried to terminate before, and that executed some operation.

$$\begin{aligned}
RequestTerm \triangleq & \wedge \exists t \in (TID \setminus termReq) : && t \text{ has not tried to terminate yet.} \\
& \wedge part[t] \neq \langle \rangle \\
& \wedge termReq' = termReq \cup \{t\} \\
& \wedge UNCHANGED \langle part \rangle
\end{aligned}$$

The disjunction of Transaction Manager's actions.

- Add a resource manager as a participant.
- Try to terminate a transaction.

$$\begin{aligned}
TMActions \triangleq & \wedge \vee AddRM \\
& \vee RequestTerm \\
& \wedge UNCHANGED \langle rvars, evars, svars, msgs, ovars \rangle
\end{aligned}$$

Log Service

Starts the incarnate procedure.

$$\begin{aligned}
IncarnateStart(p, r) \triangleq & \\
& \wedge pid2rm[p] = NoRM && p \text{ is neither incarnating nor trying} \\
& \wedge \vee rm2pid(r) = NoPID && \text{and } r \text{ is not incarnated} \\
& \vee rm2pid(r) \neq NoPID \wedge suspect[p][rm2pid(r)] && \text{or its process is suspected.} \\
& \wedge pid2rm' = [pid2rm \text{ EXCEPT } ![p] = Incarnating] \\
& \wedge msgs' = msgs \cup \{[type \mapsto \text{"Recover"}, pid \mapsto p, rm \mapsto r]\} \\
& \wedge UNCHANGED \langle incarns \rangle
\end{aligned}$$

Ends the incarnate procedure.

$$\begin{aligned}
IncarnateEnd(p, r) \triangleq & \\
& LET ApplyUpdates(u) \triangleq && \\
& \quad LET apply[i \in 1 .. Len(u)] \triangleq && \\
& \quad \quad IF i = Len(u) THEN ApplyUpdate(u[i]) && \\
& \quad \quad \quad ELSE ApplyUpdate(u[i]) \wedge apply[i + 1] && \\
& IN \quad IF u = \langle \rangle THEN TRUE && \\
& \quad \quad ELSE apply[1] && \\
& IN \quad \wedge pid2rm[p] = Incarnating && \\
& \quad \wedge \exists m \in msgs : && \\
& \quad \quad \wedge m.type = \text{"Recovered"} && \\
& \quad \quad \wedge m.pid = p && \\
& \quad \quad \wedge m.rm = r && \\
& \quad \quad \wedge pid2rm' = [pid2rm \text{ EXCEPT } ![p] = r] && \\
& \quad \quad \wedge incarns' = [incarns \text{ EXCEPT } ![p] = m.inc] && \\
& \quad \quad \wedge ApplyUpdates(m.upd) && \\
& \quad \wedge UNCHANGED \langle msgs \rangle
\end{aligned}$$

Executed by process p to give up incarnating r , when another process incarnates it.

$$\begin{aligned}
Desincarnate(p, r) \triangleq & \\
& \wedge incarns[p] > 0 \\
& \wedge \exists m \in msgs :
\end{aligned}$$

$$\begin{aligned}
& \wedge m.type = \text{"Incarnated"} \\
& \wedge m.rm = r \\
& \wedge m.inc > incarns[p] \\
& \wedge incarns' = [incarns \text{ EXCEPT } ![p] = 0] \\
& \wedge \text{UNCHANGED } \langle msgs, pid2rm \rangle
\end{aligned}$$

IncarnateStub is a “stub” to the abstract log service *Incarnate* action.

$$\begin{aligned}
\text{IncarnateStub} & \triangleq \\
& \wedge \exists p \in PID \setminus badProc, r \in RM : \quad \boxed{p \text{ is good}} \\
& \quad \vee \text{IncarnateStart}(p, r) \\
& \quad \vee \text{IncarnateEnd}(p, r) \\
& \quad \vee \text{Desincarnate}(p, r) \\
& \wedge \text{UNCHANGED } \langle terminatingAt, terminatedAt, outcome \rangle
\end{aligned}$$

VoteForMyself is executed to vote on some transaction.

$$\begin{aligned}
\text{VoteForMyself}(r, t) & \triangleq \\
& \wedge t \in (termReq \setminus (terminatingAt[part[t][r]] \cup terminatedAt[part[t][r]])) \\
& \wedge terminatingAt' = [terminatingAt \text{ EXCEPT } ![part[t][r]] = @ \cup \{t\}] \quad \boxed{\text{Start the termination}} \\
& \wedge \vee msgs' = msgs \cup \{[type \mapsto \text{"Vote"}, pid \mapsto part[t][r], rm \mapsto r, \\
& \quad tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
& \quad vt \mapsto \text{"Commit"}, upd \mapsto GetUpdate(r, t)]\} \quad \boxed{\text{voting commit.}} \\
& \quad \vee msgs' = msgs \cup \{[type \mapsto \text{"Vote"}, pid \mapsto part[t][r], rm \mapsto r, \\
& \quad tr \mapsto t, tset \mapsto \{r\}, vt \mapsto \text{"Abort"}, upd \mapsto \{\}]\} \quad \boxed{\text{voting abort.}} \\
& \wedge \text{UNCHANGED } \langle terminatedAt, outcome \rangle
\end{aligned}$$

VoteForOthers is executed to vote for some slow participant.

$$\begin{aligned}
\text{VoteForOthers}(r, t) & \triangleq \\
& \wedge t \in terminatingAt[part[t][r]] \quad \boxed{rm \text{ tried to terminate } t} \\
& \wedge outcome[part[t][r]][t] = \text{"Undefined"} \quad \boxed{\text{but did not succeed yet.}} \\
& \wedge \neg \exists m \in msgs : \\
& \quad \wedge m.type = \text{"Terminated"} \\
& \quad \wedge m.tr = t \\
& \wedge \exists s \in \text{DOMAIN } part[t] : \\
& \quad \wedge suspect[part[t][r]][part[t][s]] \\
& \quad \wedge msgs' = msgs \cup \{[type \mapsto \text{"Vote"}, pid \mapsto part[t][s], rm \mapsto s, tr \mapsto t, \\
& \quad tset \mapsto \text{DOMAIN } part[t], vt \mapsto \text{"Abort"}, upd \mapsto \{\}]\} \quad \boxed{\text{voting abort.}} \\
& \wedge \text{UNCHANGED } \langle terminatingAt, terminatedAt, outcome \rangle
\end{aligned}$$

Learn action is performed when a new transaction has terminated.

$$\begin{aligned}
\text{Learn}(r, t) & \triangleq \\
& \wedge outcome[part[t][r]][t] = \text{"Undefined"} \quad \boxed{outcome[t] \text{ is undefined}} \\
& \wedge \exists m \in msgs : \quad \boxed{\text{but } t \text{ terminated.}} \\
& \quad \wedge m.type = \text{"Terminated"} \\
& \quad \wedge m.tr = t \\
& \quad \wedge outcome' = [outcome \text{ EXCEPT } ![part[t][r]][t] = m.out] \\
& \wedge terminatingAt' = [terminatingAt \text{ EXCEPT } ![part[t][r]] = @ \setminus \{t\}]
\end{aligned}$$

$$\wedge \text{terminatedAt}' = [\text{terminatedAt} \text{ EXCEPT } ![\text{part}[t][r]] = @ \cup \{t\}]$$

$$\wedge \text{UNCHANGED } \langle \text{msgs} \rangle$$

Action *TerminateStub* is executed by *rm* to step towards transaction *t*'s termination.

It is a “stub” to the abstract log service's Terminate action.

$$\text{TerminateStub} \triangleq$$

$$\wedge \exists r \in RM, t \in TID :$$

$$\wedge r \in \text{DOMAIN } \text{part}[t]$$

$$\wedge \text{part}[t][r] \notin \text{badProc}$$

$$\wedge \text{pid2rm}[\text{part}[t][r]] = r$$

$$\wedge \vee \text{VoteForMyself}(r, t)$$

$$\quad \vee \text{VoteForOthers}(r, t)$$

$$\quad \vee \text{Learn}(r, t)$$

$$\wedge \text{UNCHANGED } \langle \text{incarns}, \text{pid2rm} \rangle$$

r is a participant of *t*.
the process is still alive,
and sees itself as the *rm*.
first attempt to terminate *t*.
other attempts
Learn that it was decided.

The disjunction of Resource Manager's actions.

- Execute the incarnation procedure.
- Try to terminate a transaction.

$$\text{RMActions} \triangleq$$

$$\wedge \vee \text{IncarnateStub}$$

$$\quad \vee \text{TerminateStub}$$

$$\wedge \text{UNCHANGED } \langle \text{tvars}, \text{evars}, \text{ovars}, \text{svars} \rangle$$

$$\text{IncarnateRequest} \triangleq$$

$$\wedge \exists p \in PID, r \in RM, c \in \text{Coord}, m \in \text{msgs} :$$

$$\wedge m.\text{type} = \text{“Recover”} \wedge m.\text{pid} = p \wedge m.\text{rm} = r$$

$$\wedge \neg \exists i \in 1 \dots \text{Len}(t\text{Hist}) :$$

$$\quad t\text{Hist}[i] = [vt \mapsto \text{“Incarnate”}, \text{pid} \mapsto p, \text{rm} \mapsto r]$$

$$\wedge \text{recSet}' = [\text{recSet} \text{ EXCEPT } ![c] = @ \cup \{[vt \mapsto \text{“Incarnate”}, \text{pid} \mapsto p, \text{rm} \mapsto r]\}]$$

$$\wedge \text{bSet}' = [\text{bSet} \text{ EXCEPT } ![c] = @ \cup \{[vt \mapsto \text{“Incarnate”}, \text{pid} \mapsto p, \text{rm} \mapsto r]\}]$$

$$\wedge \text{UNCHANGED } \langle v\text{Hist}, t\text{Hist}, \text{instances}, \text{ovars}, \text{msgs} \rangle$$

p is not incarnating.

$$\text{IncarnateReply} \triangleq$$

$$\text{LET } \text{Urm}(tinc) \triangleq$$

$$\text{LET } \text{try}[i \in 0 \dots tinc] \triangleq$$

$$\text{IF } i = 0$$

$$\text{THEN } \langle \rangle$$

$$\text{ELSE } \text{try}[i - 1] \circ \text{IF } \wedge t\text{Hist}[i] \in TID$$

$$\wedge \exists v \in v\text{Hist} :$$

$$\wedge v.\text{tr} = t\text{Hist}[i]$$

$$\wedge v.\text{rm} = t\text{Hist}[tinc].\text{rm}$$

$$\text{THEN } \langle (\text{CHOOSE } v \in v\text{Hist} :$$

$$\wedge v.\text{tr} = t\text{Hist}[i]$$

$$\wedge v.\text{rm} = t\text{Hist}[tinc].\text{rm}).\text{upd} \rangle$$

$$\text{ELSE } \langle \rangle$$

Not reincarn transaction,

is committed before *tinc*,

r took part in it.

$$\begin{aligned}
& \text{IN } \text{try}[tinc] \\
\text{Inc}(tinc) & \triangleq \text{Cardinality}(\{i \in 1 \dots tinc : \\
& \quad \wedge tHist[i] \notin TID \\
& \quad \wedge tHist[i].rm = tHist[tinc].rm\}) \\
\text{IN } \exists tinc \in 1 \dots \text{Len}(tHist), c \in \text{Coord} : \\
& \quad \wedge tHist[tinc] \in \text{recSet}[c] \\
& \quad \wedge \text{recSet}' = [\text{recSet EXCEPT } ![c] = @ \setminus \{tHist[tinc]\}] \\
& \quad \wedge \text{msgs}' = \text{msgs} \cup \{[type \mapsto \text{"Recovered"}, rm \mapsto tHist[tinc].rm, upd \mapsto \text{Urm}(tinc), \\
& \quad \quad \quad inc \mapsto \text{Inc}(tinc), pid \mapsto tHist[tinc].pid], \\
& \quad \quad \quad [type \mapsto \text{"Incarnated"}, rm \mapsto tHist[tinc].rm, inc \mapsto \text{Inc}(tinc)]\} \\
& \quad \wedge \text{UNCHANGED} \langle vHist, tHist, bSet, instances, ovars \rangle \\
\text{VoteRequest} & \triangleq \\
& \quad \exists c \in \text{Coord}, m \in \text{msgs} : \\
& \quad \wedge m.type = \text{"Vote"} \\
& \quad \wedge \neg \exists ev \in (bSet[c] \cap \text{EVotes}) \cup vHist : \\
& \quad \quad \wedge ev.rm = m.rm \\
& \quad \quad \wedge ev.tr = m.tr \\
& \quad \wedge bSet' = [bSet EXCEPT ![c] = @ \cup \{[pid \mapsto m.pid, rm \mapsto m.rm, tr \mapsto m.tr, \\
& \quad \quad \quad tset \mapsto m.tset, vt \mapsto m.vt, upd \mapsto m.upd]\}] \\
& \quad \wedge \text{UNCHANGED} \langle vHist, tHist, recSet, instances, msgs, ovars \rangle \\
\text{CoordPropose} & \triangleq \\
& \quad \wedge \exists c \in \text{Coord} : \\
& \quad \quad \wedge bSet[c] \neq \{\} \\
& \quad \quad \wedge \text{Propose}(instances, bSet[c]) \\
& \quad \wedge \text{UNCHANGED} \langle svars, msgs \rangle \\
\text{CoordDecide} & \triangleq \\
& \quad \text{LET } D \triangleq \text{Decide}(instances) \\
& \quad \text{EVotesInD} \triangleq \{v \in D : \wedge v.vt \in \{\text{"Commit"}, \text{"Abort"}\} \\
& \quad \quad \wedge \neg \exists ov \in vHist : \wedge ov.rm = v.rm \\
& \quad \quad \quad \wedge ov.tr = v.tr\} \\
& \quad \text{V2V}(v) \triangleq [rm \mapsto v.rm, tr \mapsto v.tr, tset \mapsto v.tset, pid \mapsto v.pid, \\
& \quad \quad \quad vt \mapsto \text{IF } rm2pid(v.rm) = v.pid \text{ THEN } v.vt \text{ ELSE } \text{"Abort"}, \\
& \quad \quad \quad upd \mapsto \text{IF } rm2pid(v.rm) = v.pid \text{ THEN } v.upd \text{ ELSE } \{\}] \\
& \quad \text{VotesInD} \triangleq \{\text{V2V}(v) : v \in \text{EVotesInD}\} \\
& \quad \text{Set2Seq}(S) \triangleq \\
& \quad \quad \text{LET } set2seq[SS \in \text{SUBSET } S] \triangleq \\
& \quad \quad \quad \text{IF } SS = \{\} \text{ THEN } \langle \rangle \\
& \quad \quad \quad \text{ELSE LET } ss \triangleq \text{CHOOSE } ss \in SS : \text{TRUE} \\
& \quad \quad \quad \text{IN } \text{Append}(set2seq[SS \setminus \{ss\}], ss)
\end{aligned}$$

IN $set2seq[S]$

$$newCommitted \triangleq \{t \in TID : \wedge Outcome(vHist, t) = \text{"Undefined"} \\ \wedge Outcome(vHist', t) = \text{"Commit"}\}$$

$$newTermMsgs \triangleq \{[type \mapsto \text{"Terminated"}, tr \mapsto t, out \mapsto \text{"Commit"}] : \\ t \in newCommitted\} \\ \cup \\ \{[type \mapsto \text{"Terminated"}, tr \mapsto v.tr, out \mapsto \text{"Abort"}] : \\ v \in \{vv \in VotesInD : vv.vt = \text{"Abort"}\}\}$$

$$IncarnationReqInD \triangleq \{i \in D : i.vt = \text{"Incarnate"}\}$$

IN $\wedge D \neq NoProposal$
 $\wedge vHist' = vHist \cup VotesInD$
 $\wedge msgs' = msgs \cup newTermMsgs$
 $\wedge tHist' = tHist \circ (\text{IF } newCommitted \neq \{\} \text{ THEN } Set2Seq(newCommitted) \text{ ELSE } \langle \rangle)$
 $\quad \circ Set2Seq(IncarnationReqInD)$
 $\wedge instances' = instances + 1$
 $\wedge bSet' = [c \in Coord \mapsto bSet[c] \setminus D]$
 $\wedge \text{UNCHANGED } \langle recSet, ovars \rangle$

The disjunction of Coordinator's actions.

- Process requests to incarnate a resource manager.
- Process votes from resource managers.
- Handle consensus instances.

$$CoordActions \triangleq \\ \wedge \vee IncarnateRequest \vee IncarnateReply \\ \vee VoteRequest \\ \vee CoordPropose \vee CoordDecide \\ \wedge \text{UNCHANGED } \langle tvars, evars, rvars \rangle$$

Specification

The next-state action, as a disjunction of all possible action.

$$Next \triangleq \\ \vee RMActions \vee CoordActions \quad \text{Implement } RMActions \\ \vee TMActions \quad \text{Implement } TMActions \\ \vee EnvActions \quad \text{Implement } EnvActions$$

The specification.

$$Spec \triangleq Init \wedge \square [Next]_{\langle ovars \rangle}$$

Refinement Mapping

$$\begin{aligned}
rm_pid2rm &\triangleq [r \in RM \mapsto \text{IF } pid2rm[p] = \text{"Incarnating"} \text{ THEN } NoRM \text{ ELSE } pid2rm[p]] \\
rm_rm2pid &\triangleq [r \in RM \mapsto rm2pid(r)] \\
rm_LastConcSet &\triangleq \{\} \\
rm_vHist &\triangleq \{[f \in \text{DOMAIN } v \setminus \{\text{"pid"}\} \mapsto v.f] : v \in vHist\} \\
rm_tHist &\triangleq \text{LET } Test(e) \triangleq e \notin IncarnT \\
&\quad \text{IN } SelectSeq(vHist, Test)
\end{aligned}$$

B.2 Implementation Proof

To prove that the Coordinated Log Service (CLS) is, indeed, an implementation of the Log Service's specification (LS) we give a refinement mapping of the CLS's variables to the LS's, and show that the execution of CLS's actions implies the execution of one of LS's actions, or in a stuttering step. For a thorough explanation of refinement mappings the reader is referred to [8] and [1].

We substitute every expression of the specification for an overlined expression with the same name, meaning that any variable defined in its scope is replaced by an overlined one, with the same name; these overlined variables witness the implementation of the specification. We prove that these witnesses exist by defining them from the variables in the implementation, i.e., by giving a refinement mapping.

The actual refinement is defined at the end of the specification. Below we simply rename each definition to conform the overlined notation. Variables that are not redefined are the same as in the implementation.

$$\begin{aligned}
\overline{pid2rm} &\triangleq rm_pid2rm \\
\overline{rm2pid} &\triangleq rm_rm2pid \\
\overline{LastConcSet} &\triangleq rm_LastConcSet \\
\overline{vHist} &\triangleq rm_vHist \\
\overline{tHist} &\triangleq rm_tHist
\end{aligned}$$

Proposition B.1 $Spec \Rightarrow \overline{Spec}$

1. ASSUME: \overline{Init}

PROVE: \overline{Init}

PROOF: Except for $\overline{LastConcSet}$ and $\overline{rm2pid}$, all the variables are initialized in \overline{Init} exactly as their overlined counterparts in \overline{Init} . By the refinement mapping, $\overline{LastConcSet}$ is always the empty set, therefore conforming the initialization in \overline{Init} . Finally, by the definition of operator $\overline{rm2pid}$, $\overline{tHist} = \langle \rangle$ implies that $\overline{rm2pid}$ maps from all resource managers to $NoPid$.

2. ASSUME: \overline{Next}

PROVE: $\overline{Next} \vee \text{UNCHANGED } \langle \overline{svars}, \overline{tvars}, \overline{rvars}, \overline{evars} \rangle$

2.1. ASSUME: $\overline{RMActions}$

PROVE: $\overline{RMActions}$

2.1.1. ASSUME: $\overline{IncarnateStub} \wedge \text{UNCHANGED} \langle tvars, evars \rangle$
 PROVE: $\diamond \wedge \overline{Incarnate} \vee \text{UNCHANGED} \langle svars, rvars \rangle$
 $\wedge \text{UNCHANGED} \langle \overline{tvars}, \overline{evars} \rangle$

PROOF SKETCH: We want to show that the execution of $IncarnateStart(p, r)$ for some process p and resource manager r leads to the execution $IncarnateEnd(p, r)$, if p does not crash and its messages are lost. Because $IncarnateEnd(p, r)$ can only be executed if $IncarnateStart(p, r)$ was previously executed and because the pre-conditions of $\overline{Incarnate}$ shared with $IncarnateStart$ do not change until $IncarnateEnd$ is executed, and the conditions of $\overline{Incarnate}$ shared with $IncarnateEnd$ complement the set of $\overline{Incarnate}$ pre and post-conditions already true, the execution of $IncarnateEnd$ implies an $\overline{Incarnate}$ step. If just the start action is performed, then it implies a stuttering step of $Spec$.

PROOF: The $IncarnateStub$ action is a disjunction of actions

- $IncarnateStart$,
- $IncarnateEnd$, and
- $Desincarnate$

It is clear by the specification that action $IncarnateEnd(p, r)$ cannot execute for a resource manager r and process p before an $IncarnateStart(p, r)$ is executed: $IncarnateEnd(p, r)$ only executes after receiving a message $m = [type \mapsto \langle Recovered \rangle, rm \mapsto r, pid \mapsto p]$, and such message will not be sent by action $IncarnateReply$ before a vote $v = [vt \mapsto \text{"Incarnate"}, pid \mapsto p, rm \mapsto r]$ is added to $tHist$. Hence, v will only be added to $tHist$ in action $CoordDecide$, after being proposed in action $CoordPropose$. $CoordPropose$ can only propose such value if it belongs to $bSet[c]$, for some coordinator c , what can only happen if a request message for p to incarnate r is received in action $IncarnateRequest$, and such message is only sent by the execution of action $IncarnateStart(p, r)$.

When $IncarnateStart(p, r)$ is executed, it adds a "Recover" message with fields $pid = p$ and $rm = r$ to $msgs$; this is the first pre-condition for $IncarnateRequest(p, r)$ to execute. The second pre-condition is satisfied until action $CoordDecide$ adds the "Incarnate" vote for p and r to $tHist$, what can only happen if $IncarnateRequest$ has been executed first, since $IncarnateRequest$ is the only action that creates "Incarnate" votes.

When the action $IncarnateRequest$ is executed, it adds an "Incarnate" vote for p and r to $bSet[c]$, for some coordinator c , making it not empty. This is the only pre-condition for the $CoordPropose$ action execute for coordinator c , and the action is eventually executed. Since coordinators insist on proposing its $bSet[c]$ until it is empty, and only removes votes from it if they are decided in some instance, c will keep proposing the "Incarnate" vote until it is decided or c crashes. Coordinators are deterministic state machines, and can be replicated at will (their state is only based on the outcomes of consensus instances) and, therefore, as long as coordinators can recover after crashes or infinitely many of them are available, some coordinator eventually completes the execution of $IncarnateRequest$ and $IncarnateReply$.

By the consensus problem definition, **C-Progress** ensures that a decision will eventually be reached on each instance (given that the minimum number of acceptors eventually stay up long enough for the instances to finish). If p crashes, then the request for incarnation is simply discarded or is decided but will be followed by another request for the same r .

When an instance containing the vote for p to incarnate r is decided, the vote is added to $tHist$, and action $IncarnateReply$ will be enabled for all coordinators that proposed it. If c crashes before this action is performed, p will be blocked and never execute another action, as if it had crashed. Because the change made by $IncarnateStart(r, p)$ to variable $pid2rm[p]$ does not affect $\overline{pid2rm}$, this would imply that $\langle \overline{svars}, \overline{rvars} \rangle$ did not change.

Once $IncarnateReply$ is performed, p will eventually receive the "Recovered" message, unless c crashes, enabling action $IncarnateEnd$. By the definition of Urm , p will receive all the

updates performed by previous incarnations of r . By the definition of *ApplyUpdates*, the p will apply all the updates and recover the committed state r had on its previous incarnation. The pre condition of *IncarnateStub*, *IncarnatedStart*, and *IncarnatedEnd*, and the post-conditions of *IncarnateEnd* imply the pre and post-conditions of $\overline{\text{Incarnate}}$.

By the assumption, variables in $\langle \overline{tvars}, \overline{evars} \rangle$ do not change.

- 2.1.2. ASSUME: $\overline{\text{TerminateStub}} \wedge \text{UNCHANGED} \langle \overline{tvars}, \overline{evars}, \overline{ovars}, \overline{svars} \rangle$
 PROVE: $\diamond \wedge \overline{\text{Terminate}} \vee \text{UNCHANGED} \langle \overline{svars}, \overline{rvars} \rangle$
 $\wedge \text{UNCHANGED} \langle \overline{tvars}, \overline{evars} \rangle$

PROOF SKETCH: The pre-conditions of action *TerminateStub* are the same as those of $\overline{\text{Terminate}}$. Therefore, it is enough to show that each of *TerminateStub*'s sub-actions, *VoteForMyself*, *VoteForOthers*, *Learn*, and the actions they lead to, imply a step of their equivalent overlined actions.

- 2.1.2.1. ASSUME: $\overline{\text{VoteForMyself}}(r, t)$
 PROVE: $\diamond \wedge \overline{\text{VoteForMyself}}(r, t) \vee \text{UNCHANGED} \langle \overline{\text{terminatingAt}}, \overline{vHist}, \overline{tHist} \rangle$
 $\wedge \text{UNCHANGED} \langle \overline{\text{terminatedAt}} \rangle$

Because the pre-condition and the first post-condition of both actions are the same, it is enough to prove that the second post-condition of $\overline{\text{VoteForMyself}}(r, t)$, the addition of a message m to \overline{msgs} , may the execution of $\overline{\text{Vote}}(v)$, where m is “Vote” message and v is a vote, and the fields $vt, upd, tr, tset$, and rm of m and v are equal, or has no effect on variables $\langle \overline{\text{terminatingAt}}, \overline{vHist}, \overline{tHist} \rangle$.

If message m is received by some coordinator c , a vote with its contents, therefore equal to v , is added to $\overline{bSet}[c]$, enabling the action *CoordPropose*. Action *CoordPropose* will be executed with a proposal containing this vote until it is decided and added to \overline{vHist} by action *CoordDecide*, where all coordinators can see it (\overline{vHist} is changed deterministically based on the consensus outcomes, and would be the same for all coordinators if represented independently at each one.), or until c crashes. If no coordinator succeeds in getting the vote decided, then either another vote, resulting from the execution of $\overline{\text{VoteForOthers}}(r, t)$ will be decided, or all resource managers involved in transaction t will have crashed before their “Vote” messages are seen by non-faulty coordinators. It is up to the transaction manager, to then vote to abort the transaction; in the case the transaction manager also crashes and no vote for t is ever decided, t is simply forgotten, implying that $\langle \overline{svars}, \overline{rvars} \rangle$ is kept unchanged.

CoordDecide also appends newly committed transactions to \overline{tHist} : each transaction is added in a different set, as if they were not concurrent, and $\overline{\text{LastConcSet}}$ is always empty, ensuring its type invariance and constructing \overline{tHist} in a way compatible with the specification of \overline{tHist} . Because this is the only action to change \overline{vHist} and \overline{tHist} , the mapping to \overline{vHist} and \overline{tHist} is correct.

$\langle \overline{tvars}, \overline{evars} \rangle$ are kept since none of actions changes them.

- 2.1.2.2. ASSUME: $\overline{\text{VoteForOthers}}(r, t)$
 PROVE: $\diamond \wedge \overline{\text{VoteForOthers}}(r, t) \vee \text{UNCHANGED} \langle \overline{vHist}, \overline{tHist} \rangle$
 $\wedge \text{UNCHANGED} \langle \overline{\text{terminatedAt}}, \overline{\text{terminatingAt}} \rangle$

PROOF: It is true by the same arguments of step 2.1.2.1.

- 2.1.2.3. ASSUME: $\overline{\text{Learn}}(r, t)$
 PROVE: $\wedge \overline{\text{Learn}}(r, t) \vee \text{UNCHANGED} \langle \overline{\text{terminatedAt}}, \overline{\text{terminatingAt}} \rangle$
 $\wedge \text{UNCHANGED} \langle \overline{svars} \rangle$

PROOF: As action $\overline{\text{Learn}}(r, t)$ has the same post-conditions of $\overline{\text{Learn}}(r, t)$, it is enough to show that the pre-conditions of the first imply the pre-conditions of the latter. Since the reception of “Terminated” message for transaction t implies that it was sent and since it is only

sent by action *CoordDecide* if the transaction has terminated, the reception of such message implies that the $\overline{Outcome(vHist, t)} \neq \text{"Undefined"}$.

2.1.2.4. Q.E.D.

2.1.3. Q.E.D.

2.2. $TMActions \Rightarrow \overline{TMActions}$

2.2.1. $AddRM \Rightarrow \overline{AddRM}$

PROOF: Trivially true, since the their definitions are equal.

2.2.2. $RequestTerm \Rightarrow \overline{RequestTerm}$

PROOF: Trivially true, since the their definitions are equal.

2.2.3. $UNCHANGED \langle rvars, evars, svars, msgs, ovars \rangle \Rightarrow UNCHANGED \langle \overline{svars}, \overline{rvars}, \overline{evars} \rangle$

PROOF: Clearly true since either the left-hand side of the expression contains all variables in the spec.

2.2.4. Q.E.D.

2.3. $EnvActions \Rightarrow \overline{EnvActions}$

PROOF: Trivially true, since their definitions are equal.

2.4. Q.E.D.

3. Q.E.D.

□

C Uncoordinated Implementation

C.1 Specification

MODULE *UnCoordLogService*

This module specifies the log service's Uncoordinated Implementation.

EXTENDS *Naturals*, *FiniteSets*, *Sequences*,
LogServiceConstants,
Consensus

The implementation's constants:

MPL RM's multi-programming level.
VoteSort Function that sorts votes according to the commit counter.

CONSTANTS *MPL*, *MultiProgrammingLevel*
VoteSort(-)

The environment's variables:

badProc processes that crashed
suspect process X process suspicions.

VARIABLES *badProc*, *Environments's*
suspect *Who suspects whom?*

The Resource Managers and processes that might incarnate resource managers' variables:

terminatedAt transactions terminated at rm.
terminatingAt transactions terminating at a process.
pid2rm PID \mapsto RM.
instances PID \mapsto consensus instance to be used.
countdown number of consensus instances yet to be closed.
incarnation status of the incarnation procedure.

VARIABLES *terminatedAt*,
terminatingAt,
pid2rm,
instances,
countdown,
incarnation

The Transaction Managers's variables:

termReq transactions requested to terminate.
part participants on each transaction, and incarnating process.
inst RM X TID \mapsto IID.
nextInst RM \mapsto the next IID to be attributed to a transaction.

VARIABLES *termReq*,
part,
inst,
nextInst

Defining some aliases.

rvars \triangleq \langle *terminatingAt*, *terminatedAt*, *pid2rm*, *instances*,

$$\begin{aligned}
& \text{countdown, incarnation} \rangle \\
tvars & \triangleq \langle part, termReq, nextInst, inst \rangle \\
evars & \triangleq \langle badProc, suspect \rangle \\
ovars & \triangleq \langle decision \rangle \\
avars & \triangleq \langle rvars, tvars, evars, ovars \rangle
\end{aligned}$$

Types

Vote is extended with commit counter.

$$\begin{aligned}
EVotes & \triangleq [rm : RM, tr : TID \cup \{NoTID\}, tset : \text{SUBSET } RM, \\
& vt : \{\text{"Commit"}, \text{"Abort"}\}, upd : Update, cn : Nat]
\end{aligned}$$

Votes to incarnate a resource manager.

$$\text{IncarnT} \triangleq [vt : \{\text{"Incarnate"}\}, pid : PID]$$

Invariants

Type invariant.

$$\begin{aligned}
\text{TypeInvariant} & \triangleq \\
& \wedge \text{terminatingAt} \in [PID \rightarrow \text{SUBSET } TID] \\
& \wedge \text{terminatedAt} \in [RM \rightarrow \text{SUBSET } TID] \\
& \wedge \text{pid2rm} \in [PID \rightarrow RM \cup \{NoRM\}] \\
& \wedge \text{instances} \in [PID \rightarrow Nat] \\
& \wedge \text{incarnation} \in [PID \rightarrow \{\text{"NotInc"}, \\
& \quad \text{"Inc1"}, \text{"Inc2"}, \text{"Inc3"}, \text{"Inc4"}, \text{"Inc"}, \\
& \quad \text{"DelInc"}\}] \\
& \wedge \text{countdown} \in [PID \rightarrow Nat] \\
& \wedge \text{part} \in [TID \rightarrow \text{UNION } \{[S \rightarrow PID] : S \in \text{SUBSET } RM\}] \\
& \wedge \text{termReq} \in \text{SUBSET } TID \\
& \wedge \text{nextInst} \in [RM \rightarrow Nat] \\
& \wedge \text{inst} \in [RM \rightarrow [TID \rightarrow Nat \cup \{NoIID\}]] \\
& \wedge \text{badProc} \in \text{SUBSET } PID \\
& \wedge \text{suspect} \in [PID \rightarrow [PID \rightarrow \text{BOOLEAN}]] \\
& \wedge \text{ConsensusTypeInv}
\end{aligned}$$

Initial State

$$\begin{aligned}
\text{Init} & \triangleq \\
& \wedge \text{terminatingAt} = [p \in PID \mapsto \{\}] && \text{No transaction terminating.} \\
& \wedge \text{terminatedAt} = [r \in RM \mapsto \{\}] && \text{No transaction terminated.} \\
& \wedge \text{pid2rm} = [p \in PID \mapsto NoRM] && \text{No } RM \text{ is incarnated.} \\
& \wedge \text{instances} = [p \in PID \mapsto 0] && \text{No consensus instance started.} \\
& \wedge \text{incarnation} = [p \in PID \mapsto \text{"NotInc"}] && \text{No process is recovering.} \\
& \wedge \text{countdown} = [p \in PID \mapsto MPL - 1] && \text{No process is recovering.} \\
& \wedge \text{part} = [t \in TID \mapsto [e \in \{\} \mapsto \{\}]] && \text{No transaction started.}
\end{aligned}$$

$\wedge termReq$	$= \{\}$	No termination request issued.
$\wedge nextInst$	$= [r \in RM \mapsto 0]$	No operation executed.
$\wedge inst$	$= [r \in RM \mapsto [t \in TID \mapsto NoIID]]$	No instance used.
$\wedge badProc$	$= \{\}$	No bad process.
$\wedge suspect$	$= [p \in PID \mapsto [q \in PID \mapsto FALSE]]$	No suspicion.
$\wedge ConsensusInit$		

Operators

The *pid* of the process currently incarnating *r*.

$rm2pid(r) \triangleq$

IF $\exists i \in Nat :$

$\wedge Decide(\langle r, i \rangle) \neq NoProposal$

If there is a vote

$\wedge Decide(\langle r, i \rangle).vt = \text{"Incarnate"}$

to incarnate *r*,

THEN CHOOSE $p \in PID :$

choose the process *p*

$\exists i \in Nat :$

$\wedge Decide(\langle r, i \rangle) \neq NoProposal$

$\wedge Decide(\langle r, i \rangle).vt = \text{"Incarnate"}$

who voted for last.

$\wedge Decide(\langle r, i \rangle).pid = p$

$\wedge \neg \exists j \in Nat : \wedge j > i$

$\wedge Decide(\langle r, j \rangle) \neq NoProposal$

$\wedge Decide(\langle r, j \rangle).vt = \text{"Incarnate"}$

ELSE *NoPID*

Choose *NoPID* otherwise.

Operator *Outcome*(*r*, *t*) gives the termination status of transaction *t*, considering the decisions learned by *r*.

$Outcome(r, t) \triangleq$

IF $\wedge Decide(\langle r, inst[r][t] \rangle) \neq NoProposal$

If *r*'s vote

$\wedge \vee Decide(\langle r, inst[r][t] \rangle).vt = \text{"Incarnate"}$

was to change incarnation

$\vee Decide(\langle r, inst[r][t] \rangle).vt = \text{"Abort"}$

or to *Abort*

$\vee \wedge Decide(\langle r, inst[r][t] \rangle).vt = \text{"Commit"}$

or to commit but

$\wedge \exists or \in Decide(\langle r, inst[r][t] \rangle).tset :$

another participant

$\wedge Decide(\langle or, inst[or][t] \rangle) \neq NoProposal$

has a voted to

$\wedge \vee Decide(\langle or, inst[or][t] \rangle).vt = \text{"Incarnate"}$

change incarnation

$\vee Decide(\langle or, inst[or][t] \rangle).vt = \text{"Abort"}$

or to *Abort*

THEN **"Abort"**

then it is *Abort*.

ELSE IF $\vee Decide(\langle r, inst[r][t] \rangle) = NoProposal$

If *r*'s vote has not been decided

$\vee \wedge Decide(\langle r, inst[r][t] \rangle) \neq NoProposal$

or it has

$\wedge Decide(\langle r, inst[r][t] \rangle).vt = \text{"Commit"}$

and is *Commit*, but

$\wedge \exists or \in Decide(\langle r, inst[r][t] \rangle).tset :$

someone else's

$Decide(\langle or, inst[or][t] \rangle) = NoProposal$

has not.

THEN **"Undefined"**

then it is *Undefined*.

ELSE **"Commit"**

Otherwise it is *Commit*.

Actions

Environment

This action crashes a good process.

$$\begin{aligned} \text{Crash} \triangleq & \wedge \exists p \in (PID \setminus \text{badProc}) : \text{badProc}' = \text{badProc} \cup \{p\} \\ & \wedge \text{UNCHANGED} \langle \text{suspect} \rangle \end{aligned}$$

This action changes the suspicion status.

$$\begin{aligned} \text{ChangeSuspicion} \triangleq & \wedge \exists p \in (PID \setminus \text{badProc}), q \in PID : \\ & \wedge p \neq q \\ & \wedge \text{suspect}' = [\text{suspect} \text{ EXCEPT } ![p][q] = \neg@] \\ & \wedge \text{UNCHANGED} \langle \text{badProc} \rangle \end{aligned}$$

EnvActions is the disjunction of the environment's actions.

- process crash.
- change suspicions.

$$\begin{aligned} \text{EnvActions} \triangleq & \wedge \text{Crash} \vee \text{ChangeSuspicion} \\ & \wedge \text{UNCHANGED} \langle \text{tvars}, \text{rvars}, \text{ovars} \rangle \end{aligned}$$

Transaction Manager

Adds a resource manager to a non-terminated transaction.

The transaction manager will only succeed in the adding a resource manager r to a transaction t if r was not crashed by the time it was contacted, and replied to the operation request.

AddRM \triangleq

$$\text{LET } \text{processingAt}(r) \triangleq \{tr \in TID : \wedge r \in \text{DOMAIN } \text{part}[tr] \\ \wedge tr \notin \text{terminatedAt}[r]\}$$

$$\begin{aligned} \text{Smaller}(r) \triangleq & \text{LET } \text{instS} \triangleq \{\text{inst}[r][t] : t \in \text{processingAt}(r)\} \\ & \text{IN CHOOSE } i \in \text{instS} : \forall j \in \text{instS} : i \leq j \end{aligned}$$

$$\begin{aligned} \text{CanProcess}(r) \triangleq & \text{processingAt}(r) \neq \{\} \\ \Rightarrow & \text{nextInst}[r] - \text{Smaller}(r) < \text{MPL} \end{aligned}$$

$$\begin{aligned} \text{IN } & \wedge \exists t \in (TID \setminus \text{termReq}), & & t \text{ has not had termination requested.} \\ & r \in \{r \in RM : \text{rm2pid}(r) \neq \text{NoPID}\} : & & r \text{ has been incarnated,} \\ & \wedge \text{incarnation}[\text{rm2pid}(r)] = \text{"Inc"} & & \text{and is contactable,} \\ & \wedge \text{rm2pid}(r) \notin \text{badProc} & & \text{and is not a participant of } t \\ & \wedge r \notin \text{DOMAIN } \text{part}[t] & & \text{and has free slots.} \\ & \wedge \text{CanProcess}(r) & & \text{Now it is a participant.} \\ & \wedge \text{part}' = [\text{part} \text{ EXCEPT } ![t] = & & \\ & \quad [e \in \text{DOMAIN } @ \cup \{r\} \mapsto \text{IF } e \in \text{DOMAIN } @ & & \\ & \quad \quad \text{THEN } @[e] & & \\ & \quad \quad \text{ELSE } \text{rm2pid}(r)]] & & \\ & \wedge \text{inst}' = [\text{inst} \text{ EXCEPT } ![r][t] = \text{nextInst}[r]] & & \\ & \wedge \text{nextInst}' = [\text{nextInst} \text{ EXCEPT } ![r] = @ + 1] & & \\ & \wedge \text{UNCHANGED} \langle \text{termReq} \rangle \end{aligned}$$

Request the termination of a transaction.

\wedge UNCHANGED \langle terminatingAt, terminatedAt, pid2rm, instances,
countdown, nextInst \rangle

$IncEnd(p, r) \triangleq$

LET $CVrm \triangleq \{Decide(\langle r, i \rangle) :$ Votes for instances that committed.
 $i \in \{j \in 1 \dots instances[p] - 1 :$
 $\wedge Decide(\langle r, j \rangle).vt = \text{“Commit”}$
 $\wedge Outcome(r, Decide(\langle r, j \rangle).tr) = \text{“Commit”}\}$

$ApplyUpdates(u) \triangleq$

LET $apply[i \in 1 \dots Len(u)] \triangleq$
 IF $i = Len(u)$ THEN $ApplyUpdate(u[i])$
 ELSE $ApplyUpdate(u[i]) \wedge apply[i + 1]$

IN IF $u = \langle \rangle$ THEN TRUE
 ELSE $apply[1]$

IN $\wedge incarnation[p] = \text{“Inc4”}$ p is on the last step of incarnate
to incarnate r ,
 $\wedge pid2rm[p] = r$ and is moving forward;
 $\wedge incarnation' = [incarnation \text{ EXCEPT } ![p] = \text{“Inc”}]$ and setting the first instance accordingly.
 $\wedge ApplyUpdates(VoteSort(CVrm))$ it is applying the updates,
 $\wedge nextInst' = [nextInst \text{ EXCEPT } ![r] = instances[p]]$
 \wedge UNCHANGED \langle terminatingAt, terminatedAt, pid2rm, instances,
countdown, ovars \rangle

Desincarnate is executed when p sees that another process has incarnated r .

$Desincarnate(p, r) \triangleq$

$\wedge incarnation[p] = \text{“Inc”}$ p is incarnating
 $\wedge pid2rm[p] = r$ r
 $\wedge \exists i \in instances[p] \dots nextInst[r] :$ but there is a new vote
 $\wedge Decide(\langle r, i \rangle) \neq NoProposal$
 $\wedge Decide(\langle r, i \rangle).vt = \text{“Incarnate”}$ to incarnate r
 $\wedge Decide(\langle r, i \rangle).pid \neq p$ from another process.
 $\wedge incarnation' = [incarnation \text{ EXCEPT } ![p] = \text{“Delnc”}]$
 \wedge UNCHANGED \langle terminatingAt, terminatedAt, pid2rm, instances,
countdown, nextInst, ovars \rangle

IncarnateStub is the disjunction of the steps needed to incarnate or desincarnate a resource manager.

It is a “stub” to the abstract log service’s Incarnate action.

$IncarnateStub \triangleq$

$\exists p \in PID \setminus badProc,$ p is good
 $r \in RM :$
 $\vee IncStart(p, r)$ and wants to incarnate
 $\vee Inc1Step(p, r) \vee Inc1BreakOrLoop(p, r)$
 $\vee Inc2Step(p, r) \vee Inc2BreakOrLoop(p, r)$
 $\vee Inc3Step(p, r)$
 $\vee IncEnd(p, r)$
 $\vee Desincarnate(p, r)$ or desincarnate r .

Executed by r to issue its own vote for transaction t .

$$\begin{aligned}
 & \text{VoteForMyself}(r, t) \triangleq \\
 & \text{LET } \text{CommittedCounter} \triangleq \text{Cardinality}(\{i \in \text{Nat} : \text{Decide}(\langle r, i \rangle) \neq \text{NoProposal}\}) \\
 & \text{IN } \wedge t \in (\text{termReq} \setminus (\text{terminatingAt}[part[t][r]] \cup \text{terminatedAt}[r])) \\
 & \quad \wedge \text{terminatingAt}' = [\text{terminatingAt} \text{ EXCEPT } ![part[t][r]] = @ \cup \{t\}] \quad \text{Try to terminate} \\
 & \quad \wedge \vee \text{Propose}(\langle r, inst[r][t] \rangle, [rm \mapsto r, tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
 & \quad \quad \quad vt \mapsto \text{"Commit"}, upd \mapsto \text{GetUpdate}(r, t), \quad \text{voting commit.} \\
 & \quad \quad \quad cn \mapsto \text{CommittedCounter}]) \\
 & \quad \vee \text{Propose}(\langle r, inst[r][t] \rangle, [rm \mapsto r, tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
 & \quad \quad \quad vt \mapsto \text{"Abort"}, upd \mapsto \{\}, cn \mapsto 0]) \quad \text{voting abort.} \\
 & \wedge \text{UNCHANGED } \langle \text{terminatedAt}, \text{instances}, \text{countdown}, \text{incarnation} \rangle
 \end{aligned}$$

Executed by r to issue votes on behalf of other resource managers.

$$\begin{aligned}
 & \text{VoteForOthers}(r, t) \triangleq \\
 & \quad \wedge t \in \text{terminatingAt}[part[t][r]] \quad r \text{ tried to terminate } t \\
 & \quad \wedge \text{Outcome}(r, t) = \text{"Undefined"} \quad \text{but did not succeed yet} \\
 & \quad \wedge \exists or \in \text{DOMAIN } part[t] : \quad \text{because some participant did not vote} \\
 & \quad \wedge \text{suspect}[part[t][r]][part[t][or]] \quad \setminus * \text{ and I don't want to wait} \\
 & \quad \quad \wedge \text{Propose}(\langle or, inst[or][t] \rangle, [rm \mapsto or, tr \mapsto t, tset \mapsto \text{DOMAIN } part[t], \\
 & \quad \quad \quad vt \mapsto \text{"Abort"}, upd \mapsto \{\}, cn \mapsto 0]) \quad \text{so vote Abort on its behalf.} \\
 & \wedge \text{UNCHANGED } \langle \text{terminatingAt}, \text{terminatedAt}, \text{instances}, \text{countdown}, \text{incarnation} \rangle
 \end{aligned}$$

Executed by r to learn that transaction t terminated.

$$\begin{aligned}
 & \text{Learn}(r, t) \triangleq \\
 & \quad \wedge \text{Outcome}(r, t) \neq \text{"Undefined"} \quad \text{The instance is finished} \\
 & \quad \wedge t \in \text{terminatingAt}[part[t][r]] \setminus \text{terminatedAt}[r] \quad \text{but } r \text{ did not know} \\
 & \quad \wedge \text{terminatingAt}' = [\text{terminatingAt} \text{ EXCEPT } ![part[t][r]] = @ \setminus \{t\}] \quad \text{so take note.} \\
 & \quad \wedge \text{terminatedAt}' = [\text{terminatedAt} \text{ EXCEPT } ![r] = @ \cup \{t\}] \\
 & \quad \wedge \text{UNCHANGED } \langle \text{instances}, \text{countdown}, \text{incarnation}, \text{ovars} \rangle
 \end{aligned}$$

TerminateStub is executed by a resource manager to step towards the termination of some transaction.

It is a “stub” to the abstract log service’s *Terminate* action.

$$\begin{aligned}
 & \text{TerminateStub} \triangleq \\
 & \quad \wedge \exists r \in \text{RM}, t \in \text{TID} : \\
 & \quad \quad \wedge r \in \text{DOMAIN } part[t] \quad r \text{ is a participant of } t. \\
 & \quad \quad \wedge part[t][r] \notin \text{badProc} \quad \text{The process is still alive.} \\
 & \quad \quad \wedge incarnation[part[t][r]] = \text{"Inc"} \quad \text{and sees itself as } rm \text{'s incarnation.} \\
 & \quad \quad \wedge \vee \text{VoteForMyself}(r, t) \quad \text{first attempt to terminate } t. \\
 & \quad \quad \quad \vee \text{VoteForOthers}(r, t) \quad \text{other attempts} \\
 & \quad \quad \quad \vee \text{Learn}(r, t) \quad \text{Learn that it was decided.} \\
 & \quad \wedge \text{UNCHANGED } \langle pid2rm, nextInst \rangle
 \end{aligned}$$

The disjunction of resource managers actions.

$$\begin{aligned}
 & \text{RMActions} \triangleq \wedge \text{IncarnateStub} \vee \text{TerminateStub} \\
 & \quad \wedge \text{UNCHANGED } \langle part, \text{termReq}, inst, \text{evars} \rangle
 \end{aligned}$$

Specification

$$\begin{aligned}
 \text{Next} &\triangleq \bigvee \text{RMActions} && \text{Implement RMActions} \\
 &\bigvee \text{TMActions} && \text{Implement TMActions} \\
 &\bigvee \text{EnvActions} && \text{Implement EnvActions} \\
 \text{Spec} &\triangleq \text{Init} \wedge \square[\text{Next}]_{\langle \text{avars} \rangle}
 \end{aligned}$$

Refinement Mapping

$$\text{rm_pid2rm} \triangleq [p \in \text{PID} \mapsto \text{IF incarnation}[p] = \text{"Inc"} \text{ THEN pid2rm}[p] \text{ ELSE NoRM}]$$

$$\text{rm_rm2pid} \triangleq [r \in \text{RM} \mapsto \text{rm2pid}(r)]$$

$$\text{rm_LastConcSet} \triangleq$$

$$\begin{aligned}
 &\text{LET LastToCommitAt}(r) \triangleq \\
 &\quad \text{CHOOSE } t \in \text{TID} : && \text{Last to commit as some resource manager.} \\
 &\quad \quad \wedge \text{Outcome}(r, t) = \text{"Commit"} \\
 &\quad \quad \wedge \neg \exists ot \in \text{TID} : \\
 &\quad \quad \quad \wedge \text{Outcome}(r, ot) = \text{"Commit"} \\
 &\quad \quad \quad \wedge \text{Decide}(\langle r, \text{inst}[r][ot] \rangle).cn > \text{Decide}(\langle r, \text{inst}[r][t] \rangle).cn
 \end{aligned}$$

$$\begin{aligned}
 \text{LastToCommit} &\triangleq && \text{Last to commit on each resource managers.} \\
 &\{ \text{LastToCommitAt}(r) : r \in \{ \\
 &\quad \quad \quad \text{or} \in \text{RM} : \\
 &\quad \quad \quad \exists t \in \text{TID} : \\
 &\quad \quad \quad \quad \text{Outcome}(r, t) = \text{"Commit"} \} \}
 \end{aligned}$$

$$\begin{aligned}
 \text{DistLastToCommit} &\triangleq \\
 &\{ t \in \text{LastToCommit} : && \text{Last to commit globally.} \\
 &\quad \neg \exists r \in \text{RM} : \\
 &\quad \quad \exists ot \in \text{TID} : \\
 &\quad \quad \quad \wedge \text{Outcome}(r, ot) = \text{"Commit"} \\
 &\quad \quad \quad \wedge \text{Decide}(\langle r, \text{inst}[r][ot] \rangle).cn > \text{Decide}(\langle r, \text{inst}[r][t] \rangle).cn \\
 &\quad \}
 \end{aligned}$$

$$\begin{aligned}
 \text{ConcToLast} &\triangleq \{ ot \in \text{TID} : && \text{Concurrent to the last to commit.} \\
 &\quad \exists r \in \text{RM}, t \in \text{DistLastToCommit} : \\
 &\quad \quad \wedge \text{Decide}(r, \text{inst}[r][ot]).cn = \text{Decide}(r, \text{inst}[r][t]).cn \}
 \end{aligned}$$

IN *ConcToLast*

$$\text{rm_terminatedAt} \triangleq$$

$$\begin{aligned}
 &\text{LET IncInstOf}(p) \triangleq \\
 &\quad \text{CHOOSE } i \in \text{IID} : \wedge i[1] = \text{pid2rm}[p] && \text{The instance that incarnated } p \\
 &\quad \quad \wedge \text{Decide}(i).vt = \text{"Incarnate"} \\
 &\quad \quad \wedge \text{Decide}(i).pid = p
 \end{aligned}$$

$$\begin{aligned}
DeIncInstOf(p) &\triangleq \\
&\text{CHOOSE } i \in IID : \wedge i[1] = pid2rm[p] \quad \text{The instance that desincarnated } p \\
&\quad \wedge i[2] > IncInstOf(p)[2] \\
&\quad \wedge Decide(i).vt = \text{"Incarnate"} \\
&\quad \wedge \neg \exists j \in IID : \wedge j[2] > IncInstOf(p)[2] \\
&\quad \quad \wedge j[2] < i[2] \\
&\quad \quad \wedge Decide(j).vt = \text{"Incarnate"} \\
FirstInstOf(p) &\triangleq \langle pid2rm[p], IncInstOf(p)[2] + MPL \rangle \quad p\text{'s first instance.} \\
LastInstOf(p) &\triangleq \langle pid2rm[p], DeIncInstOf(p)[2] + MPL - 1 \rangle \quad p\text{'s last instance.} \\
TermTransOf(p) &\triangleq \\
&\{ t : \wedge Outcome(pid2rm[p], t) = \text{"Commit"} \\
&\quad \wedge inst[pid2rm[p]][t][2] > FirstInstOf(p)[2] \\
&\quad \wedge \exists i \in IID : \wedge i[1] = pid2rm[p] \quad \text{The instance that desincarnated } p \\
&\quad \wedge i[2] > IncInstOf(p)[2] \\
&\quad \wedge Decide(i).vt = \text{"Incarnate"} \\
&\quad \Rightarrow inst[pid2rm[p]][t][2] \leq LastInstOf(p)[2] \\
&\} \\
IN \ [p \in PID \mapsto \text{IF } pid2rm[p] = NoRM \text{ THEN } \{\} \text{ ELSE } TermTransOf(p)]
\end{aligned}$$

$$\begin{aligned}
rm_vHist &\triangleq \\
&\text{LET } TheVote(i) \triangleq \\
&\quad \text{IF } Decide(i).vt \in \{\text{"Abort"}, \text{"Incarnate"}\} \\
&\quad \text{THEN } [rm \mapsto i.[i], tr \mapsto \text{CHOOSE } t \in TID : inst[i[1]][t] = i, \\
&\quad \quad tset \mapsto \text{DOMAIN } part[t], vt \mapsto \text{"Abort"}, upd \mapsto \{\}] \\
&\quad \text{ELSE LET } v \triangleq Decide(i) \\
&\quad \quad IN \ [rm \mapsto v.rm, tr \mapsto v.tr, tset \mapsto v.tset, vt \mapsto v.vt, upd \mapsto v.upd] \\
IN \ \{TheVote(i) : i \in \{oi \in IID : \wedge \exists t \in TID : inst[i[1]][t] \neq NoIID \\
\quad \wedge Decide(i) \neq NoProposal\}\}
\end{aligned}$$

$$\begin{aligned}
rm_tHist &\triangleq \text{CHOOSE } s \in Seq(\{t \in TID : \exists r \in RM : Outcome(r, ot) = \text{"Commit"}\}) : \\
&\quad \forall i, j \in \text{DOMAIN } s : \\
&\quad \quad \forall r \in (\text{DOMAIN } part[i] \cap \text{DOMAIN } part[j]) : \\
&\quad \quad \quad Decide(\langle r, inst[r][i] \rangle).cn < Decide(\langle r, inst[r][j] \rangle).cn \\
&\quad \quad \Rightarrow i < j
\end{aligned}$$

C.2 Implementation Proof

The refinement mapping is given below.

$$\overline{pid2rm} \triangleq rm_pid2rm$$

$$\overline{rm2pid} \triangleq rm_rm2pid$$

$$\overline{LastConcSet} \triangleq rm_LastConcSet$$

$$\overline{terminatedAt} \triangleq rm_terminatedAt$$

$$\overline{vHist} \triangleq rm_vHist$$

$$\overline{tHist} \triangleq rm_tHist$$

Proposition C.1 $Spec \Rightarrow \overline{Spec}$

1. ASSUME: \overline{Init}

PROVE: \overline{Init}

PROOF: Variables $\overline{terminatingAt}$, \overline{part} , $\overline{termReq}$, $\overline{badProc}$ and $\overline{suspect}$ are initialized in \overline{Init} exactly as specified in \overline{Init} . By the refinement mapping, $\overline{LastConcSet}$ is initially the empty set and $\overline{terminatedAt}$ maps each process to the empty set, conforming with their initialization in \overline{Init} . By the definition of operator $rm2pid$ and because no value has neither been proposed nor decided, $\overline{rm2pid}$ maps from all resource managers to $NoPID$, as specified in \overline{Init} . By the refinement mapping and because no has been decided at the initial state, $\overline{vHist} = \{\}$ and $\overline{tHist} = \langle \rangle$, according to \overline{Init} .

2. ASSUME: \overline{Next}

PROVE: $\overline{Next} \vee \text{UNCHANGED} \langle \overline{svars}, \overline{tvars}, \overline{rvars}, \overline{evars} \rangle$

2.1. ASSUME: $\overline{RMActions}$

PROVE: $\overline{RMActions}$

2.1.1. ASSUME: $\overline{IncarnateStub} \wedge \text{UNCHANGED} \langle \overline{part}, \overline{termReq}, \overline{inst}, \overline{evars} \rangle$

PROVE: $\diamond \wedge \overline{Incarnate} \vee \text{UNCHANGED} \langle \overline{svars}, \overline{rvars} \rangle$

$\wedge \text{UNCHANGED} \langle \overline{tvars}, \overline{evars} \rangle$

PROOF: The $\overline{Incarnate}$ action is implemented by a series of steps, implemented by actions

- $\overline{IncStart}(p, r)$,
- $\overline{Inc1Step}(p, r)$,
- $\overline{Inc1BreakOrLoop}(p, r)$,
- $\overline{Inc2Step}(p, r)$,
- $\overline{Inc2BreakOrLoop}(p, r)$,
- $\overline{Inc3Step}(p, r)$,
- $\overline{IncEnd}(p, r)$, and
- $\overline{Desincarnate}(p, r)$,

where p is a process and r the resource manager p is trying to incarnate. These actions execute in a specific order, according to the state of the variable $\overline{incarnation}[p]$: $\overline{IncStart}(p, r)$ requires $\overline{incarnation}[p] = \text{"NotInc"}$, and changes it to "Inc1". $\overline{Inc1Step}(p, r)$ requires $\overline{incarnation}[p] = \text{"Inc1"}$, and does not change it, possibly being executed until it changes to "Inc2". Action $\overline{Inc1BreakOrLoop}(p, r)$ also has $\overline{incarnation}[p] = \text{"Inc1"}$ as a pre-condition, but changes it to "Inc2" if some conditions are met. Action $\overline{Inc2Step}(p, r)$ requires $\overline{incarnation}[p] = \text{"Inc2"}$, and executes until it turns to "Inc3". $\overline{Inc2BreakOrLoop}(p, r)$ has $\overline{incarnation}[p] = \text{"Inc2"}$ as pre-condition and changes $\overline{incarnating}[p]$ to "Inc3". The change of $\overline{incarnating}[p]$ from "Inc3" to "Inc4" is done by action $\overline{Inc3Step}(p, r)$ and, finally, action $\overline{IncEnd}(p, r)$ requires $\overline{incarnation}[p] = \text{"Inc4"}$ and changes it to "Inc", when the incarnation

procedure is over. Once the first step is finished, the process q that previously incarnated r either crashes or eventually learns about p , by the termination property of consensus, and executes action $Desincarnate(q, r)$.

The initial steps in this series serves the purpose of determining the updates executed by previous incarnations of r , ensuring that the consensus instances possibly used by previous incarnations of r are terminated. Because none of the actions, except for the last, changes any of the overlined variables, their execution implies in stuttering steps of \overline{Spec} . When the last action, $IncEnd(p, r)$, is executed, an equivalent $\overline{Incarnate}$ step is executed;

The specification does not ensure that $Inc1BreakOrLoop$ will ever succeed in incrementing the variable $instances[p]$, but once it is done, the process incarnating r will crash or eventually learn that it was replaced, since there is a maximum number of consensus instances it can start before learning the outcome of the previous ones (MPL), by the **C-Termination** property of consensus. If the process previously incarnating r is not crashed, then it may proceed executing transactions, and p will eventually learn that it made a mistake and crash itself (although this behavior is not specified, it is clearly correct, as p is allowed to crash at any moment). For the same reason, once $incarnation[p]$ equals “Inc2”, p will crash or eventually execute the other actions of $IncarnateStub$, learning the updates of all transactions executed by previous incarnations of r . This is performed by simply enforcing the termination of the consensus instances used on those incarnations and by the other resource managers involved in such transactions. Due to the agreement property of consensus, the outcome of each committed transaction is seen by any process that executes the incarnation procedure.

The order p applies these updates is determined by the $VoteSort$ operator, that orders updates according to the cn field contained in their respective votes. Because non-concurrent transactions are guaranteed to have different cn , ordering the updates according to cn ensures that updates of non-concurrent transactions are ordered in the order they committed.

2.1.2. ASSUME: $\overline{TerminateStub} \wedge \text{UNCHANGED} \langle part, termReq, inst, evars \rangle$

PROVE: $\diamond \wedge \overline{Terminate} \vee \text{UNCHANGED} \langle svars, rvars \rangle$
 $\wedge \text{UNCHANGED} \langle tvars, evars \rangle$

PROOF: The first two pre-conditions are the same for both actions. The third pre-condition of $\overline{TerminateStub}$ is true iff t is in the range of transactions in which r was incarnated by $part[t][r]$ and the action $Desincarnate$ has not been executed by the process. If $Desincarnate$ is not enabled, then the incarnating process cannot have executed action $IncEnd$ and the $part[t][r]$ is still incarnating r . Otherwise, although this pre-condition is true, the action will not succeed adding a vote to $vHist$, since the consensus instance it would use has been already decided by the second step of the incarnation of process of the process starting to incarnate r .

The action $VoteForMyself(r, t)$ has equivalent pre-conditions in both specifications and similar first post-condition. By the consensus properties, the second post-condition ensures that r 's vote for t is added to \overline{vHist} , if $part[t][r]$ is not suspected of having crashed, as specified in \overline{Vote} . Otherwise it implies a stuttering step of \overline{Spec} . By the definition of \overline{vHist} , deciding on a vote implies in adding a vote to \overline{vHist} . If the vote leads to the “Commit” of a transaction, then $\overline{LastConcSet}$ becomes the set of transactions for which votes have been issued with the same commit counter as those for the transaction just committed.

As the case for $VoteForMyself$ and $\overline{VoteForMyself}$, $VoteForOthers$ has conditions that are equivalent to those in $\overline{VoteForOthers}$. By the same argumentation as in the previous paragraph, the execution of a $Propose$ ensures that either a vote for the other resource manager will be added to $vHist$, if $part[t][r]$ does not crash, or it will imply in a stuttering step of \overline{Spec} .

A *Learn* step clearly implies a \overline{Learn} step, since all conditions for the latter are also stated for the first.

2.1.3. Q.E.D.

2.2. $TMActions \Rightarrow \overline{TMActions}$

2.2.1. $AddRM \Rightarrow \overline{AddRM}$

PROOF: Trivially true, since all of \overline{AddRM} pre and post-conditions are also conditions for *AddRM*.

2.2.2. $RequestTerm \Rightarrow \overline{RequestTerm}$

PROOF: Trivially true, since their definitions are equal.

2.3. $EnvActions \Rightarrow \overline{EnvActions}$

PROOF: Trivially true, since their definitions are equal.

2.4. Q.E.D.

3. Q.E.D.

□