



# Paxos Consensus, Deconstructed and Abstracted

Álvaro García-Pérez<sup>1</sup>  , Alexey Gotsman<sup>1</sup>, Yuri Meshman<sup>1</sup>,  
and Ilya Sergey<sup>2</sup> 

<sup>1</sup> IMDEA Software Institute, Madrid, Spain

{alvaro.garcia.perez,alexey.gotsman,yuri.meshman}@imdea.org

<sup>2</sup> University College London, London, UK

i.sergey@ucl.ac.uk

**Abstract.** Lamport’s Paxos algorithm is a classic consensus protocol for state machine replication in environments that admit crash failures. Many versions of Paxos exploit the protocol’s intrinsic properties for the sake of gaining better run-time performance, thus widening the gap between the original description of the algorithm, which was proven correct, and its real-world implementations. In this work, we address the challenge of specifying and verifying complex Paxos-based systems by (a) devising composable specifications for implementations of Paxos’s single-decree version, and (b) engineering disciplines to reason about protocol-aware, semantics-preserving optimisations to single-decree Paxos. In a nutshell, our approach elaborates on the deconstruction of single-decree Paxos by Boichat et al. We provide novel non-deterministic specifications for each module in the deconstruction and prove that the implementations refine the corresponding specifications, such that the proofs of the modules that remain unchanged can be reused across different implementations. We further reuse this result and show how to obtain a verified implementation of Multi-Paxos from a verified implementation of single-decree Paxos, by a series of novel protocol-aware transformations of the network semantics, which we prove to be behaviour-preserving.

## 1 Introduction

Consensus algorithms are an essential component of the modern fault-tolerant deterministic services implemented as message-passing distributed systems. In such systems, each of the distributed nodes contains a replica of the system’s state (*e.g.*, a database to be accessed by the system’s clients), and certain nodes may propose values for the next state of the system (*e.g.*, requesting an update in the database). Since any node can crash at any moment, all the replicas have to keep copies of the state that are consistent with each other. To achieve this, at each update to the system, all the non-crashed nodes run an instance of a *consensus protocol*, uniformly deciding on its outcome. The safety requirements for consensus can be thus stated as follows: “only a single value is decided uniformly by all non-crashed nodes, it never changes in the future, and the decided value has been proposed by some node participating in the protocol” [16].

The Paxos algorithm [15,16] is the classic consensus protocol, and its single-decree version (SD-Paxos for short) allows a set of distributed nodes to reach an agreement on the outcome of a *single* update. Optimisations and modifications to SD-Paxos are common. For instance, the multi-decree version, often called Multi-Paxos [15,27], considers multiple slots (*i.e.*, multiple positioned updates) and decides upon a result for *each* slot, by running a slot-specific instance of an SD-Paxos. Even though it is customary to think of Multi-Paxos as of a series of independent SD-Paxos instances, in reality the implementation features multiple protocol-aware optimisations, exploiting intrinsic dependencies between separate single-decree consensus instances to achieve better throughput. To a great extent, these and other optimisations to the algorithm are pervasive, and verifying a modified version usually requires to devise a new protocol definition and a proof from scratch. New versions are constantly springing (*cf.* Sect. 5 of [27] for a comprehensive survey) widening the gap between the description of the algorithms and their real-world implementations.

We tackle the challenge of *specifying* and *verifying* these distributed algorithms by contributing two verification techniques for consensus protocols.

Our first contribution is a family of composable specifications for Paxos' core subroutines. Our starting point is the deconstruction of SD-Paxos by Boichat *et al.* [2,3], allowing one to consider a distributed consensus instance as a *shared-memory concurrent program*. We introduce novel specifications for Boichat *et al.*'s modules, and let them be non-deterministic. This might seem as an unorthodox design choice, as it *weakens* the specification. To show that our specifications are still *strong enough*, we restore the top-level *deterministic* abstract specification of the consensus, which is convenient for client-side reasoning. The weakness introduced by the non-determinism in the specifications has been impelled by the need to prove that the implementations of Paxos' components *refine* the specifications we have ascribed [9]. We prove the refinements modularly via the Rely/Guarantee reasoning with prophecy variables and explicit linearisation points [11,26]. On the other hand, this weakness becomes a virtue when better understanding the volatile nature of Boichat *et al.*'s abstractions and of the Paxos algorithm, which may lead to newer modifications and optimisations.

Our second contribution is a methodology for verifying composite consensus protocols by reusing the proofs of their constituents, targeting specifically Multi-Paxos. We distill protocol-aware system optimisations into a separate semantic layer and show how to obtain the realistic Multi-Paxos implementation from SD-Paxos by a *series of transformations* to the *network semantics* of the system, as long as these transformations preserve the behaviour observed by clients. We then provide a family of such transformations along with the formal conditions allowing one to compose them in a behaviour-preserving way.

We validate our approach for construction of modularly verified consensus protocols by providing an executable proof-of-concept implementation of Multi-Paxos with a high-level shared memory-like interface, obtained via a series of behaviour-preserving network transformations. The full proofs of lemmas and

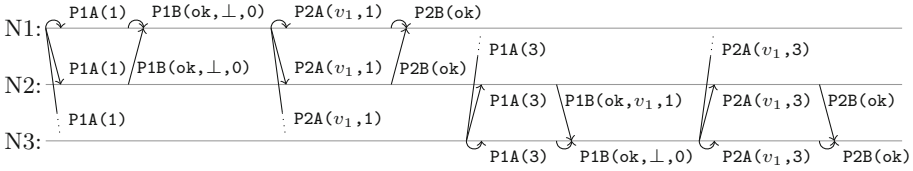


Fig. 1. A run of SD-Paxos.

theorems from our development, as well as some boilerplate definitions, are given in the appendices of the supplementary extended version of this paper.<sup>1</sup>

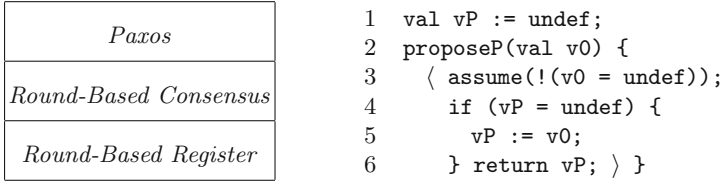
## 2 The Single-Decree Paxos Algorithm

We start with explaining SD-Paxos through an intuitive scenario. In SD-Paxos, each node in the system can adopt the roles of *proposer* or *acceptor*, or both. A value is decided when a *quorum* (i.e., a majority of acceptors) accepts the value proposed by some proposer. Now consider a system with three nodes N1, N2 and N3, where N1 and N3 are both proposers and acceptors, and N2 is an acceptor, and assume N1 and N3 propose values  $v_1$  and  $v_3$ , respectively.

The algorithm works in two phases. In Phase 1, a proposer polls every acceptor in the system and tries to convince a quorum to promise that they will later accept its value. If the proposer succeeds in Phase 1 then it moves to Phase 2, where it requests the acceptors to fulfil their promises in order to get its value decided. In our example, it would seem in principle possible that N1 and N3 could respectively convince two different quorums—one consisting of N1 and N2, and the other consisting of N2 and N3—to go through both phases and to respectively accept their values. This would happen if the communication between N1 and N3 gets lost and if N2 successively grants the promise and accepts the value of N1, and then does the same with N3. This scenario breaks the safety requirements for consensus because both  $v_1$  and  $v_3$ —which can be different—would get decided. However, this cannot happen. Let us explain why.

The way SD-Paxos enforces the safety requirements is by distinguishing each attempt to decide a value with a unique *round*, where the rounds are totally ordered. Each acceptor stores its current round, initially the least one, and only grants a promise to proposers with a round greater or equal than its current round, at which moment the acceptor switches to the proposer’s round. Figure 1 depicts a possible run of the algorithm. Assume that rounds are natural numbers, that the acceptors’ current rounds are initially 0, and that the nodes N1 and N3 attempt to decide their values with rounds 1 and 3 respectively. In Phase 1, N1 tries to convince a quorum to switch their current round to 1 (messages P1A(1)). The message to N3 gets lost and the quorum consisting of N1 and N2 switches round and promises to only accept values at a round greater or

<sup>1</sup> Find the extended version online at <https://arxiv.org/abs/1802.05969>.



**Fig. 2.** Deconstruction of SD-Paxos (left) and specification of module *Paxos* (right).

equal than 1. Each acceptor that switches to the proposer’s round sends back to the proposer its stored value and the round at which this value was accepted, or an undefined value if the acceptor never accepted any value yet (messages P1B(ok, ⊥, 0), where ⊥ denotes a default undefined value). After Phase 1, N1 picks as a candidate value the one accepted at the greatest round from those returned by the acceptors in the quorum, or its proposed value if all acceptors returned an undefined value. In our case, N1 picks its value  $v_1$ . In Phase 2, N1 requests the acceptors to accept the candidate value  $v_1$  at round 1 (messages P2A( $v_1, 1$ )). The message to N3 gets lost, and N1 and N2 accept value  $v_1$ , which gets decided (messages P2B(ok)).

Now N3 goes through Phase 1 with round 3 (messages P1A(3)). Both N2 and N3 switch to round 3. N2 answers N3 with its stored value  $v_1$  and with the round 1 at which  $v_1$  was accepted (message P1B(ok,  $v_1, 1$ )), and N3 answers itself with an undefined value, as it has never accepted any value yet (message P1B(ok, ⊥, 0)). This way, if some value has been already decided upon, *any* proposer that convinces a quorum to switch to its round would receive the decided value from some of the acceptors in the quorum (recall that two quorums have a non-empty intersection). That is, N3 picks the  $v_1$  returned by N2 as the candidate value, and in Phase 2 it manages that the quorum N2 and N3 accepts  $v_1$  at round 3 (messages P2A( $v_1, 3$ ) and P2B(ok)). N3 succeeds in making a new decision, but the decided value remains the same, and, therefore, the safety requirements of a consensus protocol are satisfied.

### 3 The Faithful Deconstruction of SD-Paxos

We now recall the faithful deconstruction of SD-Paxos in [2,3], which we take as the reference architecture for the implementations that we aim to verify. We later show how each module of the deconstruction can be verified separately.

The deconstruction is depicted on the left of Fig. 2, which consists of modules *Paxos*, *Round-Based Consensus* and *Round-Based Register*. These modules correspond to the ones in Fig. 4 of [2], with the exception of *Weak Leader Election*. We assume that a correct process that is trusted by every other correct process always exists, and omit the details of the leader election. Leaders take the role of proposers and invoke the interface of *Paxos*. Each module uses the interface provided by the module below it.

```

1 read(int k) {
2   int j; val v; int kW; val maxV;
3   int maxKW; set of int Q; msg m;
4   for (j := 1, j <= n, j++)
5   { send(j, [RE, k]); }
6   maxKW := 0; maxV := undef; Q := {};
7   do { (j, m) := receive();
8     switch (m) {
9       case [ackRE, @k, v, kW]:
10        Q := Q ∪ {j};
11        if (kW >= maxKW)
12        { maxKW := kW; maxV := v; }
13        case [nackRE, @k]:
14          return (false, _);
15      } if (|Q| = ⌈(n+1)/2⌉)
16        { return (true, maxV); } }
17  while (true); }

18 write(int k, val vW) {
19   int j; set of int Q; msg m;
20   for (j := 1, j <= n, j++)
21   { send(j, [WR, k, vW]); }
22   Q := {};
23   do { (j, m) := receive();
24     switch (m) {
25       case [ackWR, @k]:
26        Q := Q ∪ {j};
27       case [nackWR, @k]:
28        return false;
29     } if (|Q| = ⌈(n+1)/2⌉)
30       { return true; } }
31  while (true); }

```

**Fig. 3.** Implementation of *Round-Based Register* (read and write).

The entry module *Paxos* implements SD-Paxos. Its specification (right of Fig. 2) keeps a variable  $vP$  that stores the decided value (initially undefined) and provides the operation  $\text{propose}_P$  that takes a proposed value  $v0$  and returns  $vP$  if some value was already decided, or otherwise it returns  $v0$ . The code of the operation runs *atomically*, which we emphasise via angle brackets  $\langle \dots \rangle$ . We define this specification so it meets the safety requirements of a consensus, therefore, any implementation whose entry point refines this specification will have to meet the same safety requirements.

In this work we present both specifications and implementations in pseudo-code for an imperative WHILE-like language with basic arithmetic and primitive types, where  $\text{val}$  is some user-defined type for the values decided by Paxos, and  $\text{undef}$  is a literal that denotes an undefined value. The pseudo-code is self-explanatory and we restraint ourselves from giving formal semantics to it, which could be done in standard fashion if so wished [30]. At any rate, the pseudo-code is ultimately a vehicle for illustration and we stick to this informal presentation.

The implementation of the modules is depicted in Figs. 3, 4 and 5. We describe the modules following a bottom-up approach, which better fits the purpose of conveying the connection between the deconstruction and SD-Paxos. We start with module *Round-Based Register*, which offers operations  $\text{read}$  and  $\text{write}$  (Fig. 3) and implements the replicated processes that adopt the role of acceptors (Fig. 4). We adapt the wait-free, crash-stop implementation of *Round-Based Register* in Fig. 5 of [2] by adding loops for the explicit reception of each individual message and by counting acknowledgement messages one by one. Processes are identified by integers from 1 to  $n$ , where  $n$  is the number of processes in the system. Proposers and acceptors exchange read and write requests, and their corresponding acknowledgements and non/acknowledgements. We assume a type  $\text{msg}$  for messages and let the message vocabulary to be as follows.

```

1 process Acceptor(int j) {
2   val v := undef; int r := 0; int w := 0;
3   start() {
4     int i; msg m; int k;
5     do { (i, m) := receive();
6       switch (m) {
7         case [RE, k]:
8           if (k < r) { send(i, [nackRE, k]); }
9           else { < r := k; send(i, [ackRE, k, v, w]); } }
10        case [WR, k, vW]:
11          if (k < r) { send(i, [nackWR, k]); }
12          else { < r := k; w := k; v := vW; send(i, [ackWR, k]); } }
13      } }
14   while (true); } }

```

Fig. 4. Implementation of *Round-Based Register* (acceptor).

Read requests  $[RE, k]$  carry the proposer's round  $k$ . Write requests  $[WR, k, v]$  carry the proposer's round  $k$  and the proposed value  $v$ . Read acknowledgements  $[ackRE, k, v, k']$  carry the proposer's round  $k$ , the acceptor's value  $v$ , and the round  $k'$  at which  $v$  was accepted. Read non-acknowledgements  $[nackRE, k]$  carry the proposer's round  $k$ , and so do carry write acknowledgements  $[ackWR, k]$  and write non/acknowledgements  $[nackWR, K]$ .

In the pseudo-code, we use  $\_$  for a wildcard that could take any literal value. In the pattern-matching primitives, the literals specify the pattern against which an expression is being matched, and operator  $@$  turns a variable into a literal with the variable's value. Compare the case  $[ackRE, @k, v, kW]$ : in Fig. 3, where the value of  $k$  specifies the pattern and  $v$  and  $kW$  get some values assigned, with the case  $[RE, k]$ : in Fig. 4, where  $k$  gets some value assigned.

We assume the network ensures that messages are neither created, modified, deleted, nor duplicated, and that they are always delivered but with an arbitrarily large transmission delay.<sup>2</sup> Primitive `send` takes the destination  $j$  and the message  $m$ , and its effect is to send  $m$  from the current process to the process  $j$ . Primitive `receive` takes no arguments, and its effect is to receive at the current process a message  $m$  from origin  $i$ , after which it delivers the pair  $(i, m)$  of identifier and message. We assume that `send` is non-blocking and that `receive` blocks and suspends the process until a message is available, in which case the process awakens and resumes execution.

Each acceptor (Fig. 4) keeps a value  $v$ , a current round  $r$  (called the *read round*), and the round  $w$  at which the acceptor's value was last accepted (called the *write round*). Initially,  $v$  is `undef` and both  $r$  and  $w$  are 0.

Phase 1 of SD-Paxos is implemented by operation `read` on the left of Fig. 3. When a proposer issues a `read`, the operation requests each acceptor's promise to only accept values at a round greater or equal than  $k$  by sending  $[RE, k]$

<sup>2</sup> We allow creation and duplication of  $[RE, k]$  messages in Sect. 5, where we obtain Multi-Paxos from SD-Paxos by a series of transformations of the network semantics.

```

1 proposeRC(int k, val v0) {
2   bool res; val v;
3   (res, v) := read(k);
4   if (res) {
5     if (v = undef) { v := v0; }
6     res := write(k, v);
7     if (res) { return (true, v); } }
8   return (false, _); }

1 proposeP(val v0) {
2   int k; bool res; val v;
3   k := pid();
4   do { (res, v) :=
5     proposeRC(k, v0);
6     k := k + n;
7   } while (!res);
8   return v; }

```

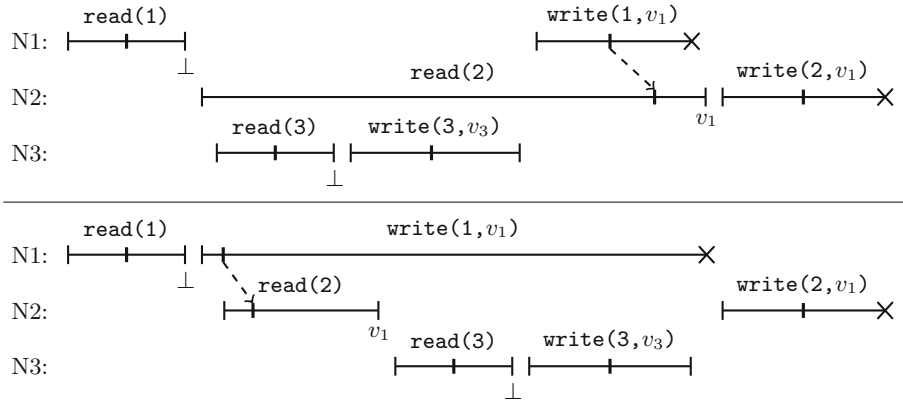
**Fig. 5.** Implementation of *Round-Based Consensus* (left) and *Paxos* (right)

(lines 4–5). When an acceptor receives a  $[\text{RE}, k]$  (lines 5–7 of Fig. 4) it acknowledges the promise depending on its read round. If  $k$  is strictly less than  $r$  then the acceptor has already made a promise to another proposer with greater round and it sends  $[\text{ackRE}, k]$  back (line 8). Otherwise, the acceptor updates  $r$  to  $k$  and acknowledges by sending  $[\text{ackRE}, k, v, w]$  (line 9). When the proposer receives an acknowledgement (lines 8–10 of Fig. 3) it counts acknowledgements up (line 10) and calculates the greatest write round at which the acceptors acknowledging so far accepted a value, and stores this value in  $\text{maxV}$  (lines 11–12). If a majority of acceptors acknowledged, the operation succeeds and returns  $(\text{true}, \text{maxV})$  (lines 15–16). Otherwise, if the proposer received some  $[\text{ackRE}, k]$  the operation fails, returning  $(\text{false}, \_)$  (lines 13–14).

Phase 2 of SD-Paxos is implemented by operation `write` on the right of Fig. 3. After having collected promises from a majority of acceptors, the proposer picks the candidate value  $vW$  and issues a `write`. The operation requests each acceptor to accept the candidate value by sending  $[\text{WR}, k, vW]$  (lines 20–21). When an acceptor receives  $[\text{WR}, k, vW]$  (line 10 of Fig. 4) it accepts the value depending on its read round. If  $k$  is strictly less than  $r$ , then the acceptor never promised to accept at such round and it sends  $[\text{ackWR}, k]$  back (line 11). Otherwise, the acceptor fulfils its promise and updates both  $w$  and  $r$  to  $k$  and assigns  $vW$  to its value  $v$ , and acknowledges by sending  $[\text{ackWR}, k]$  (line 12). Finally, when the proposer receives an acknowledgement (lines 23–25 of Fig. 3) it counts acknowledgements up (line 26) and checks whether a majority of acceptors acknowledged, in which case  $vW$  is decided and the operation succeeds and returns `true` (lines 29–30). Otherwise, if the proposer received some  $[\text{ackWR}, k]$  the operation fails and returns `false` (lines 27–28).<sup>3</sup>

Next, we describe module *Round-Based Consensus* on the left of Fig. 5. The module offers an operation `proposeRC` that takes a round  $k$  and a proposed value  $v0$ , and returns a pair  $(\text{res}, v)$  of Boolean and value, where `res` informs of the success of the operation and  $v$  is the decided value in case `res` is `true`. We have taken the implementation from Fig. 6 in [2] but adapted to our pseudocode conventions. *Round-Based Consensus* carries out Phase 1 and Phase 2 of

<sup>3</sup> For the implementation to be correct with our shared-memory-concurrency approach, the update of the data in acceptors must happen atomically with the sending of acknowledgements in lines 9 and 12 of Fig. 4.



**Fig. 6.** Two histories in which a failing `write` contaminates some acceptor.

SD-Paxos as explained in Sect. 2. The operation `proposeRC` calls `read` (line 3) and if it succeeds then chooses a candidate value between the proposed value  $v_0$  or the value  $v$  returned by `read` (line 5). Then, the operation calls `write` with the candidate value and returns `(true, v)` if `write` succeeds, or fails and returns `(false, _)` (line 8) if either the `read` or the `write` fails.

Finally, the entry module *Paxos* on the right of Fig. 5 offers an operation `proposeP` that takes a proposed value  $v_0$  and returns the decided value. We assume that the system primitive `pid()` returns the process identifier of the current process. We have come up with this straightforward implementation of operation `proposeP`, which calls `proposeRC` with increasing round until the call succeeds, starting at a round equal to the process identifier `pid()` and increasing it by the number of processes  $n$  in each iteration. This guarantees that the round used in each invocation to `proposeRC` is unique.

**The Challenge of Verifying the Deconstruction of Paxos.** Verifying each module of the deconstruction separately is cumbersome because of the distributed character of the algorithm and the nature of a linearisation proof. A process may not be aware of the information that will flow from itself to other processes, but this future information flow may dictate whether some operation has to be linearised at the present. Figure 6 illustrates this challenge.

Let N1, N2 and N3 adopt both the roles of acceptors and proposers, which propose values  $v_1$ ,  $v_2$  and  $v_3$  with rounds 1, 2 and 3 respectively. Consider the history on the top of the figure. N2 issues a read with round 2 and gets acknowledgements from all but one acceptors in a quorum. (Let us call this one acceptor A.) None of these acceptors have accepted anything yet and they all return  $\perp$  as the last accepted value at round 0. In parallel, N3 issues a read with round 3 (third line in the figure) and gets acknowledgements from a quorum in which A does not occur. This read succeeds as well and returns `(true, undef)`.



```

1 (bool × val) ptp[1..n] := undef;
2 val abs_vP := undef; single bool abs_resP[1..n] := undef;
3 proposeP(val v0) {
4   int k; bool res; val v; assume(!(v0 = undef));
5   k := pid(); ptp[pid()] := (true, v0);
6   do { (res, v) := proposeRC(k, v0);
7     if (res) {
8       for (i := 1, i <= n, i++) {
9         if (ptp[i] = (true, v)) { lin(i); ptp[i] := (false, v); } }
10        if (!(v = v0)) { lin(pid()); ptp[pid()] := (false, v0); } } }
11   k := k + n; }
12   while (!res); return v; }

```

Fig. 7. Instrumented implementation of *Paxos*.

Then N3 issues a write with round 3 and value  $v_3$ . Again, it gets acknowledgements from a quorum in which A does not occur, and the write succeeds deciding value  $v_3$  and returns **true**. Later on, and in real time order with the write by N3 but in parallel with the read by N2, node N1 issues a write with round 1 and value  $v_1$  (first line in the figure). This write is to fail because the value  $v_3$  was already decided with round 3. However, the write manages to “contaminate” acceptor A with value  $v_1$ , which now acknowledges N2 and sends  $v_1$  as its last accepted value at round 1. Now N2 has gotten acknowledgements from a quorum, and since the other acceptors in the quorum returned 0 as the round of their last accepted value, the read will catch value  $v_1$  accepted at round 1, and the operation succeeds and returns  $(\mathbf{true}, v_1)$ . This history linearises by moving N2’s read after N1’s write, and by respecting the real time order for the rest of the operations. (The linearisation ought to respect the information flow order between N1 and N2 as well, *i.e.*, N1 contaminates A with value  $v_1$ , which is read by N2.)

In the figure, a segment ending in an  $\times$  indicates that the operation fails. The value returned by a successful read operation is depicted below the end of the segment. The linearisation points are depicted with a thick vertical line, and the dashed arrow indicates that two operations are in the information flow order.

The variation of this scenario on the bottom of Fig. 6 is also possible, where N1’s write and N2’s read happen concurrently, but where N2’s read is shifted backwards to happen before in real time order with N3’s read and write. Since N1’s write happens before N2’s read in the information flow order, then N1’s write has to inexorably linearise before N3’s operations, which are the ones that will “steal” N1’s valid round.

These examples give us three important hints for designing the specifications of the modules. First, after a decision is committed it is *not enough* to store only the decided value, since a posterior write may contaminate some acceptor with a value different from the decided one. Second, a read operation *may succeed* with some round even if by that time other operation has already succeeded with a higher round. And third, a write with a valid round *may fail* if its round will be “stolen” by a concurrent operation. The non-deterministic specifications that we introduce next allow one to model execution histories as the ones in Fig. 6.

## 4 Modularly Verifying SD-Paxos

In this section, we provide non-deterministic specifications for *Round-Based Consensus* and *Round-Based Register* and show that each implementation refines its specification [9]. To do so, we instrument the implementations of all the modules with *linearisation-point* annotations and use Rely/Guarantee reasoning [26].

This time we follow a top-down order and start with the entry module *Paxos*.

**Module *Paxos*.** In order to prove that the implementation on the right of Fig. 5 refines its specification on the right of Fig. 2, we introduce the instrumented implementation in Fig. 7, which uses the helping mechanism for external linearisation points of [18]. We assume that each proposer invokes `proposeP` with a unique proposed value. The auxiliary pending thread pool `ptp[n]` is an array of pairs of Booleans and values of length  $n$ , where  $n$  is the number of processes in the system. A cell `ptp[i]` containing a pair `(true, v)` signals that the process  $i$  proposed value  $v$  and the invocation `proposeP(v)` by process  $i$  awaits to be linearised. Once this invocation is linearised, the cell `ptp[i]` is updated to the pair `(false, v)`. A cell `ptp[i]` containing `undef` signals that the process  $i$  never proposed any value yet. The array `abs_resP[n]` of Boolean single-assignment variables stores the abstract result of each proposer’s invocation. A linearisation-point annotation `lin(i)` takes a process identifier  $i$  and performs atomically the abstract operation invoked by proposer  $i$  and assigns its result to `abs_resP[i]`. The abstract state is modelled by variable `abs_vP`, which corresponds to variable `vP` in the specification on the right of Fig. 2. One invocation of `proposeP` may help linearise other invocations as follows. The linearisation point is together with the invocation to `proposeRC` (line 6). If `proposeRC` committed with some value  $v$ , the instrumented implementation traverses `ptp` and linearises all the proposers which were proposing value  $v$  (the proposer may linearise itself in this traversal) (lines 8–9). Then, the current proposer linearises itself if its proposed value  $v_0$  is different from  $v$  (line 10), and the operation returns  $v$  (line 12). All the annotations and code in lines 6–10 are executed inside an atomic block, together with the invocation to `proposeRC(k, v_0)`.

**Theorem 1.** *The implementation of Paxos on the right of Fig. 5 linearises with respect to its specification on the right of Fig. 2.*

**Module *Round-Based Consensus*.** The top of Fig. 8 shows the non-deterministic module’s specification. Global variable `vRC` is the decided value, initially `undef`. Global variable `roundRC` is the highest round at which some value was decided, initially 0; a global set of values `valsRC` (initially empty) contains values that may have been proposed by proposers. The specification is non-deterministic in that local value `vD` and Boolean `b` are unspecified, which we model by assigning random values to them. We assume that the current process identifier is  $((k - 1) \bmod n) + 1$ , which is consistent with how rounds are assigned to each process and incremented in the code of `proposeP` on the right of Fig. 5. If the unspecified value `vD` is neither in the set `valsRC` nor equal to `v_0` then the operation returns `(false, _)` (line 11). This models that the operation fails

```

1  val vRC := undef; int roundRC := 0; set of val valsRC := {};
2  proposeRC(int k, val v0) {
3    ⟨ val vD := random(); bool b := random();
4      assume(!(v0 = undef)); assume(pid() = ((k - 1) mod n) + 1);
5      if (vD ∈ (valsRC ∪ {v0})) {
6        valsRC := valsRC ∪ {vD};
7        if (b && (k ≥ roundRC)) { roundRC := k;
8                                if (vRC = undef) { vRC := vD; }
9                                return (true, vRC); }
10       else { return (false, _); } }
11     else { return (false, _); } } }

1  val abs_vRC := undef; int abs_roundRC := 0;
2  set of val abs_valsRC := {};
3  proposeRC(int k, val v0) {
4    single (bool × val) abs_resRC := undef; bool res; val v;
5    assume(!(v0 = undef)); assume(pid() = ((k - 1) mod n) + 1);
6    ⟨ (res, v) := read(k); if (res = false) { linRC(undef, _); } ⟩
7    if (res) { if (v = undef) { v := v0; }
8              ⟨ res := write(k, v); if (res) { linRC(v, true); }
9              else { linRC(v, false); } ⟩
10             if (res) { return (true, v); } }
11    return (false, _); }

```

**Fig. 8.** Specification (top) and instrumented implementation (bottom) of *Round-Based Consensus*.

without contaminating any acceptor. Otherwise, the operation may contaminate some acceptor and the value  $vD$  is added to the set  $valsRC$  (line 6). Now, if the unspecified Boolean  $b$  is false, then the operation returns  $(false, \_)$  (lines 7 and 10), which models that the round will be stolen by a posterior operation. Finally, the operation succeeds if  $k$  is greater or equal than  $roundRC$  (line 7), and  $roundRC$  and  $vRC$  are updated and the operation returns  $(true, vRC)$  (lines 7–9).

In order to prove that the implementation in Fig. 5 linearises with respect to the specification on the top of Fig. 8, we use the instrumented implementation on the bottom of the same figure, where the abstract state is modelled by variables  $abs\_vRC$ ,  $abs\_roundRC$  and  $abs\_valsRC$  in lines 1–2, the local single-assignment variable  $abs\_resRC$  stores the result of the abstract operation, and the linearisation-point annotations  $linRC(vD, b)$  take a value and a Boolean parameters and invoke the non-deterministic abstract operation and disambiguate it by assigning the parameters to the unspecified  $vD$  and  $b$  of the specification. There are two linearisation points together with the invocations of `read` (line 6) and `write` (line 8). If `read` fails, then we linearise forcing the unspecified  $vD$  to be `undef` (line 6), which ensures that the abstract operation fails without adding any value to  $abs\_valsRC$  nor updating the round  $abs\_roundRC$ . Otherwise, if `write` succeeds with value  $v$ , then we linearise forcing the unspecified value  $vD$  and Boolean  $b$  to be  $v$  and `true` respectively (line 8). This ensures that

```

1  read(int k) {
2    ( val vD := random();
3      bool b := random(); val v;
4        assume(vD ∈ valsRR);
5        assume(pid() =
6          ((k - 1) mod n) + 1);
7        if (b) {
8          if (k >= roundRR) {
9            roundRR := k;
10           if (!(vRR = undef)) {
11             v := vRR; }
12           else { v := vD; } }
13         else { v := vD; }
14         return (true, v); }
15       else { return (false, _); } } }
16  val vRR := undef;
17  int roundRR := 0;
18  set of val valsRR := {undef};
19
20  write(int k, val vW) {
21    ( bool b := random();
22      assume(!(vW = undef));
23      assume(pid() =
24        ((k - 1) mod n) + 1);
25      valsRR := valsRR ∪ {vW};
26      if (b && (k >= roundRR)) {
27        roundRR := k;
28        vRR := vW;
29        return true; }
30      else { return false; } } }

```

**Fig. 9.** Specification of *Round-Based Register*.

the abstract operation succeeds and updates the round `abs_roundRC` to `k` and assigns `v` to the decided value `abs_vRC`. If `write` fails then we linearise forcing the unspecified `vD` and `b` to be `v` and `false` respectively (line 9). This ensures that the abstract operation fails.

**Theorem 2.** *The implementation of Round-Based Consensus in Fig. 5 linearises with respect to its specification on the top of Fig. 8.*

**Module *Round-Based Register*.** Figure 9 shows the module’s non-deterministic specification. Global variable `vRR` represents the decided value, initially `undef`. Global variable `roundRR` represents the current round, initially 0, and global set of values `valsRR`, initially containing `undef`, stores values that may have been proposed by some proposer. The specification is non-deterministic in that method `read` has unspecified local Boolean `b` and local value `vD` (we assume that `vD` is `valsRR`), and method `write` has unspecified local Boolean `b`. We assume the current process identifier is  $((k - 1) \bmod n) + 1$ .

Let us explain the specification of the `read` operation. The operation can succeed regardless of the proposer’s round `k`, depending on the value of the unspecified Boolean `b`. If `b` is `true` and the proposer’s round `k` is valid (line 8), then the read round is updated to `k` (line 9) and the operation returns `(true, v)` (line 14), where `v` is the read value, which coincides with the decided value if some decision was committed already or with `vD` otherwise. Now to the specification of operation `write`. The value `vW` is always added to the set `valsRR` (line 25). If the unspecified Boolean `b` is false (the round will be stolen by a posterior operation) or if the round `k` is non-valid, then the operation returns `false` (lines 26 and 30). Otherwise, the current round is updated to `k`, and the decided value `vRR` is updated to `vW` and the operation returns `true` (lines 27–29).

In order to prove that the implementation in Figs. 3 and 4 linearises with respect to the specification in Fig. 9, we use the instrumented implementation in

Figs. 10 and 11, which uses prophecy variables [1,26] that “guess” whether the execution of the method will reach a particular program location or not. The instrumented implementation also uses external linearisation points. In particular, the code of the acceptors may help to linearise some of the invocations to `read` and `write`, based on the prophecies and on auxiliary variables that count the number of acknowledgements sent by acceptors after each invocation of a `read` or a `write`. The next paragraphs elaborate on our use of prophecy variables and on our helping mechanism.

Variables `abs_vRR`, `abs_roundRR` and `abs_valsRR` in Fig. 10 model the abstract state. They are initially set to `undef`, 0 and the set containing `undef` respectively. Variable `abs_res_r[k]` is an infinite array of single-assignment pairs of Boolean and value that model the abstract results of the invocations to `read`. (Think of an infinite array as a map from integers to some type; we use the array notation for convenience.) Similarly, variable `abs_res_w[k]` is an infinite array of single-assignment Booleans that models the abstract results of the invocations to `write`. All the cells in both arrays are initially `undef` (*e.g.* the initial maps are empty). Variables `count_r[k]` and `count_w[k]` are infinite arrays of integers that model the number of acknowledgements sent (but not necessarily received yet) from acceptors in response to respectively read or write requests. All cells in both arrays are initially 0. The variable `proph_r[k]` is an infinite array of single-assignment pairs `bool × val`, modelling the prophecy for the invocations of `read`, and variable `proph_w[k]` is an infinite array of single-assignment Booleans modelling the prophecy for the invocations of `write`.

The linearisation-point annotations `linRE(k, vD, b)` for `read` take the proposer’s round `k`, a value `vD` and a Boolean `b`, and they invoke the abstract operation and disambiguate it by assigning the parameters to the unspecified `vD` and `b` of the specification on the left of Fig. 9. At the beginning of a `read(k)` (lines 11–14 of Fig. 10), the prophecy `proph_r[k]` is set to `(true, v)` if the invocation reaches `PL: RE_SUCC` in line 26. The `v` is defined to coincide with `maxV` at the time when that location is reached. That is, `v` is the value accepted at the greatest round by the acceptors acknowledging so far, or undefined if no acceptor ever accepted any value. If the operation reaches `PL: RE_FAIL` in line 24 instead, the prophecy is set to `(false, _)`. (If the method never returns, the prophecy is left `undef` since it will never linearise.) A successful `read(k)` linearises in the code of the acceptor in Fig. 11, when the  $\lceil (n+1)/2 \rceil$ th acceptor sends `[ackRE, k, v, w]`, and only if the prophecy is `(true, v)` and the operation was not linearised before (lines 10–14). We force the unspecified `vD` and `b` to be `v` and `true` respectively, which ensures that the abstract operation succeeds and returns `(true, v)`. A failing `read(k)` linearises at the `return` in the code of `read` (lines 23–24 of Fig. 10), after the reception of `[nackRE, k]` from one acceptor. We force the unspecified `vD` and `b` to be `undef` and `false` respectively, which ensures that the abstract operation fails.

The linearisation-point annotations `linWR(k, vW, b)` for `write` take the proposer’s round `k` and value `vW`, and a Boolean `b`, and they invoke the abstract operation and disambiguate it by assigning the parameter to the unspecified `b`

```

1  val abs_vRR := undef; int abs_roundRR := 0;
2  set of val abs_valsRR := {undef};
3  single val abs_res_r[1..∞] := undef;
4  single val abs_res_w[1..∞] := undef;
5  int count_r[1..∞] := 0; int count_w[1..∞] := 0;
6  single (bool × val) proph_r[1..∞] := undef;
7  single bool proph_w[i..∞] := undef;
8  read(int k) {
9    int j; val v; set of int Q; int maxKW; val maxV; msg m;
10   assume(pid() = ((k - 1) mod n) + 1);
11   ⟨ if (operation reaches PL: RE_SUCC and define v = maxV at that time) {
12     proph_r[k] := (true, v); }
13   else { if (operation reaches PL: RE_FAIL) {
14     proph_r[k] := (false, _); } } }
15   for (j := 1, j <= n, j++) { send(j, [RE, k]); }
16   maxKW := 0; maxV := undef; Q := {};
17   do { (j, m) := receive();
18     switch (m) {
19       case [ackRE, @k, v, kW]:
20         Q := Q ∪ {j};
21         if (kW >= maxKW) { maxKW := kW; maxV := v; }
22       case [nackRE, @k]:
23         ⟨ linRE(k, undef, false); proph_r[k] := undef;
24           return (false, _); ⟩ // PL: RE_FAIL
25     } if (|Q| = ⌈(n+1)/2⌉) {
26       return (true, maxV); } } // PL: RE_SUCC
27   while (true); }
28 write(int k, val vW) {
29   int j; set of int Q; msg m;
30   assume(!(vW = undef)); assume(pid() = ((k - 1) mod n) + 1);
31   ⟨ if (operation reaches PL: WR_SUCC) { proph_w[k] := true; }
32   else { if (operation reaches PL: WR_FAIL) {
33     proph_w[k] := false; } } }
34   for (j := 1, j <= n, j++) { send(j, [WR, k, vW]); }
35   Q := {};
36   do { (j, m) := receive();
37     switch (m) {
38       case [ackWR, @k]:
39         Q := Q ∪ {j};
40       case [nackWR, @k]:
41         ⟨ if (count_w[k] = 0) {
42           linWR(k, vW, false); proph_w[k] := undef; }
43           return false; ⟩ // PL: WR_FAIL
44     } if (|Q| = ⌈(n+1)/2⌉) {
45       return true; } } // PL: WR_SUCC
46   while (true); }

```

Fig. 10. Instrumented implementation of read and write methods.

```

1 process Acceptor(int j) {
2   val v := undef; int r := 0; int w := 0;
3   start() {
4     int i; msg m; int k;
5     do { (i, m) := receive();
6       switch (m) {
7         case [RE, k]:
8           if (k < r) { send(i, [nackRE, k]); }
9           else { r := k;
10              if (abs_res_r[k] = undef) {
11                if (proph_r[k] = (true, v)) {
12                  if (count_r[k] =  $\lceil (n+1)/2 \rceil - 1$ ) {
13                    linRE(k, v, true); } } }
14              count_r[k]++; send(i, [ackRE, k, v, w]); } } }
15         case [WR, k, vW]:
16           if (k < r) { send(j, i, [nackWR, k]); }
17           else { r := k; w := k; v := vW;
18              if (abs_res_w[k] = undef) {
19                if (!(proph_w[k] = undef)) {
20                  if (proph_w[k]) {
21                    if (count_w[k] =  $\lceil (n+1)/2 \rceil - 1$ ) {
22                      linWR(k, vW, true); } } }
23                  else { linWR(k, vW, false); } } }
24              count_w[k]++; send(j, i, [ackWR, k]); } } }
25       } }
26   while (true); } }

```

**Fig. 11.** Instrumented implementation of acceptor processes.

of the specification on the right of Fig. 9. At the beginning of a write( $k, vW$ ) (lines 31–33 of Fig. 10), the prophecy  $\text{proph\_r}[k]$  is set to **true** if the invocation reaches PL:WR\_SUCC in line 45, or to **false** if it reaches PL:WR\_FAIL in line 43 (or it is left **undef** if the method never returns). A successfully write( $k, vW$ ) linearises in the code of the acceptor in Fig. 11, when the  $\lceil (n+1)/2 \rceil$ th acceptor sends  $[\text{ackWR}, k]$ , and only if the prophecy is **true** and the operation was not linearised before (lines 17–24). We force the unspecified  $b$  to be **true**, which ensures that the abstract operation succeeds deciding value  $vW$  and updates  $\text{roundRR}$  to  $k$ . A failing write( $k, vW$ ) may linearise either at the return in its own code (lines 41–43 of Fig. 10) if the proposer received one  $[\text{nackWR}, k]$  and no acceptor sent any  $[\text{ackWR}, k]$  yet, or at the code of the acceptor, when the first acceptor sends  $[\text{ackWR}, k]$ , and only if the prophecy is **false** and the operation was not linearised before. In both cases, we force the unspecified  $b$  to be **false**, which ensures that the abstract operation fails.

**Theorem 3.** *The implementation of Round-Based Register in Figs. 10 and 11 linearises with respect to its specification in Fig. 9.*

## 5 Multi-Paxos via Network Transformations

We now turn to more complicated distributed protocols that build upon the idea of Paxos consensus. Our ultimate goal is to reuse the verification result from the Sects. 3 and 4, as well as the high-level round-based register interface. In this section, we will demonstrate how to reason about an implementation of Multi-Paxos as of an array of *independent* instances of the *Paxos* module defined previously, despite the subtle dependencies between its sub-components, as present in Multi-Paxos’s “canonical” implementations [5, 15, 27]. While an abstraction of Multi-Paxos to an array of independent shared “single-shot” registers is almost folklore, what appears to be inherently difficult is to verify a Multi-Paxos-based consensus (*wrt.* to the array-based abstraction) by means of *reusing* the proof of a SD-Paxos. All proofs of Multi-Paxos we are aware of are, thus, *non-modular* with respect to underlying SD-Paxos instances [5, 22, 24], *i.e.*, they require one to redesign the invariants of the *entire* consensus protocol.

This proof modularity challenge stems from the optimised nature of a classical Multi-Paxos protocol, as well as its real-world implementations [6]. In this part of our work is to distil such protocol-aware optimisations into a separate *network semantics layer*, and show that each of them refines the semantics of a Cartesian product-based view, *i.e.*, exhibits the very same client-observable behaviours. To do so, we will establish the refinement between the optimised implementations of Multi-Paxos and a simple Cartesian product abstraction, which will allow to extend the register-based abstraction, explored before in this paper, to what is considered to be a canonical amortised Multi-Paxos implementation.

### 5.1 Abstract Distributed Protocols

We start by presenting the formal definitions of encoding distributed protocols (including Paxos), their message vocabularies, protocol-based network semantics, and the notion of an observable behaviours.

**Protocols and Messages.** Figure 12 provides basic definitions of the distributed protocols and their components. Each protocol  $p$  is a tuple  $\langle \Delta, \mathcal{M}, \mathcal{S}_{\text{int}}, \mathcal{S}_{\text{rcv}}, \mathcal{S}_{\text{snd}} \rangle$ .  $\Delta$  is a set of local states, which can be assigned to each of the participating nodes, also determining the node’s role via an additional tag,<sup>4</sup> if necessary (*e.g.*, an acceptor and a proposer states in Paxos are different).  $\mathcal{M}$  is a “message vocabulary”, determining the set of messages that can be used for communication between the nodes.

Protocols	$\mathcal{P} \ni p \triangleq \langle \Delta, \mathcal{M}, \mathcal{S} \rangle$
Configurations	$\Sigma \ni \sigma \triangleq \text{Nodes} \rightarrow \Delta$
Internal steps	$\mathcal{S}_{\text{int}} \in \Delta \times \Delta$
Receive-steps	$\mathcal{S}_{\text{rcv}} \in \Delta \times \mathcal{M} \times \Delta$
Send-steps	$\mathcal{S}_{\text{snd}} \in \Delta \times \Delta \times \wp(\mathcal{M})$

**Fig. 12.** States and transitions.

<sup>4</sup> We leave out implicit the consistency laws for the state, that are protocol-specific.



$$\begin{array}{c}
\text{STEPINT} \\
\frac{n \in \text{dom}(\sigma) \quad \delta = \sigma(n) \quad \langle \delta, \delta' \rangle \in p.\mathcal{S}_{\text{int}} \quad \sigma' = \sigma[n \mapsto \delta']}{\langle \sigma, M \rangle \xrightarrow[\text{int}]{p} \langle \sigma', M \rangle} \\
\\
\text{STEPSEND} \\
\frac{n \in \text{dom}(\sigma) \quad \delta = \sigma(n) \quad \langle \delta, \delta', \text{ms} \rangle \in p.\mathcal{S}_{\text{snd}} \quad \sigma' = \sigma[n \mapsto \delta'] \quad M' = M \cup \text{ms}}{\langle \sigma, M \rangle \xrightarrow[\text{snd}]{p} \langle \sigma', M' \rangle} \\
\\
\text{STEPRECEIVE} \\
\frac{m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad \delta = \sigma(m.\text{to}) \quad \langle \delta, m, \delta' \rangle \in p.\mathcal{S}_{\text{rcv}} \quad m' = m[\text{active} \mapsto \text{False}] \quad \sigma' = \sigma[n \mapsto \delta'] \quad M' = M \setminus \{m\} \cup \{m'\}}{\langle \sigma, M \rangle \xrightarrow[\text{rcv}]{p} \langle \sigma', M' \rangle}
\end{array}$$

**Fig. 13.** Transition rules of the simple protocol-aware network semantics

Messages can be thought of as JavaScript-like dictionaries, pairing unique fields (isomorphic to strings) with their values. For the sake of a uniform treatment, we assume that each message  $m \in \mathcal{M}$  has at least two fields, *from* and *to* that point to the source and the destination node of a message, correspondingly. In addition to that, for simplicity we will assume that each message carries a Boolean field *active*, which is set to **True** when the message is sent and is set to **False** when the message is received by its destination node. This flag is required to keep history information about messages sent in the past, which is customary in frameworks for reasoning about distributed protocols [10, 23, 28]. We assume that a “message soup”  $M$  is a multiset of messages (*i.e.* a set with zero or more copies of each message) and we consider that each copy of the same message in the multiset has its own “identity”, and we write  $m \neq m'$  to represent that  $m$  and  $m'$  are not the same copy of a particular message.

Finally,  $\mathcal{S}_{\{\text{int}, \text{rcv}, \text{snd}\}}$  are step-relations that correspond to the internal changes in the local state of a node ( $\mathcal{S}_{\text{int}}$ ), as well as changes associated with sending ( $\mathcal{S}_{\text{snd}}$ ) and receiving ( $\mathcal{S}_{\text{rcv}}$ ) messages by a node, as allowed by the protocol. Specifically,  $\mathcal{S}_{\text{int}}$  relates a local node state before and after the allowed internal change;  $\mathcal{S}_{\text{rcv}}$  relates the initial state and an incoming message  $m \in \mathcal{M}$  with the resulting state;  $\mathcal{S}_{\text{snd}}$  relates the internal state, the output state and the set of atomically sent messages. For simplicity we will assume that  $\text{id} \subseteq \mathcal{S}_{\text{int}}$ .

In addition, we consider  $\Delta_0 \subseteq \Delta$ —the set of the allowed *initial* states, in which the system can be present at the very beginning of its execution. The global state of the network  $\sigma \in \Sigma$  is a map from node identifiers ( $n \in \mathbf{Nodes}$ ) to local states from the set of states  $\Delta$ , defined by the protocol.

**Simple Network Semantics.** The simple initial operational semantics of the network ( $\xrightarrow{p} \subseteq (\Sigma \times \wp(\mathcal{M})) \times (\Sigma \times \wp(\mathcal{M}))$ ) is parametrised by a protocol  $p$  and relates the initial *configuration* (*i.e.*, the global state and the set of messages) with the resulting configuration. It is defined via as a reflexive closure of the union of three relations  $\xrightarrow[\text{int}]{p} \cup \xrightarrow[\text{rcv}]{p} \cup \xrightarrow[\text{snd}]{p}$ , their rules are given in Fig. 13.

The rule `STEPINT` corresponds to a node  $n$  picked non-deterministically from the domain of a global state  $\sigma$ , executing an internal transition, thus changing its local state from  $\delta$  to  $\delta'$ . The rule `STEPRECEIVE` non-deterministically picks a  $m$  message from a message soup  $M \subseteq \mathcal{M}$ , changes the state using the protocol's receive-step relation  $p.\mathcal{S}_{\text{rcv}}$  at the corresponding host node  $to$ , and updates its local state accordingly in the common mapping ( $\sigma[to \mapsto \delta']$ ). Finally, the rule `STEPSEND`, non-deterministically picks a node  $n$ , executes a send-step, which results in updating its local state emission of a set of messages  $\text{ms}$ , which is added to the resulting soup. In order to “bootstrap” the execution, the initial states from the set  $\Delta_0 \subseteq \Delta$  are assigned to the nodes.

We next define the observable protocol behaviours *wrt.* the simple network semantics as the prefix-closed set of all system's configuration traces.

**Definition 1. (Protocol behaviours)**

$$\mathcal{B}_p = \bigcup_{m \in \mathbb{N}} \left\{ \langle \langle \sigma_0, M_0 \rangle, \dots, \langle \sigma_m, M_m \rangle \rangle \left| \begin{array}{l} \exists \delta_0^{n \in N} \in \Delta_0, \sigma_0 = \biguplus_{n \in N} [n \mapsto \delta_0^n] \wedge \\ \langle \sigma_0, M_0 \rangle \xrightarrow{p} \dots \xrightarrow{p} \langle \sigma_m, M_m \rangle \end{array} \right. \right\}$$

That is, the set of behaviours captures all possible configurations of initial states for a fixed set of nodes  $N \subseteq \text{Nodes}$ . In this case, the set of nodes  $N$  is an implicit parameter of the definition, which we fix in the remainder of this section.

*Example 1 (Encoding SD-Paxos).* An abstract distributed protocol for SD-Paxos can be extracted from the pseudo-code of Sect. 3 by providing a suitable small-step operational semantics à la Winskel [30]. We restraint ourselves from giving such formal semantics, but in Appendix D of the extended version of the paper we outline how the distributed protocol would be obtained from the given operational semantics and from the code in Figs. 3, 4 and 5.

**5.2 Out-of-Thin-Air Semantics**

We now introduce an intermediate version of a simple protocol-aware semantics that generates messages “out of thin air” according to a certain predicate  $\mathcal{P} \subseteq \Delta \times \mathcal{M}$ , which determines whether the network generates a certain message without exercising the corresponding send-transition. The rule is as follows:

$$\frac{\text{OTASEND} \quad n \in \text{dom}(\sigma) \quad \delta = \sigma(n) \quad \mathcal{P}(\delta, m) \quad M' = M \cup \{m\}}{\langle \sigma, M \rangle \xrightarrow[\text{ota}]{p, \mathcal{P}} \langle \sigma, M' \rangle}$$

That is, a random message  $m$  can be sent at any moment in the semantics described by  $\xrightarrow{p} \cup \xrightarrow[\text{ota}]{p, \mathcal{P}}$ , given that the node  $n$ , “on behalf of which” the message is sent is in a state  $\delta$ , such that  $\mathcal{P}(\delta, m)$  holds.

*Example 2.* In the context of Single-Decree Paxos, we can define  $\mathcal{P}$  as follows:

$$\mathcal{P}(\delta, m) \triangleq m.\text{content} = [\text{RE}, k] \wedge \delta.\text{pid} = n \wedge \delta.\text{role} = \text{Proposer} \wedge k \leq \delta.\text{kP}$$

In other words, if a node  $n$  is a *Proposer* currently operating with a round  $\delta.\text{kP}$ , the network semantics can always send another request “on its behalf”, thus generating the message “out-of-thin-air”. Importantly, the last conjunct in the definition of  $\mathcal{P}$  is in terms of  $\leq$ , rather than equality. This means that the predicate is intentionally loose, allowing for sending even “stale” messages, with expired rounds that are smaller than what  $n$  currently holds (no harm in that!).

By definition of single-decree Paxos protocol, the following lemma holds:

**Lemma 1 (OTA refinement).**  $\mathcal{B} \xrightarrow{p} \bigcup_{\text{ota}} \xrightarrow{p, \mathcal{P}} \subseteq \mathcal{B}_p$ , where  $p$  is an instance of the module *Paxos*, as defined in Sect. 3 and in Example 1.

### 5.3 Slot-Replicating Network Semantics

With the basic definitions at hand, we now proceed to describing alternative network behaviours that make use of a specific protocol  $p = \langle \Delta, \mathcal{M}, \mathcal{S}_{\text{int}}, \mathcal{S}_{\text{rcv}}, \mathcal{S}_{\text{snd}} \rangle$ , which we will consider to be fixed for the remainder of this section, so we will be at times referring to its components (*e.g.*,  $\mathcal{S}_{\text{int}}$ ,  $\mathcal{S}_{\text{rcv}}$ , *etc.*) without a qualifier.

$\begin{array}{l} \text{SRSTEPINT} \\ i \in I \quad n \in \text{dom}(\sigma) \\ \delta = \sigma(n)[i] \quad \langle \delta, \delta' \rangle \in p.\mathcal{S}_{\text{int}} \\ \sigma' = \sigma[n[i] \mapsto \delta'] \\ \hline \langle \sigma, M \rangle \xrightarrow[\text{int}]{x} \langle \sigma', M \rangle \end{array}$	$\begin{array}{l} \text{SRSTEPSEND} \\ i \in I \quad n \in \text{dom}(\sigma) \\ \delta = \sigma(n)[i] \quad \langle \delta, \delta', \text{ms} \rangle \in p.\mathcal{S}_{\text{snd}} \\ \sigma' = \sigma[n[i] \mapsto \delta'] \quad M' = M \cup \text{ms}[slot \mapsto i] \\ \hline \langle \sigma, M \rangle \xrightarrow[\text{snd}]{x} \langle \sigma', M' \rangle \end{array}$
$\begin{array}{l} \text{SRSTEPRECEIVE} \\ m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad \delta = \sigma(m.\text{to})[m.\text{slot}] \quad \langle \delta, m, \delta' \rangle \in p.\mathcal{S}_{\text{rcv}} \\ m' = m[\text{active} \mapsto \text{False}] \quad \sigma' = \sigma(n)[m.\text{slot} \mapsto \delta'] \quad M' = M \setminus \{m\} \cup \{m'\} \\ \hline \langle \sigma, M \rangle \xrightarrow[\text{rcv}]{x} \langle \sigma', M' \rangle \end{array}$	

**Fig. 14.** Transition rules of the slot-replicating network semantics.

Figure 14 describes a semantics of a *slot-replicating* (SR) network that exercises multiple copies of the *same* protocol instance  $p_i$  for  $i \in I$ , some, possibly infinite, set of indices, to which we will be also referring as *slots*. Multiple copies of the protocol are incorporated by enhancing the messages from  $p$ ’s vocabulary  $\mathcal{M}$  with the corresponding indices, and implementing the on-site dispatch of the indexed messages to corresponding protocol instances at each node. The local protocol state of each node is, thus, no longer a single element being updated,

but rather an *array*, mapping  $i \in I$  into  $\delta_i$ —the corresponding local state component. The small-step relation for SR semantics is denoted by  $\xRightarrow{\times}$ . The rule SRSTEPINT is similar to STEPINT of the simple semantics, with the difference that it picks not only a node but also an index  $i$ , thus referring to a specific component  $\sigma(n)[i]$  as  $\delta$  and updating it correspondingly ( $\sigma(n)[i] \mapsto \delta'$ ). For the remaining transitions, we postulate that the messages from  $p$ 's vocabulary  $p.\mathcal{M}$  are enhanced to have a dedicated field *slot*, which indicates a protocol copy at a node, to which the message is directed. The receive-rule SRSTEPRECEIVE is similar to STEPRECEIVE but takes into the account the value of  $m.slot$  in the received message  $m$ , thus redirecting it to the corresponding protocol instance and updating the local state appropriately. Finally, the rule SRSTEPSEND can be now executed for any slot  $i \in I$ , reusing most of the logic of the initial protocol and otherwise mimicking its simple network semantic counterpart STEPSEND.

Importantly, in this semantics, for two different slots  $i, j$ , such that  $i \neq j$ , the corresponding “projections” of the state behave *independently* from each other. Therefore, transitions and messages in the protocol instances indexed by  $i$  at different nodes *do not interfere* with those indexed by  $j$ . This observation can be stated formally. In order to do so we first defined the behaviours of slot-replicating networks and their projections as follows:

**Definition 2 (Slot-replicating protocol behaviours).**

$$\mathcal{B}_{\times} = \bigcup_{m \in \mathbb{N}} \left\{ \langle \langle \sigma_0, M_0 \rangle, \dots, \langle \sigma_m, M_m \rangle \rangle \left| \begin{array}{l} \exists \delta_0^{n \in N} \in \Delta_0, \\ \sigma_0 = \biguplus_{n \in N} [n \mapsto \{i \mapsto \delta_0^n \mid i \in I\}] \wedge \\ \langle \sigma_0, M_0 \rangle \xrightarrow{p} \dots \xrightarrow{p} \langle \sigma_m, M_m \rangle \end{array} \right. \right\}$$

That is, the slot-replicated behaviours are merely behaviours with respect to networks, whose nodes hold *multiple instances* of the same protocol, indexed by slots  $i \in I$ . For a slot  $i \in I$ , we define *projection*  $\mathcal{B}_{\times}|_i$  as a set of global state traces, where each node’s local states is restricted only to its  $i$ th component. The following simulation lemma holds naturally, connecting the state-replicating network semantics and simple network semantics.

**Lemma 2 (Slot-replicating simulation).** *For all  $I, i \in I, \mathcal{B}_{\times}|_i = \mathcal{B}_p$ .*

*Example 3 (Slot-replicating semantics and Paxos).* Given our representation of Paxos using roles (acceptors/proposers) encoded via the corresponding parts of the local state  $\delta$ , we can construct a “naïve” version of Multi-Paxos by using the SR semantics for the protocol. In such, every slot will correspond to a SD-Paxos instance, not interacting with any other slots. From the practical perspective, such an implementation is rather non-optimal, as it does not exploit dependencies between rounds accepted at different slots.

## 5.4 Widening Network Semantics

We next consider a version of the SR semantics, extended with a new rule for handling received messages. In the new semantics, dubbed *widening*, a node, upon receiving a message  $m \in T$ , where  $T \subseteq p.\mathcal{M}$ , for a slot  $i$ , *replicates* it for all slots from the index set  $I$ , for the very same node. The new rule is as follows:

$$\frac{\text{WSTEPRECEIVET} \quad \begin{array}{l} m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad \delta = \sigma(m.\text{to})[m.\text{slot}] \\ \langle \delta, m, \delta' \rangle \in p.\mathcal{S}_{\text{rev}} \quad m' = m[\text{active} \mapsto \text{False}] \quad \sigma' = \sigma(n)[m.\text{slot} \mapsto \delta'] \\ \text{ms} = \text{if } (m \in T) \text{ then } \{m' \mid m' = m[\text{slot} \mapsto j], j \in I\} \text{ else } \emptyset \end{array}}{\langle \sigma, M \rangle \xrightarrow[\text{rcv}]{\nabla} \langle \sigma', (M \setminus \{m\}) \cup \{m'\} \cup \text{ms} \rangle}$$

At first, this semantics seems rather unreasonable: it might create more messages than the system can “consume”. However, it is possible to prove that, under certain conditions on the protocol  $p$ , the set of behaviours observed under this semantics (*i.e.*, with SRSTEPRECEIVE replaced by WSTEPRECEIVET) is *not larger* than  $\mathcal{B}_\times$  as given by Definition 2. To state this formally we first relate the set of “triggering” messages  $T$  from WSTEPRECEIVET to a specific predicate  $\mathcal{P}$ .

**Definition 3 (OTA-compliant message sets).** The set of messages  $T \subseteq p.\mathcal{M}$  is OTA-compliant with the predicate  $\mathcal{P}$  iff for any  $b \in \mathcal{B}_p$  and  $\langle \sigma, M \rangle \in b$ , if  $m \in M$ , then  $\mathcal{P}(\sigma(m.\text{from}), m)$ .

In other words, the protocol  $p$  is relaxed enough to “justify” the presence of  $m$  in the soup at *any* execution, by providing the predicate  $\mathcal{P}$ , relating the message to the corresponding sender’s state. Next, we use this definition to slot-replicating and widening semantics via the following definition.

**Definition 4 ( $\mathcal{P}$ -monotone protocols).** A protocol  $p$  is  $\mathcal{P}$ -monotone iff for any,  $b \in \mathcal{B}_\times$ ,  $\langle \sigma, M \rangle \in b$ ,  $m, i = m.\text{slot}$ , and  $j \neq i$ , if  $\mathcal{P}(\sigma(m.\text{from})[i], \natural m)$  then we have that  $\mathcal{P}(\sigma(m.\text{from})[j], \natural m)$ , where  $\natural m$  “removes” the *slot* field from  $m$ .

Less formally, Definition 4 ensures that in a slot-replicated product  $\times$  of a protocol  $p$ , different components cannot perform “out of sync” *wrt.*  $\mathcal{P}$ . Specifically, if a node in  $i$ th projection is related to a certain message  $\natural m$  via  $\mathcal{P}$ , then any other projection  $j$  of the same node will be  $\mathcal{P}$ -related to this message, as well.

*Example 4.* This is a “non-example”. A version of slot-replicated SD-Paxos, where we allow for arbitrary increments of the round *per-slot* at a same proposer node (*i.e.*, out of sync), would not be monotone *wrt.*  $\mathcal{P}$  from Example 2. In contrast, a slot-replicated product of SD-Paxos instances with fixed rounds is monotone *wrt.* the same  $\mathcal{P}$ .

**Lemma 3.** *If  $T$  from WSTEPRECEIVET is OTA-compliant with predicate  $\mathcal{P}$ , such that  $\mathcal{B} \xrightarrow[p]{\cup} \xrightarrow[\text{ota}]{p, \mathcal{P}} \subseteq \mathcal{B} \xrightarrow[p]{\cup}$  and  $p$  is  $\mathcal{P}$ -monotone, then  $\mathcal{B} \xrightarrow[\nabla]{\cup} \subseteq \mathcal{B} \xrightarrow[\times]{\cup}$ .*

*Example 5 (Widening semantics and Paxos).* The SD-Paxos instance as described in Sect. 3 satisfies the refinement condition from Lemma 3. By taking  $T = \{m \mid m = \{\text{content} = [\text{RE}, \mathbf{k}]; \dots\}\}$  and using Lemma 3, we obtain the refinement between widened semantics and SR semantics of Paxos.

## 5.5 Optimised Widening Semantics

Our next step towards a realistic implementation of Multi-Paxos out of SD-Paxos instances is enabled by an observation that in the widening semantics, the replicated messages are *always* targeting the same node, to which the initial message  $m \in T$  was addressed. This means that we can optimise the receive-step, making it possible to execute multiple receive-transitions of the core protocol in batch. The following rule `OWSTEPRECEIVET` captures this intuition formally:

$$\frac{\text{OWSTEPRECEIVET} \quad m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad \langle \sigma', \text{ms} \rangle = \text{receiveAndAct}(\sigma, n, m)}{\langle \sigma, M \rangle \xrightarrow[\text{rcv}]{\nabla^*} \langle \sigma', M \setminus \{m\} \cup \{m[\text{active} \mapsto \text{False}]\} \cup \text{ms} \rangle}$$

where  $\text{receiveAndAct}(\sigma, n, m) \triangleq \langle \sigma', \text{ms} \rangle$ , such that  $\text{ms} = \bigcup_j \{m[\text{slot} \mapsto j] \mid m \in \text{ms}_j\}$ ,  $\forall j \in I, \delta = \sigma(m.\text{to})[j] \wedge \langle \delta_j, \sharp m, \delta_j^1 \rangle \in p.\mathcal{S}_{\text{rcv}} \wedge \langle \delta_j^1, \delta_j^2 \rangle \in p.\mathcal{S}_{\text{int}}^* \wedge \langle \delta_j^2, \delta_j^3, \text{ms}_j \rangle \in p.\mathcal{S}_{\text{snd}}$ ,  $\forall j \in I, \sigma'(m.\text{to})[j] = \delta_j^3$ .

In essence, the rule `OWSTEPRECEIVET` blends several steps of the widening semantics together for a single message: (a) it first receives the message and replicates it for all slots at a destination node; (b) performs receive-steps for the message’s replicas at each slot; (c) takes a number of internal steps, allowed by the protocol’s  $\mathcal{S}_{\text{int}}$ ; and (d) takes a send-transition, eventually sending all emitted message, instrumented with the corresponding slots.

*Example 6.* Continuing Example 5, with the same parameters, the optimising semantics will execute the transitions of an acceptor, *for all slots*, triggered by receiving a single `[RE, k]` message for a particular slot, sending back *all* the results for all the slots, which might either agree to accept the value or reject it.

The following lemma relates the optimising and the widening semantics.

**Lemma 4 (Refinement for OW semantics).** *For any  $b \in \mathcal{B}_{\xrightarrow{\nabla^*}}$  there exists  $b' \in \mathcal{B}_{\xrightarrow{\nabla}}$ , such that  $b$  can be obtained from  $b'$  by replacing sequences of configurations  $[\langle \sigma_k, M_k \rangle, \dots, \langle \sigma_{k+m}, M_{k+m} \rangle]$  that have just a single node  $n$ , whose local state is affected in  $\sigma_k, \dots, \sigma_{k+m}$ , by  $[\langle \sigma_k, M_k \rangle, \langle \sigma_{k+m}, M_{k+m} \rangle]$ .*

That is, behaviours in the optimised semantics are the same as in the widening semantics, modulo some sequences of locally taken steps that are being “compressed” to just the initial and the final configurations.

## 5.6 Bunching Semantics

As the last step towards Multi-Paxos, we introduce the final network semantics that optimises executions according to  $\xrightarrow{\nabla^*}$  described in previous section even further by making a simple addition to the message vocabulary of a slot-replicated SD-Paxos—*bunched messages*. A bunched message simply packages

$$\begin{array}{c}
 \text{BSTEPRECVB} \\
 \frac{m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad \langle \sigma', \text{ms} \rangle = \text{receiveAndAct}(\sigma, n, m) \quad M' = M \setminus \{m\} \cup \{m[\text{active} \mapsto \text{False}]\} \quad m' = \text{bunch}(\text{ms}, m.\text{to}, m.\text{from})}{\langle \sigma, M \rangle \xrightarrow[\text{rcv}]{B} \langle \sigma', M' \cup \{m'\} \rangle}
 \end{array}
 \qquad
 \begin{array}{c}
 \text{BSTEPRECVU} \\
 \frac{m \in M \quad m.\text{active} \quad m.\text{to} \in \text{dom}(\sigma) \quad m.\text{msgs} = \text{ms} \quad M' = M \setminus \{m\} \cup \text{ms}}{\langle \sigma, M \rangle \xrightarrow[\text{rcv}]{B} \langle \sigma, M' \rangle}
 \end{array}$$

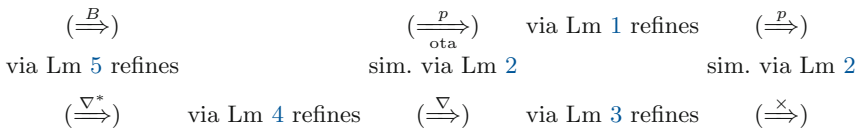
where  $\text{bunch}(\text{ms}, n_1, n_2) = \{\text{msgs} = \text{ms}; \text{from} = n_1; \text{to} = n_2; \text{active} = \text{True}\}$ .

**Fig. 15.** Added rules of the Bunching Semantics

together several messages, obtained typically as a result of a “compressed” execution via the optimised semantics from Sect. 5.5. We define two new rules for packaging and “unpackaging” certain messages in Fig. 15. The two new rules can be added to enhance either of the versions of the slot-replicating semantics shown before. In essence, the only effect they have is to combine the messages resulting in the execution of the corresponding steps of an optimised widening (via BSTEPRECVB), and to unpackage the messages  $\text{ms}$  from a bunching message, adding them back to the soup (BSTEPRECVU). The following natural refinement result holds:

**Lemma 5.** *For any  $b \in \mathcal{B} \xrightarrow{B}$  there exists  $b' \in \mathcal{B} \xrightarrow{\nabla^*}$ , such that  $b'$  can be obtained from  $b$  by replacing all bunched messages in  $b$  by their  $\text{msgs}$ -component.*

The rule BSTEPRECVU enables effective local caching of the bunched messages, so they are processed *on demand* on the recipient side (i.e., by the per-slot proposers), allowing the implementation to *skip* an entire round of Phase 1.



**Fig. 16.** Refinement between different network semantics.

```

1 proposeM(val^ v, val v0) {
2   < assume(!(v0 = undef));
3   if (*v = undef) { *v := v0; }
4   return *v; }
5 val vM[1..∞] := undef;
6 getR(int s) { return &(vM[s]); }
7 proposeM(getR(1), v);
8 proposeM(getR(2), v);

```

**Fig. 17.** Specification of *Multi-Paxos* and interaction via a *register provider*.

## 5.7 The Big Picture

What exactly have we achieved by introducing the described above family of semantics? As illustrated in Fig. 16, all behaviours of the leftmost-topmost, bunching semantics, which corresponds precisely to an implementation of Multi-Paxos with an “amortised” Phase 1, can be transitively related to the corresponding behaviours in the rightmost, vanilla slot-replicated version of a simple semantics (via the correspondence from Lemma 1) by constructing the corresponding refinement mappings [1], delivered by the proofs of Lemmas 3–5.

From the perspective of Rely/Guarantee reasoning, which was employed in Sect. 4, the refinement result from Fig. 16 justifies the replacement of a semantics on the right of the diagram by one to the left of it, as all program-level assertions will remain substantiated by the corresponding system configurations, as long as they are *stable* (i.e., resilient *wrt.* transitions taken by nodes different from the one being verified), which they are in our case.

## 6 Putting It All Together

We culminate our story of faithfully deconstructing and abstracting Paxos via a round-based register, as well as recasting Multi-Paxos via a series of network transformations, by showing how to *implement* the register-based abstraction from Sect. 3 in tandem with the network semantics from Sect. 5 in order to deliver provably correct, yet efficient, implementation of Multi-Paxos.

The crux of the composition of the two results—a register-based abstraction of SD-Paxos and a family of semantics-preserving network transformations—is a convenient interface for the end client, so she could interact with a consensus instance via the `proposeM` method in lines 1–4 of Fig. 17, no matter with which particular slot of a Multi-Paxos implementation she is interacting. To do so, we propose to introduce a *register provider*—a service that would give a client a “reference” to the consensus object to interact with. Lines 6–7 of Fig. 17 illustrate the interaction with the service provider, where the client requests two specific slots, 1 and 2, of Multi-Paxos by invoking `getR` and providing a slot parameter. In both cases the client proposes the very same value  $v$  in the two instances that run the same machinery. (Notice that, except for the reference to the consensus object, `proposeM` is identical to the `proposeP` on the right of Fig. 2, which we have verified *wrt.* linearisability in Sect. 3.)

The implementation of Multi-Paxos that we have in mind resembles the one in Figs. 3, 4 and 5 of Sect. 3, but where all the global data is provided by the register provider and passed by reference. What differs in this implementation with respect to the one in Sect. 3 and is hidden from the client is the semantics of the network layer used by the bottom layer (*cf.* left part of Fig. 2) of the register-based implementation. The Multi-Paxos instances run (without changing the register’s code) over this network layer, which “overloads” the meaning of the `send/receive` primitives from Figs. 3 and 4 to follow the bunching network semantics, described in Sect. 5.6.



**Theorem 4.** *The implementation of Multi-Paxos that uses a register provider and bunching network semantics refines the specification in Fig. 17.*

We implemented the register/network semantics in a proof-of-concept prototype written in Scala/Akka.<sup>5</sup> We relied on the abstraction mechanisms of Scala, allowing us to implement the register logic, verified in Sect. 4, separately from the network middle-ware, which has provided a family of Semantics from Sect. 5. Together, they provide a family of provably correct, modularly verified *distributed* implementations, coming with a simple *shared memory-like* interface.

## 7 Related Work

**Proofs of Linearisability via Rely/Guarantee.** Our work builds on the results of Boichat *et al.* [3], who were first to propose to a systematic deconstruction of Paxos into read/write operations of a *round-based register* abstraction. We extend and harness those abstractions, by intentionally introducing more non-determinism into them, which allows us to provide the first modular (*i.e.*, mutually independent) proofs of Proposer and Acceptor using Rely/Guarantee with linearisation points and prophecies. While several logics have been proposed recently to prove linearisability of concurrent implementations using Rely/Guarantee reasoning [14, 18, 19, 26], none of them considers message-passing distributed systems or consensus protocols.

**Verification of Paxos-Family Algorithms.** Formal verification of different versions of Paxos-family protocols *wrt.* inductive invariants and liveness has been a focus of multiple verification efforts in the past fifteen years. To name just a few, Lamport has specified and verified Fast Paxos [17] using TLA+ and its accompanying model checker [32]. Chand *et al.* used TLA+ to specify and verify Multi-Paxos implementation, similar to the one we considered in this work [5]. A version of SD-Paxos has been verified by Kellomaki using the PVS theorem prover [13]. Jaskelioff and Merz have verified Disk Paxos in Isabelle/HOL [12]. More recently, Rahli *et al.* formalised an executable version of Multi-Paxos in EventML [24], a dialect of NuPRL. Dragoi *et al.* [8] implemented and verified SD-Paxos in the PSYNC framework, which implements a partially synchronised model [7], supporting automated proofs of system invariants. Padon *et al.* have proved the system invariants and the consensus property of both simple Paxos and Multi-Paxos using the verification tool IVY [22, 23].

Unlike all those verification efforts that consider (Multi-/Disk/Fast/...)Paxos as a *single monolithic protocol*, our approach provides the first *modular* verification of single-decree Paxos using Rely/Guarantee framework, as well as the first verification of Multi-Paxos that directly reuses the proof of SD-Paxos.

<sup>5</sup> The code is available at <https://github.com/certichain/protocol-combinators>.

**Compositional Reasoning about Distributed Systems.** Several recent works have partially addressed modular formal verification of distributed systems. The IronFleet framework by Hawblitzel *et al.* has been used to verify both safety and liveness of a real-world implementation of a Paxos-based replicated state machine library and a lease-based shared key-value store [10]. While the proof is structured in a modular way by composing specifications in a way similar to our decomposition in Sects. 3 and 4, that work does not address the linearisability and does not provide composition of proofs about complex protocols (*e.g.*, Multi-Paxos) from proofs about its subparts

The Verdi framework for deductive verification of distributed systems [29, 31] suggests the idea of *Verified System Transformers* (VSTs), as a way to provide *vertical composition* of distributed system implementation. While Verdi’s VSTs are similar in its purpose and idea to our network transformations, they *do not* exploit the properties of the protocol, which was crucial for us to verify Multi-Paxos’s implementation.

The DIESEL framework [25, 28] addresses the problem of *horizontal composition* of distributed protocols and their client applications. While we do not compose Paxos with any clients in this work, we believe its register-based specification could be directly employed for verifying applications that use Paxos as its sub-component, which is what is demonstrated by our prototype implementation.

## 8 Conclusion and Future Work

We have proposed and explored two complementary mechanisms for modular verification of Paxos-family consensus protocols [15]: (a) non-deterministic register-based specifications in the style of Boichat *et al.* [3], which allow one to decompose the proof of protocol’s linearisability into separate independent “layers”, and (b) a family of protocol-aware transformations of network semantics, making it possible to reuse the verification efforts. We believe that the applicability of these mechanisms spreads beyond reasoning about Paxos and its variants and that they can be used for verifying other consensus protocols, such as Raft [21] and PBFT [4]. We are also going to employ network transformations to verify implementations of Mencius [20], and accommodate more protocol-specific optimisations, such as implementation of master leases and epoch numbering [6].

**Acknowledgements.** We thank the ESOP 2018 reviewers for their feedback. This work by was supported by ERC Starting Grant H2020-EU 714729 and EPSRC First Grant EP/P009271/1.

## References

1. Abadi, M., Lamport, L.: The existence of refinement mappings. In: LICS, pp. 165–175. IEEE Computer Society (1988)
2. Boichat, R., Dutta, P., Frølund, S., Guerraoui, R.: Deconstructing Paxos (2001). OAIPMH server at infoscience.epfl.ch, record 52373. <http://infoscience.epfl.ch/record/52373>

3. Boichat, R., Dutta, P., Frølund, S., Guerraoui, R.: Deconstructing Paxos. *SIGACT News* **34**(1), 47–67 (2003)
4. Castro, M., Liskov, B.: Practical Byzantine fault tolerance. In: *OSDI*, pp. 173–186. USENIX Association (1999)
5. Chand, S., Liu, Y.A., Stoller, S.D.: Formal verification of multi-Paxos for distributed consensus. In: Fitzgerald, J., Heitmeyer, C., Gnesi, S., Philippou, A. (eds.) *FM 2016*. LNCS, vol. 9995, pp. 119–136. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-48989-6\\_8](https://doi.org/10.1007/978-3-319-48989-6_8)
6. Chandra, T., Griesemer, R., Redstone, J.: Paxos made live: an engineering perspective. In: *PODC*, pp. 398–407. ACM (2007)
7. Charron-Bost, B., Merz, S.: Formal verification of a consensus algorithm in the heard-of model. *Int. J. Softw. Inform.* **3**(2–3), 273–303 (2009)
8. Dragoi, C., Henzinger, T.A., Zufferey, D.: PSync: a partially synchronous language for fault-tolerant distributed algorithms. In: *POPL*, pp. 400–415. ACM (2016)
9. Filipovic, I., O’Hearn, P.W., Rinetzky, N., Yang, H.: Abstraction for concurrent objects. *Theor. Comput. Sci.* **411**(51–52), 4379–4398 (2010)
10. Hawblitzel, C., Howell, J., Kapritsos, M., Lorch, J.R., Parno, B., Roberts, M.L., Setty, S.T.V., Zill, B.: IronFleet: proving practical distributed systems correct. In: *SOSP*, pp. 1–17. ACM (2015)
11. Herlihy, M., Wing, J.M.: Linearizability: a correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.* **12**(3), 463–492 (1990)
12. Jaskelioff, M., Merz, S.: Proving the correctness of disk Paxos. *Archive of Formal Proofs* (2005)
13. Kellomäki, P.: An annotated specification of the consensus protocol of Paxos using superposition in PVS. Technical report 36, Tampere University of Technology, Institute of Software Systems (2004)
14. Khyzha, A., Gotsman, A., Parkinson, M.: A generic logic for proving linearizability. In: Fitzgerald, J., Heitmeyer, C., Gnesi, S., Philippou, A. (eds.) *FM 2016*. LNCS, vol. 9995, pp. 426–443. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-48989-6\\_26](https://doi.org/10.1007/978-3-319-48989-6_26)
15. Lamport, L.: The part-time parliament. *ACM Trans. Comput. Syst.* **16**(2), 133–169 (1998)
16. Lamport, L.: Paxos made simple. *SIGACT News* **32**, 18–25 (2001)
17. Lamport, L.: Fast Paxos. *Distrib. Comput.* **19**(2), 79–103 (2006)
18. Liang, H., Feng, X.: Modular verification of linearizability with non-fixed linearization points. In: *PLDI*, pp. 459–470. ACM (2013)
19. Liang, H., Feng, X.: A program logic for concurrent objects under fair scheduling. In: *POPL*, pp. 385–399. ACM (2016)
20. Mao, Y., Junqueira, F.P., Marzullo, K.: Mencius: building efficient replicated state machine for WANs. In: *OSDI*, pp. 369–384. USENIX Association (2008)
21. Ongaro, D., Ousterhout, J.K.: In search of an understandable consensus algorithm. In: *2014 USENIX Annual Technical Conference*, pp. 305–319 (2014)
22. Padon, O., Losa, G., Sagiv, M., Shoham, S.: Paxos made EPR: decidable reasoning about distributed protocols. *PACMPL* **1**(OOPSLA), 108:1–108:31 (2017)
23. Padon, O., McMillan, K.L., Panda, A., Sagiv, M., Shoham, S.: Ivy: safety verification by interactive generalization. In: *PLDI*, pp. 614–630. ACM (2016)
24. Rahli, V., Guaspari, D., Bickford, M., Constable, R.L.: Formal specification, verification, and implementation of fault-tolerant systems using EventML. In: *AVOCS. EASST* (2015)
25. Sergey, I., Wilcox, J.R., Tatlock, Z.: Programming and proving with distributed protocols. *PACMPL* **2**(POPL), 28:1–28:30 (2018)

26. Vafeiadis, V.: Modular fine-grained concurrency verification. Ph.D. thesis, University of Cambridge (2007)
27. van Renesse, R., Altinbuken, D.: Paxos made moderately complex. *ACM Comput. Surv.* **47**(3), 42:1–42:36 (2015)
28. Wilcox, J.R., Sergey, I., Tatlock, Z.: Programming language abstractions for modularly verified distributed systems. In: *SNAPL. LIPIcs*, vol. 71, pp. 19:1–19:12. Schloss Dagstuhl (2017)
29. Wilcox, J.R., Woos, D., Panchekha, P., Tatlock, Z., Wang, X., Ernst, M.D., Anderson, T.E.: Verdi: a framework for implementing and formally verifying distributed systems. In: *PLDI*, pp. 357–368. ACM (2015)
30. Winskel, G.: *The Formal Semantics of Programming Languages*. The MIT Press, Cambridge (1993)
31. Woos, D., Wilcox, J.R., Anton, S., Tatlock, Z., Ernst, M.D., Anderson, T.E.: Planning for change in a formal verification of the Raft consensus protocol. In: *CPP*, pp. 154–165. ACM (2016)
32. Yu, Y., Manolios, P., Lamport, L.: Model checking TLA<sup>+</sup> specifications. In: Pierre, L., Kropf, T. (eds.) *CHARME 1999. LNCS*, vol. 1703, pp. 54–66. Springer, Heidelberg (1999). [https://doi.org/10.1007/3-540-48153-2\\_6](https://doi.org/10.1007/3-540-48153-2_6)

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

