

July 4, 2000

SRC Research
Report

163

Disk Paxos

Eli Gafni and Leslie Lamport

COMPAQ

Systems Research Center
130 Lytton Avenue
Palo Alto, California 94301

<http://www.research.digital.com/SRC/>

Systems Research Center

The charter of SRC is to advance both the state of knowledge and the state of the art in computer systems. From our establishment in 1984 by Digital Equipment Corporation (now Compaq), we have performed basic and applied research to support the company's business objectives. Our interests span scalable systems (including hardware, networks, distributed systems, and programming languages and technology), the Internet (including the web, and internet appliances), and human/computer interaction.

Our strategy is to test the technical and practical value of our ideas by building hardware and software prototypes and using them as daily tools. Interesting systems are too complex to be evaluated solely in the abstract; extended use allows us to investigate their properties in depth. This experience is useful in the short term in refining our designs, and invaluable in the long term in advancing our knowledge. Most of the major advances in information systems have come through this strategy, including personal computing, distributed systems, and the Internet.

We also perform complementary work of a more mathematical flavor. Some of it is in established fields of theoretical computer science, such as the analysis of algorithms, computational geometry, and logics of programming. Other work explores new ground motivated by problems that arise in our systems research.

We have a strong commitment to communicating our results; exposing and testing our ideas in the research and development communities leads to improved understanding. Our research report series supplements publication in professional journals and conferences while our technical note series complements research reports and journal/conference publication by allowing timely dissemination of recent research findings. We seek users for our prototype systems among those with whom we have common interests, and we encourage collaboration with university researchers.

Disk Paxos

Eli Gafni and Leslie Lamport

July 4, 2000

Author Affiliation

Eli Gafni is a member of the Computer Science Department at UCLA. He can be reached at `eli@cs.ucla.edu`.

©Digital Equipment Corporation 2000

This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of the Systems Research Center of Digital Equipment Corporation in Palo Alto, California; an acknowledgment of the authors and individual contributors to the work; and all applicable portions of the copyright notice. Copying, reproducing, or republishing for any other purpose shall require a license with payment of fee to the Systems Research Center. All rights reserved.

Author's Abstract

We present an algorithm, called Disk Paxos, for implementing a reliable distributed system with a network of processors and disks. Like the original Paxos algorithm, Disk Paxos maintains consistency in the presence of arbitrary non-Byzantine faults. Progress can be guaranteed as long as a majority of the disks are available, even if all processors but one have failed.

Contents

1	Introduction	1
2	The State-Machine Approach	2
3	An Informal Description of Disk Synod	3
3.1	The Algorithm	3
3.2	Why the Algorithm Works	5
3.3	Deriving Classic Paxos from Disk Paxos	7
4	Conclusion	8
4.1	Implementation Considerations	8
4.2	Concluding Remarks	9
	Bibliography	10
	Appendix	11
A.1	The Specification of Consensus	12
A.2	The Disk Synod Algorithm	15
A.3	An Assertional Proof	21
A.4	Proofs	28
A.4.1	Lemma I2c	29
A.4.2	Lemma BksOf	30
A.4.3	Lemma I2d	30
A.4.4	Lemma I2e	35
A.4.5	Lemma I2f	37
A.4.6	Theorem R2b	42

1 Introduction

Fault tolerance requires redundant components. Maintaining consistency in the event of a system partition makes it impossible for a two-component system to make progress if either component fails. There are innumerable fault-tolerant algorithms for implementing distributed systems, but all that we know of equate *component* with *processor*. But there are other types of components that one might replicate instead. In particular, modern networks can now include disk drives as independent components. Because commodity disks are cheaper than computers, it is attractive to use them as the replicated components for achieving fault tolerance. Commodity disks differ from processors in that they are not programmable, so we can't just substitute disks for processors in existing algorithms.

We present here an algorithm called *Disk Paxos* for implementing an arbitrary fault-tolerant system with a network of processors and disks. It maintains consistency in the event of any number of non-Byzantine failures. That is, the algorithm tolerates faulty processors that pause for arbitrarily long periods, fail completely, and possibly restart; and it tolerates lost and delayed messages. Disk Paxos guarantees progress if the system is stable and there is at least one nonfaulty processor that can read and write a majority of the disks. Stability means that each processor is either nonfaulty or has failed completely, and nonfaulty processors can access nonfaulty disks.

Disk Paxos is a variant of the classic Paxos algorithm [3, 10, 12], a simple, efficient algorithm that has been used in practical distributed systems [13, 16]. Classic Paxos can be viewed as an implementation of Disk Paxos in which there is one disk per processor, and a disk can be accessed directly only by its processor.

In the next section, we recall how to reduce the problem of implementing an arbitrary distributed system to the consensus problem. Section 3 informally describes Disk Synod, the consensus algorithm used by Disk Paxos. It includes a sketch of an incomplete correctness proof and explains the relation between Disk Synod and the Synod protocol of classic Paxos. Section 4 briefly discusses some implementation details and contains the conventional concluding remarks. An appendix gives formal specifications of the consensus problem and the Disk Synod algorithm, and sketches a rigorous correctness proof.

2 The State-Machine Approach

The state-machine approach [5, 14] is a general method for implementing an arbitrary distributed system. The system is designed as a deterministic state machine that executes a sequence of commands, and a consensus algorithm ensures that, for each n , all processors agree on the n^{th} command. This reduces the problem of building an arbitrary system to solving the consensus problem. In the consensus problem, each processor p starts with an input value $input[p]$, and all processors output the same value, which equals $input[p]$ for some p . A solution should be:

Consistent All values output are the same.

Nonblocking If the system is stable and a nonfaulty processor can communicate with a majority of disks, then the processor will eventually output a value.

It has long been known that a consistent, nonblocking consensus algorithm requires a three-phase commit protocol [15], with *voting*, *prepare to commit*, and *commit* phases. Nonblocking algorithms that use fewer phases don't guarantee consistency. For example, the group communication algorithms of Isis [2] permit two processors belonging to the current group to disagree on whether a message was broadcast in a previous group to which they both belonged. This algorithm cannot, by itself, guarantee consistency because disagreement about whether a message had been broadcast can result in disagreement about the output value.

The classic Paxos algorithm [3, 10, 12] achieves its efficiency by using a three-phase commit protocol, called the *Synod* algorithm, in which the value to be committed is not chosen until the second phase. When a new leader is elected, it executes the first phase just once for the entire sequence of consensus algorithms performed for all later system commands. Only the last two phases are performed separately for each individual command.

In the *Disk Synod* algorithm, the consensus algorithm used by Disk Paxos, each processor has an assigned block on each disk. The algorithm has two phases. In each phase, a processor writes to its own block and reads each other processor's block on a majority of the disks.¹ Only the last phase needs to be executed anew for each command. So, in the normal steady-state case, a leader chooses a state-machine command by executing a single write to each of its blocks and a single read of every other processor's blocks.

¹There is also an extra phase that a processor executes when recovering from a failure.

The classic result of Fischer, Lynch, and Patterson [4] implies that a purely asynchronous nonblocking consensus algorithm is impossible. So, real-time clocks must be introduced. The typical industry approach is to use an *ad hoc* algorithm based on timeouts to elect a leader, and then have the leader choose the output. It is easy to devise a leader-election algorithm that works when the system is stable, which means that it works most of the time. It is very hard to make one that always works correctly even when the system is unstable. Both classic Paxos and Disk Paxos also assume a real-time algorithm for electing a leader. However, the leader is used only to ensure progress. Consistency is maintained even if there are multiple leaders. Thus, if the leader-election algorithm fails because the network is unstable, the system can fail to make progress; it cannot become inconsistent. The system will again make progress when it becomes stable and a single leader is elected.

3 An Informal Description of Disk Synod

We now informally describe the Disk Synod algorithm and explain why it works. We also discuss its relation to classic Paxos's Synod Protocol. Remember that, in normal operation, only a single leader will be executing the algorithm. The other processors do nothing; they simply wait for the leader to inform them of the outcome. However, the algorithm must preserve consistency even when it is executed by multiple processors, or when the leader fails before announcing the outcome and a new leader is chosen.

3.1 The Algorithm

We assume that each processor p starts with an input value $input[p]$.² As in Paxos's Synod algorithm, a processor executes a sequence of numbered ballots, with increasing ballot numbers. A ballot number is a positive integer, and different processors use different ballot numbers. For example, if the processors are numbered from 1 through N , then processor i could use ballot numbers i , $i + N$, $i + 2N$, etc. A ballot has two phases:

Phase 1 Choose a value v .

Phase 2 Try to commit v .

In either phase, a processor aborts its ballot if it learns that another processor has begun a higher-numbered ballot. In that case, the processor may

²If processor p fails, it can restart with a new value of $input[p]$.

then choose a higher ballot number and start a new ballot. If the processor completes phase 2 without aborting—that is, without learning of a higher-numbered ballot—then value v is *committed* and the processor can output it. Since a processor does not choose the value to be committed until phase 2, phase 1 can be performed once for any number of separate instances of the algorithm.

To ensure consistency, we must guarantee that two different values cannot be successfully committed—either by different processors or by the same processor in two different ballots. To ensure that the algorithm is nonblocking, we must guarantee that, if there is only a single processor p executing it, then p will eventually commit a value.

In practice, when a processor successfully commits a value, it will write on its disk block that the value was committed and also broadcast that fact to the other processors. If a processor learns that a value has been committed, it will abort its ballot and simply output the value. It is obvious that this optimization preserves correctness; we will not consider it further.

To execute the algorithm, a processor p maintains a record $dblock[p]$ containing the following three components:

mbal The current ballot number.

bal The largest ballot number for which p reached phase 2.

inp The value p tried to commit in ballot number *bal*.

Initially, *bal* equal 0, *inp* equals a special value *NotAnInput* that is not a possible input value, and *mbal* is any ballot number. We let $disk[d][p]$ be the block on disk d in which processor p writes $dblock[p]$. We assume that reading and writing a block are atomic operations.

Processor p executes phase 1 or 2 of a ballot as follows. For each disk d , it tries first to write $dblock[p]$ to $disk[d][p]$ and then to read $disk[d][q]$ for all other processors q . It aborts the ballot if, for any d and q , it finds $disk[d][q].mbal > dblock[p].mbal$. The phase completes when p has written and read a majority of the disks, without reading any block whose *mbal* component is greater than $dblock[p].mbal$. When it completes phase 1, p chooses a new value of $dblock[p].inp$, sets $dblock[p].bal$ to $dblock[p].mbal$ (its current ballot number), and begins phase 2. When it completes phase 2, p has committed $dblock[p].inp$.

To complete our description of the two phases, we now describe how processor p chooses the value of $dblock[p].inp$ that it tries to commit in phase 2. Let $blocksSeen$ be the set consisting of $dblock[p]$ and all the records $disk[d][q]$

read by p in phase 1. Let $nonInitBlks$ be the subset of $blocksSeen$ consisting of those records whose inp field is not *NotAnInput*. If $nonInitBlks$ is empty, then p sets $dblock[p].inp$ to its own input value $input[p]$. Otherwise, it sets $dblock[p].inp$ to $bk.inp$ for some record bk in $nonInitBlks$ having the largest value of $bk.bal$.

Finally, we describe what processor p does when it recovers from a failure. In this case, p reads its own block $disk[d][p]$ from a majority of disks d . It then sets $dblock[p]$ to any block bk it read having the maximum value of $bk.mbal$, and it starts a new ballot by increasing $dblock[p].mbal$ and beginning phase 1.

3.2 Why the Algorithm Works

Suppose processor p can read and write a majority of the disks, and all processors other than p stop executing the algorithm. In this case, p will eventually choose a ballot number greater than the $mbal$ field of all blocks on the disks it can read, and its ballot will succeed. Hence, this algorithm is nonblocking, in the sense explained above.

We now explain, intuitively, why the Disk Synod algorithm maintains consistency. First, we consider the following shared-memory version of the algorithm that uses single-writer, multiple-reader regular registers.³ Instead of writing to disk, processor p writes $dblock[p]$ to a shared register; and it reads the values of $dblock[q]$ for other processors q from the registers. A processor chooses its bal and inp values for phase 2 the same way as before, except that it reads just one $dblock$ value for each other processor, rather than one from each disk. We assume for now that processors do not fail.

To prove consistency, we must show that, for any processors p and q , if p finishes phase 2 and commits the value v_p and q finishes phase 2 and commits the value v_q , then $v_p = v_q$. Let b_p and b_q be the respective ballot numbers on which these values are committed. Without loss of generality, we can assume $b_p \leq b_q$. Moreover, using induction on b_q , we can assume that, if any processor r starts phase 2 for a ballot b_r with $b_p \leq b_r < b_q$, then it does so with $dblock[r].inp = v_p$.

When reading in phase 2, p cannot have seen the value of $dblock[q].mbal$ written by q in phase 1—otherwise, p would have aborted. Hence p 's read of $dblock[q]$ in phase 2 did not follow q 's phase 1 write. Because reading follows writing in each phase, this implies that q 's phase 1 read of $dblock[p]$

³A regular register is one in which a read that does not overlap a write returns the register's current value, and a read that overlaps one or more writes returns either the register's previous value or one of the values being written [6].

must have followed p 's phase 2 write. Hence, q read the current (final) value of $dblock[p]$ in phase 1—a record with bal field b_p and inp field v_p . Let bk be any other block that q read in its phase 1. Since q did not abort, $b_q > bk.mbal$. Since $bk.mbal \geq bk.bal$ for any block bk , this implies $b_q > bk.bal$. By the induction assumption, we obtain that, if $bk.bal \geq b_p$, then $bk.inp = v_p$. Since this is true for all blocks bk read by q in phase 1, and since q read the final value of $dblock[p]$, the algorithm implies that q must set $dblock[q].inp$ to v_p for phase 2, proving that $v_p = v_q$.

To obtain the Disk Synod algorithm from the shared-memory version, we use a technique due to Attiya, Bar-Noy, and Dolev [1] to implement a single-writer, multiple reader register with a network of disks. To write a value, a processor writes the value together with a version number to a majority of the disks. To read, a processor reads a majority of the disks and takes the value with the largest version number. Since two majorities of disks contain at least one disk in common, a read must obtain either the last version for which the write was completed, or else a later version. Hence, this implements a regular register. With this technique, we transform the shared-memory version into a version for a network of processors and disks.

The actual Disk Synod algorithm simplifies the algorithm obtained by this transformation in two ways. First, the version number is not needed. The $mbal$ and bal values play the role of a version number. Second, a processor p need not choose a single version of $dblock[q]$ from among the ones it reads from disk. Because $mbal$ and bal values do not decrease, earlier versions have no effect.

So far, we have ignored processor failures. There is a trivial way to extend the shared-memory algorithm to allow processor failures. A processor recovers by simply reading its $dblock$ value from its register and starting a new ballot. A failed process then acts like one in which a processor may start a new ballot at any time. We can show that this generalized version is also correct. However, in the actual disk algorithm, a processor can fail while it is writing. This can leave its disk blocks in a state in which no value has been written to a majority of the disks. Such a state has no counterpart in the shared-memory version. There seems to be no easy way to derive the recovery procedure from a shared-memory algorithm. The proof of the complete Disk Synod algorithm, with failures, is much more complicated than the one for the simple shared-memory version. Trying to write the kind of behavioral proof given above for the simple algorithm leads to the kind of complicated, error-prone reasoning that we have learned to avoid. Instead, we sketch a rigorous assertional proof in the appendix.

3.3 Deriving Classic Paxos from Disk Paxos

In the usual view of a distributed fault-tolerant system, a processor performs actions and maintains its state in local memory, using stable storage to recover from failures. An alternative view is that a processor maintains the state of its stable storage, using local memory only to cache the contents of stable storage. Identifying disks with stable storage, a traditional distributed system is then a network of disks and processors in which each disk belongs to a separate processor; other processors can read a disk only by sending messages to its owner.

Let us now consider how to implement Disk Synod on a network of processors that each has its own disk. To perform phase 1 or 2, a processor p would access a disk d by sending a message containing $dblock[p]$ to disk d 's owner q . Processor q could write $dblock[p]$ to $disk[d][p]$, read $disk[d][r]$ for all $r \neq p$, and send the values it read back to p . However, examining the Disk Synod algorithm reveals that there's no need to send back all that data. All p needs are (i) to know if its $mbal$ field is larger than any other block's $mbal$ field and, if it is, (ii) the bal and inp fields for the block having the maximum bal field. Hence, q need only store on disk three values: the bal and inp fields for the block with maximum bal field, and the maximum $mbal$ field of all disk blocks. Of course, q would have those values cached in its memory, so it would actually write to disk only if any of those values are changed.

A processor must also read its own disk blocks to recover from a failure. Suppose we implement Disk Synod by letting p write to its own disk before sending messages to any other processor. This ensures that its own disk has the maximum value of $disk[d][p].mbal$ among all the disks d . Hence, to restart after a failure, p need only read its block from its own disk. In addition to the $mbal$, bal , and inp value mentioned above, p would also keep the value of $dblock[p]$ on its disk.

We can now compare this algorithm with classic Paxos's Synod protocol [10]. The $mbal$, bal , and inp components of $dblock[p]$ are just $lastTried[p]$, $nextBal[p]$, and $prevVote[p]$ of the Synod Protocol. Phase 1 of the Disk Synod algorithm corresponds to sending the *NextBallot* message and receiving the *LastVote* responses in the Synod Protocol. Phase 2 corresponds to sending the *BeginBallot* and receiving the *Voted* replies.⁴ The Synod Protocol's *Success* message corresponds to the optimization mentioned above

⁴In the Synod Protocol, a processor q does not bother sending a response if p sends it a disk block with a value of $mbal$ smaller than one already on disk. Sending back the maximum $mbal$ value is an optimization mentioned in [10].

of recording on disk that a value has been committed.

This version of the Disk Synod algorithm differs from the Synod Protocol in two ways. First, the Synod Protocol's *NextBallot* message contains only the *mbal* value; it does not contain *bal* and *inp* values. To obtain the Synod Protocol, we would have to modify the Disk Synod algorithm so that, in phase 1, it writes only the *mbal* field of its disk block and leaves the *bal* and *inp* fields unchanged. The algorithm remains correct, with essentially the same proof, under this modification. However, the modification makes the algorithm harder to implement with real disks.

The second difference between this version of the Disk Synod algorithm and the Synod Protocol is in the restart procedure. A disk contains only the aforementioned *mbal*, *bal*, and *inp* values. It does not contain a separate copy of its owner's *dblock* value. The Synod Protocol can be obtained from the following variant of the Disk Synod algorithm. Let *bk* be the block *disk*[*d*][*p*] with maximum *bal* field read by processor *p* in the restart procedure. Processor *p* can begin phase 1 with *bal* and *inp* values obtained from any disk block *bk'*, written by any processor, such that *bk'.bal* \geq *bk.bal*. It can be shown that the Disk Synod algorithm remains correct under this modification too.

4 Conclusion

4.1 Implementation Considerations

Implicit in our description of the Disk Synod algorithm are certain assumptions about how reading and writing are implemented when disks are accessed over a network. If operations sent to the disks may be lost, a processor *p* must receive an acknowledgment from disk *d* that its write to *disk*[*d*][*p*] succeeded. This may require *p* to explicitly read its disk block after writing it. If operations may arrive at the disk in a different order than they were sent, *p* will have to wait for the acknowledgment that its write to disk *d* succeeded before reading other processors' blocks from *d*. Moreover, some mechanism is needed to ensure that a write from an earlier ballot does not arrive after a write from a later one, overwriting the later value with the earlier one. How this is achieved will be system dependent. (It is impossible to implement any fault-tolerant system if writes to disk can linger arbitrarily long in the network and cause later values to be overwritten.)

Recall that, in Disk Paxos, a sequence of instances of the Disk Synod algorithm is used to commit a sequence of commands. In a straightforward implementation of Disk Paxos, processor *p* would write to its disk blocks the

value of $dblock[p]$ for the current instance of Disk Synod, plus the sequence of all commands that have already been committed. The sequence of all commands that have ever been committed is probably too large to fit on a single disk block. However, the complete sequence can be stored on multiple disk blocks. All that must be kept in the same disk block as $dblock[p]$ is a pointer to the head of the queue. For most applications, it is not necessary to remember the entire sequence of commands [10, Section 3.3.2]. In many cases, all the data that must be kept will fit in a single disk block.

In the application for which Disk Paxos was devised (a future Compaq product), the set of processors is not known in advance. Each disk contains a directory listing the processors and the locations of their disk blocks. Before reading a disk, a processor reads the disk's directory. To write a disk's directory, a processor must acquire a lock for that disk by executing a real-time mutual exclusion algorithm based on Fischer's protocol [7]. A processor joins the system by adding itself to the directory on a majority of disks.

4.2 Concluding Remarks

We have presented Disk Paxos, an efficient implementation of the state machine approach in a system in which processors communicate by accessing ordinary (nonprogrammable) disks. In the normal case, the leader commits a command by writing its own block and reading every other processor's block on a majority of the shared disks. This is clearly the minimal number of disk accesses needed.

Disk Paxos was motivated by the recent development of the Storage Area Network (SAN)—an architecture consisting of a network of computers and disks in which all disks can be accessed by each computer. Commodity disks are cheaper than computers, so using redundant disks for fault tolerance is more economical than using redundant computers. Moreover, since disks do not run application-level programs, they are less likely to crash than computers.

Because commodity disks are not programmable, we could not simply substitute disks for processors in the classic Paxos algorithm. Instead we took the ideas of classic Paxos and transplanted them to the SAN environment. What we obtained is almost, but not quite, a generalization of classic Paxos. Indeed, when Disk Paxos is instantiated to a single disk, we obtain what may be called Shared-Memory Paxos. Algorithms for shared memory are usually more succinct and clear than their message passing counterparts. Thus, Disk Paxos can be considered yet another revisiting of classic

Paxos that exposes its underlying ideas by removing the message-passing clutter. Perhaps other distributed algorithms can also be made more clear by recasting them in a shared-memory setting.

References

- [1] Hagit Attiya, Amotz Bar-Noy, and Danny Dolev. Sharing memory robustly in message-passing systems. *Journal of the ACM*, 42(1):124–142, January 1995.
- [2] Kenneth Birman, André Schiper, and Pat Stephenson. Lightweight causal and atomic group multicast. *ACM Transactions on Computer Systems*, 9(3):272–314, August 1991.
- [3] Roberto De Prisco, Butler Lampson, and Nancy Lynch. Revisiting the Paxos algorithm. In Marios Mavronicolas and Philippas Tsigas, editors, *Proceedings of the 11th International Workshop on Distributed Algorithms (WDAG 97)*, volume 1320 of *Lecture Notes in Computer Science*, pages 111–125, Saarbrücken, Germany, 1997. Springer-Verlag.
- [4] Michael J. Fischer, Nancy Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, April 1985.
- [5] Leslie Lamport. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, July 1978.
- [6] Leslie Lamport. On interprocess communication. *Distributed Computing*, 1:77–101, 1986.
- [7] Leslie Lamport. A fast mutual exclusion algorithm. *ACM Transactions on Computer Systems*, 5(1):1–11, February 1987.
- [8] Leslie Lamport. The temporal logic of actions. *ACM Transactions on Programming Languages and Systems*, 16(3):872–923, May 1994.
- [9] Leslie Lamport. How to write a proof. *American Mathematical Monthly*, 102(7):600–608, August-September 1995.
- [10] Leslie Lamport. The part-time parliament. *ACM Transactions on Computer Systems*, 16(2):133–169, May 1998.

- [11] Leslie Lamport. Specifying concurrent systems with TLA⁺. In Manfred Broy and Ralf Steinbrüggen, editors, *Calculational System Design*, pages 183–247, Amsterdam, 1999. IOS Press.
- [12] Butler W. Lampson. How to build a highly available system using consensus. In Ozalp Babaoglu and Keith Marzullo, editors, *Distributed Algorithms*, volume 1151 of *Lecture Notes in Computer Science*, pages 1–17, Berlin, 1996. Springer-Verlag.
- [13] Edward K. Lee and Chandramohan Thekkath. Petal: Distributed virtual disks. In *Proceedings of the Seventh International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-VII)*, pages 84–92, New York, October 1996. ACM Press.
- [14] Fred B. Schneider. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys*, 22(4):299–319, December 1990.
- [15] Marion Dale Skeen. *Crash Recovery in a Distributed Database System*. PhD thesis, University of California, Berkeley, May 1982.
- [16] Chandramohan Thekkath, Timothy Mann, and Edward K. Lee. Frangipani: A scalable distributed file system. In *Proceedings of the 16th ACM Symposium on Operating Systems Principles*, pages 224–237, New York, October 1997. ACM Press.
- [17] Yuan Yu, Panagiotis Manolios, and Leslie Lamport. Model checking TLA⁺ specifications. In Laurence Pierre and Thomas Kropf, editors, *Correct Hardware Design and Verification Methods*, volume 1703 of *Lecture Notes in Computer Science*, pages 54–66, Berlin, Heidelberg, New York, September 1999. Springer-Verlag. 10th IFIP wg 10.5 Advanced Research Working Conference, CHARME '99.

Appendix

We now give a precise specification of the consensus problem solved by the Disk Synod algorithm and of the algorithm itself. The specification is written in TLA⁺ [11], a formal language that combines the temporal logic of actions (TLA) [8], set theory, and first-order logic with notation for making

definitions and encapsulating them in modules. In the course of writing the specifications, we try to explain any TLA⁺ notation whose meaning is not self-evident. These specifications have been debugged with the aid of the TLC model checker [17].⁵

We prove only consistency of the algorithm. We feel that the nonblocking property is sufficiently obvious not to need a formal proof. We therefore do not specify or reason about liveness properties. This means that we make hardly any use of temporal logic.

A.1 The Specification of Consensus

We now formally specify the consensus problem. We assume N processors, numbered 1 through N . Each processor p has two registers: an input register $input[p]$ that initially equals some element of a set $Inputs$ of possible input values, and an output register $output[p]$ that initially equals a special value $NotAnInput$ that is not an element of $Inputs$. Processor p chooses an output value by setting $output[p]$. It can also fail, which it does by setting $input[p]$ to any value in $Inputs$ and resetting $output[p]$ to $NotAnInput$. The precise condition to be satisfied is that, if some processor p ever sets $output[p]$ to some value v , then

- v must be a value that is, or at one time was, the value of $input[q]$ for some processor q
- if any processor r (including p itself) later sets $output[r]$ to some value w other than $NotAnInput$, then $w = v$.

We specify only safety. There is no liveness requirement, so the specification is satisfied if no processor ever changes $output[p]$.

TLA⁺ specifications are organized into modules. The specification of consensus is in a module named *SynodSpec*, which begins:

MODULE *SynodSpec*

EXTENDS *Naturals*

The EXTENDS statement imports the *Naturals* module, which defines the set *Nat* of natural numbers and the usual arithmetic operations. It also defines $i..j$ to be the set of natural numbers from i through j . We next declare the specification's two constants: the number N of processors, and the set

⁵The typeset versions were generated manually from the actual TLA⁺ specifications by a procedure that may have introduced errors.

Inputs of inputs; and we assert the assumption that N is a positive natural number.

CONSTANT N , *Inputs*
 ASSUME $(N \in \text{Nat}) \wedge (N > 0)$

In TLA^+ , every value is a set, so we don't have to assert that *Inputs* is a set. We next define two constants: the set *Proc* of processors, and the value *NotAnInput*. In TLA^+ , \triangleq means *is defined to equal*, and CHOOSE $x : F(x)$ equals an arbitrary value x such that $F(x)$ is true (if such an x exists).

Proc $\triangleq 1..N$
NotAnInput \triangleq CHOOSE $c : c \notin \text{Inputs}$

We next declare the variables *input* and *output*.

VARIABLES *input*, *output*

To write the specification, we introduce two internal variables: *allInput*, which equals the set of all current and past values of *input*[p], for all processors p ; and *chosen*, which records the first input value output by some processor (and hence, the value that all processors must henceforth output). These variables are internal or “hidden” variables. In TLA, such variables are bound variables of the temporal existential quantifier \exists . Since internal variables aren't part of the specification, they should not be declared in module *SynodSpec*. One way to introduce such variables in TLA^+ is to declare them in a submodule. So, we introduce a submodule called *Inner*.

┌────────────────────────────────── MODULE *Inner* ───────────────────────────────────┐
 VARIABLES *allInput*, *chosen*

Before going further, we explain some TLA^+ notation. In programming languages, the variables *input* and *output* would be arrays indexed by the *Proc*. What programmers call an array indexed by S , mathematicians call a function with domain S . TLA^+ uses the notation $[x \in S \mapsto e(x)]$ for the function f with domain S such that $f[x] = e(x)$ for all x in S . It denotes by $[S \rightarrow T]$ the set of all functions f with domain S such that $f[x] \in T$ for all $x \in S$. TLA^+ allows a conjunction or disjunction to be written as a list of formulas bulleted by \wedge or \vee . Indentation is used to eliminate parentheses.

We now define *IInit* to be the predicate describing the initial state.

$$\begin{aligned}
IInit &\triangleq \wedge \text{input} \in [Proc \rightarrow Inputs] \\
&\wedge \text{output} = [p \in Proc \mapsto NotAnInput] \\
&\wedge \text{chosen} = NotAnInput \\
&\wedge \text{allInput} = \{\text{input}[p] : p \in Proc\}
\end{aligned}$$

We next define the two actions, $Choose(p)$ and $Fail(p)$, that describe the operations that a processor p can perform. In TLA, an action is a formula with primed and unprimed variables that describes the relation between the values of the variables in a new (primed) state and their values in an old (unprimed) state. For example, in a system with the two variables x and y , the action $(x' = x + 1) \wedge (y' = y)$ corresponds to the programming-language statement $x := x + 1$. A conjunct with no primed variables is an enabling condition.

In TLA⁺, the expression $[f \text{ EXCEPT } ![x] = e]$ represents the function \widehat{f} that is the same as f except that $\widehat{f}[x] = e$. Thus, $f' = [f \text{ EXCEPT } ![c] = e]$ corresponds to the programming-language statement $f[c] := e$, except that it says nothing about variables other than f . An action must explicitly state what remains unchanged. We do this with the expression $\text{UNCHANGED } v$, which means $v' = v$. Leaving a tuple $\langle v_1, \dots, v_n \rangle$ unchanged is equivalent to leaving all its components v_i unchanged.

The $Choose(p)$ action represents the processor p choosing its output. It is enabled iff $\text{output}[p]$ equals $NotAnInput$. If chosen is $NotAnInput$, then chosen and $\text{output}[p]$ are set to any element of allInput . Otherwise, $\text{output}[p]$ is set to chosen .

$$\begin{aligned}
Choose(p) &\triangleq \\
&\wedge \text{output}[p] = NotAnInput \\
&\wedge \text{IF } \text{chosen} = NotAnInput \\
&\quad \text{THEN } \exists ip \in \text{allInput} : \wedge \widehat{\text{chosen}}' = ip \\
&\quad \quad \quad \wedge \text{output}' = [\text{output} \text{ EXCEPT } ![p] = ip] \\
&\quad \text{ELSE } \wedge \text{output}' = [\text{output} \text{ EXCEPT } ![p] = \text{chosen}] \\
&\quad \quad \quad \wedge \text{UNCHANGED } \text{chosen} \\
&\wedge \text{UNCHANGED } \langle \text{input}, \text{allInput} \rangle
\end{aligned}$$

The $Fail(p)$ action represents processor p failing. It is always enabled. It sets $\text{output}[p]$ to $NotAnInput$, sets $\text{input}[p]$ to any element of $Inputs$, and adds that element to the set allInput .

$$\begin{aligned}
Fail(p) &\triangleq \\
&\wedge \text{output}' = [\text{output} \text{ EXCEPT } ![p] = NotAnInput]
\end{aligned}$$

$$\begin{aligned}
\wedge \exists ip \in Inputs : \wedge input' &= [input \text{ EXCEPT } ![p] = ip] \\
\wedge allInput' &= allInput \cup \{ip\} \\
\wedge \text{UNCHANGED } chosen &
\end{aligned}$$

We next define the next-state action $INext$, which describes all possible steps. We then define $ISpec$, the specification with the internal variables $chosen$ and $allInput$ visible. It asserts that the initial state satisfies $IInit$, and every step either satisfies $INext$ or leaves all the variables unchanged. Formula $ISpec$ is defined to be a temporal formula, using the ordinary operator \square (always) of temporal logic, and the TLA notation that $[N]_v$ equals $N \vee (v' = v)$. These definitions end the submodule.

$$\begin{aligned}
INext &\triangleq \exists p \in Proc : Choose(p) \vee Fail(p) \\
ISpec &\triangleq IInit \wedge \square [INext]_{(input, output, chosen, allInput)}
\end{aligned}$$

Finally, we define $SynodSpec$, the complete specification, to be $ISpec$ with the variables $chosen$ and $allInput$ hidden—that is, quantified with the temporal existential quantifier \exists of TLA. The precise meaning of the TLA⁺ constructs used here is unimportant.

$$\begin{aligned}
IS(chosen, allInput) &\triangleq \text{INSTANCE } Inner \\
SynodSpec &\triangleq \exists chosen, allInput : IS(chosen, allInput)!ISpec
\end{aligned}$$

This ends module $SynodSpec$.

A.2 The Disk Synod Algorithm

The Disk Synod algorithm is specified by a module $DiskSynod$ that imports all the declarations and definitions from the $SynodSpec$ module.

$$\begin{aligned}
&\text{MODULE } DiskSynod \\
&\text{EXTENDS } SynodSpec
\end{aligned}$$

The algorithm assumes that different processors use different ballot numbers. Instead of fixing some specific assignment choice of ballot numbers, we let $Ballot(p)$ represent the set of ballot numbers that processor p can use, where $Ballot$ is an unspecified constant operator.

We have described the algorithm in terms of a majority of disks. The property of majorities we need is that any two majorities has a disk in common. If there are an even number d of disks, we can maintain that property

even if we consider certain sets containing $d/2$ disks to constitute a majority. We let *IsMajority* be an unspecified predicate so that if *IsMajority*(S) and *IsMajority*(T) is true for two sets S and T of disks, then S and T are not disjoint.

The module now declares *Ballot*, *IsMajority*, and the constant *Disk* that represents the set of disks. It also asserts the assumptions we make about them. In TLA⁺, the expression SUBSET S denotes the set of all subsets of the set S .

```

CONSTANTS Ballot( $\_$ ), Disk, IsMajority( $\_$ )
ASSUME  $\wedge \forall p \in Proc : \wedge Ballot(p) \subseteq \{n \in Nat : n > 0\}$ 
       $\wedge \forall q \in Proc \setminus \{p\} : Ballot(p) \cap Ballot(q) = \{\}$ 
       $\wedge \forall S, T \in SUBSET Disk :$ 
         $IsMajority(S) \wedge IsMajority(T) \Rightarrow (S \cap T \neq \{\})$ 

```

We next define two constants: the set *DiskBlock* of all possible records that a processor can write to its disk blocks, and the record *InitDB* that is the initial value of all disk blocks. In TLA⁺, $[f_1 \mapsto v_1, \dots, f_n \mapsto v_n]$ is the record r with fields f_1, \dots, f_n such that $r.f_i = v_i$, for all i in $1 \dots n$, and $[f_1 : S_1, \dots, f_n : S_n]$ is the set of all such records with v_i an element of the set S_i , for all i in $1 \dots n$. The set $\bigcup S$, the union of all the elements of S , is written UNION S . For example, UNION $\{A, B, C\}$ equals $A \cup B \cup C$.

```

DiskBlock  $\triangleq$  [mbal : (UNION  $\{Ballot(p) : p \in Proc\}) \cup \{0\}$ ,
               bal : (UNION  $\{Ballot(p) : p \in Proc\}) \cup \{0\}$ ,
               inp : Inputs  $\cup \{NotAnInput\}$ ]
InitDB  $\triangleq$  [mbal  $\mapsto 0$ , bal  $\mapsto 0$ , inp  $\mapsto NotAnInput$ ]

```

We now declare all the specification's variables—except for *input* and *output*, whose declarations are imported from *SynodSpec*. We have described the variables *disk* (the contents of the disks) and *dblock* above. We let *phase*[p] be the current phase of processor p , which will be set to 0 when p fails and to 3 when p chooses its output. For convenience, we let each processor start in phase 0 and begin the algorithm as if it were recovering from a failure. The variables *disksWritten* and *blocksRead* record a processor's progress in the current phase; *disksWritten*[p] is the set of disks that processor p has written, and *blocksRead*[p][d] is the set of values p has read from disk d . More precisely, *blocksRead*[p][d] is a set of records with *block* and *proc* fields, where $[block \mapsto bk, proc \mapsto q]$ is in *blocksRead*[p][d] iff p has read the value bk from *disk*[d][q] in the current phase. For convenience, we declare

vars to be the tuple of all the specification's variables. We also define the predicate *Init* that defines the initial values of all variables.

VARIABLES *disk*, *dblock*, *phase*, *disksWritten*, *blocksRead*
 $vars \triangleq \langle input, output, disk, phase, dblock, disksWritten, blocksRead \rangle$
 $Init \triangleq \wedge input \in [Proc \rightarrow Inputs]$
 $\wedge output = [p \in Proc \mapsto NotAnInput]$
 $\wedge disk = [d \in Disk \mapsto [p \in Proc \mapsto InitDB]]$
 $\wedge phase = [p \in Proc \mapsto 0]$
 $\wedge dblock = [p \in Proc \mapsto InitDB]$
 $\wedge output = [p \in Proc \mapsto NotAnInput]$
 $\wedge disksWritten = [p \in Proc \mapsto \{\}]$
 $\wedge blocksRead = [p \in Proc \mapsto [d \in Disk \mapsto \{\}]]$

We now define two operators that describe the state of a processor during the current phase: *hasRead*(*p*, *d*, *q*) is true iff *p* has read *disk*[*d*][*q*], and *allBlocksRead*(*p*) equals the set of all *disk*[*d*][*q*] values that *p* has read during the current phase. The TLA⁺ expression LET *def* IN *exp* equals expression *exp* in the context of the local definitions in *def*.

$hasRead(p, d, q) \triangleq \exists br \in blocksRead[p][d] : br.proc = q$
 $allBlocksRead(p) \triangleq$
 LET *allRdBlks* \triangleq UNION {*blocksRead*[*p*][*d*] : *d* \in *Disk*}
 IN {*br.block* : *br* \in *allRdBlks*}

We now define *InitializePhase*(*p*) to be an action that sets *disksWritten*[*p*] and *blocksRead*[*p*] to their initial values, to indicate that *p* has done no reading or writing yet in the current phase. This action will be used to define other actions that make up the next-state relation; it itself is not part of the next-state relation.

$InitializePhase(p) \triangleq$
 $\wedge disksWritten' = [disksWritten \text{ EXCEPT } ![p] = \{\}]$
 $\wedge blocksRead' = [blocksRead \text{ EXCEPT } ![p] = [d \in Disk \mapsto \{\}]]$

We now define the actions that will form part of the next-state action. These actions describe all the atomic actions of the algorithm that a processor *p* can perform. The first is *StartBallot*(*p*) in which *p* initiates a new ballot. We all *p* to do this at any time during phase 1 or 2. The action sets *phase*[*p*] to 1, increases *dblock*[*p*].*mbal*, and initializes the phase,

$$\begin{aligned}
\text{StartBallot}(p) &\triangleq \\
&\wedge \text{phase}[p] \in \{1, 2\} \\
&\wedge \text{phase}' = [\text{phase EXCEPT } ![p] = 1] \\
&\wedge \exists b \in \text{Ballot}(p) : \wedge b > \text{dblock}[p].\text{mbal} \\
&\quad \wedge \text{dblock}' = [\text{dblock EXCEPT } ![p].\text{mbal} = b] \\
&\wedge \text{InitializePhase}(p) \\
&\wedge \text{UNCHANGED } \langle \text{input}, \text{output}, \text{disk} \rangle
\end{aligned}$$

In action $\text{Phase1or2Write}(p, d)$, processor p writes $\text{disk}[d][p]$ and adds d to the set $\text{disksWritten}[p]$ of disks written by p . The action is enabled iff p is in phase 1 or 2.⁶ In the TLA⁺ expression $[f \text{ EXCEPT } ![c] = e]$, an @ appearing in e stands for $f[c]$. Thus, $x' = [x \text{ EXCEPT } ![c] = @ + 1]$ corresponds to the programming-language statement $x[c] := x[c] + 1$.

$$\begin{aligned}
\text{Phase1or2Write}(p, d) &\triangleq \\
&\wedge \text{phase}[p] \in \{1, 2\} \\
&\wedge \text{disk}' = [\text{disk EXCEPT } ![d][p] = \text{dblock}[p]] \\
&\wedge \text{disksWritten}' = [\text{disksWritten EXCEPT } ![p] = @ \cup \{d\}] \\
&\wedge \text{UNCHANGED } \langle \text{input}, \text{output}, \text{phase}, \text{dblock}, \text{blocksRead} \rangle
\end{aligned}$$

Action $\text{Phase1or2Read}(p, d, q)$ describes p reading $\text{disk}[d][q]$. It is enabled iff d is in $\text{disksWritten}[p]$, meaning that p has already written its block to disk d . (This implies that p is in phase 1 or 2.) We allow p to reread a disk block it has already read. If $\text{disk}[d][q].\text{mbal}$ is less than p 's current mbal value, then $\text{blocksRead}[p][d]$ is updated and p continues executing its ballot. Otherwise, p aborts the ballot and begins a new one. The EXCEPT construct has a more general form for “arrays of arrays”. For example, the formula $x' = [x \text{ EXCEPT } ![a][b] = e]$ corresponds to the programming-language statement $x[a][b] := e$.

$$\begin{aligned}
\text{Phase1or2Read}(p, d, q) &\triangleq \\
&\wedge d \in \text{disksWritten}[p] \\
&\wedge \text{IF } \text{disk}[d][q].\text{mbal} < \text{dblock}[p].\text{mbal} \\
&\quad \text{THEN } \wedge \text{blocksRead}' = \\
&\quad \quad [\text{blocksRead EXCEPT} \\
&\quad \quad \quad ![p][d] = @ \cup \{[\text{block} \mapsto \text{disk}[d][q], \text{proc} \mapsto q]\}] \\
&\quad \wedge \text{UNCHANGED} \\
&\quad \quad \langle \text{input}, \text{output}, \text{disk}, \text{phase}, \text{dblock}, \text{disksWritten} \rangle \\
&\quad \text{ELSE } \text{StartBallot}(p)
\end{aligned}$$

⁶We could add the enabling condition $d \notin \text{disksWritten}[p]$, but it's not necessary because the action is a no-op, leaving all variables unchanged, if p has already written its current value of $\text{dblock}[p]$ to disk d .

The action $EndPhase1or2(p)$ describes processor p successfully finishing phase 1 or 2. It is enabled when p is in phase 1 or 2 and, on a majority of the disks, p has written its block and read every other processor's block. When p finishes phase 1, it sets $dblock[p].inp$ and $dblock[p].bal$ as described in Section 3.1 and starts phase 2. When p finishes phase 2, it sets $output[p]$, sets $phase[p]$ to 3, and terminates. (However, it could still fail and start again.) The TLA⁺ EXCEPT construct applies to records as well as functions, and it can have multiple “replacements” separated by commas.

$$\begin{aligned}
EndPhase1or2(p) &\triangleq \\
&\wedge IsMajority(\{d \in disksWritten[p] : \\
&\quad \forall q \in Proc \setminus \{p\} : hasRead(p, d, q)\}) \\
&\wedge \vee \wedge phase[p] = 1 \\
&\quad \wedge dblock' = \\
&\quad \quad [dblock \text{ EXCEPT} \\
&\quad \quad \quad ![p].bal = dblock[p].mbal, \\
&\quad \quad \quad ![p].inp = \\
&\quad \quad \quad \text{LET } blocksSeen \triangleq allBlocksRead(p) \cup \{dblock[p]\} \\
&\quad \quad \quad nonInitBlks \triangleq \\
&\quad \quad \quad \quad \{bs \in blocksSeen : bs.inp \neq NotAnInput\} \\
&\quad \quad \quad \quad maxBlk \triangleq \\
&\quad \quad \quad \quad \text{CHOOSE } b \in nonInitBlks : \\
&\quad \quad \quad \quad \quad \forall c \in nonInitBlks : b.bal \geq c.bal \\
&\quad \quad \quad \text{IN IF } nonInitBlks = \{\} \text{ THEN } input[p] \\
&\quad \quad \quad \quad \text{ELSE } maxBlk.inp] \\
&\quad \wedge \text{UNCHANGED } output \\
&\quad \vee \wedge phase[p] = 2 \\
&\quad \quad \wedge output' = [output \text{ EXCEPT } ![p] = dblock[p].inp] \\
&\quad \quad \wedge \text{UNCHANGED } dblock \\
&\quad \wedge phase' = [phase \text{ EXCEPT } ![p] = @ + 1] \\
&\quad \wedge InitializePhase(p) \\
&\quad \wedge \text{UNCHANGED } \langle input, disk \rangle
\end{aligned}$$

Action $Fail(p)$ represents a failure by processor p . The action is always enabled. It chooses a new value of $input[p]$, sets $phase[p]$ to 0 and initializes $dblock[p]$, $output[p]$, $disksWritten[p]$, and $blocksRead[p]$.

$$\begin{aligned}
Fail(p) &\triangleq \\
&\wedge \exists ip \in Inputs : input' = [input \text{ EXCEPT } ![p] = ip]
\end{aligned}$$

$$\begin{aligned}
&\wedge \textit{phase}' = [\textit{phase} \text{ EXCEPT } ![p] = 0] \\
&\wedge \textit{dblock}' = [\textit{dblock} \text{ EXCEPT } ![p] = \textit{InitDB}] \\
&\wedge \textit{output}' = [\textit{output} \text{ EXCEPT } ![p] = \textit{NotAnInput}] \\
&\wedge \textit{InitializePhase}(p) \\
&\wedge \text{UNCHANGED } \textit{disk}
\end{aligned}$$

The next two actions describe failure recovery. In $\textit{Phase0Read}(p, d)$, processor p reads $\textit{disk}[d][p]$, recording the value read in $\textit{blocksRead}[p]$. Again, we allow redundant reads of the same disk block. In $\textit{EndPhase0}(p)$, processor p completes its recovery and enters phase 1, as described in Section 3.1.

$$\begin{aligned}
\textit{Phase0Read}(p, d) &\triangleq \\
&\wedge \textit{phase}[p] = 0 \\
&\wedge \textit{blocksRead}' = [\textit{blocksRead} \text{ EXCEPT} \\
&\quad ![p][d] = @ \cup \{[\textit{block} \mapsto \textit{disk}[d][p], \textit{proc} \mapsto p]\}] \\
&\wedge \text{UNCHANGED } \langle \textit{input}, \textit{output}, \textit{disk}, \textit{phase}, \textit{dblock}, \textit{disksWritten} \rangle \\
\textit{EndPhase0}(p) &\triangleq \\
&\wedge \textit{phase}[p] = 0 \\
&\wedge \textit{IsMajority}(\{d \in \textit{Disk} : \textit{hasRead}(p, d, p)\}) \\
&\wedge \exists b \in \textit{Ballot}(p) : \\
&\quad \wedge \forall r \in \textit{allBlocksRead}(p) : b > r.\textit{mbal} \\
&\quad \wedge \textit{dblock}' = [\textit{dblock} \text{ EXCEPT} \\
&\quad \quad ![p] = [(\text{CHOOSE } r \in \textit{allBlocksRead}(p) : \\
&\quad \quad \quad \forall s \in \textit{allBlocksRead}(p) : r.\textit{bal} \geq s.\textit{bal}) \\
&\quad \quad \text{EXCEPT } !.\textit{mbal} = b]] \\
&\wedge \textit{InitializePhase}(p) \\
&\wedge \textit{phase}' = [\textit{phase} \text{ EXCEPT } ![p] = 1] \\
&\wedge \text{UNCHANGED } \langle \textit{input}, \textit{output}, \textit{disk} \rangle
\end{aligned}$$

As in most TLA specifications, we define the next-state action \textit{Next} that describes all possible steps of all processors. We then define the formula $\textit{DiskSynodSpec}$, our specification of the algorithm, to assert that the initial state satisfies \textit{Init} and every step either satisfies \textit{Next} or leaves all the variables unchanged.

$$\begin{aligned}
\textit{Next} &\triangleq \exists p \in \textit{Proc} : \\
&\quad \vee \textit{StartBallot}(p) \\
&\quad \vee \exists d \in \textit{Disk} : \vee \textit{Phase0Read}(p, d) \\
&\quad \quad \vee \textit{Phase1or2Write}(p, d) \\
&\quad \quad \vee \exists q \in \textit{Proc} \setminus \{p\} : \textit{Phase1or2Read}(p, d, q) \\
&\quad \vee \textit{EndPhase1or2}(p)
\end{aligned}$$

$$\begin{aligned}
& \vee \text{Fail}(p) \\
& \vee \text{EndPhase0}(p) \\
\text{DiskSynodSpec} & \triangleq \text{Init} \wedge \square[\text{Next}]_{\text{vars}}
\end{aligned}$$

The module ends by asserting the correctness of the algorithm, which means that the algorithm's specification implies the formula SynodSpec that is its correctness condition.

THEOREM $\text{DiskSynodSpec} \Rightarrow \text{SynodSpec}$

A.3 An Assertional Proof

We now sketch a proof of the correctness of the Disk Synod algorithm—that is, a proof that DiskSynodSpec implies SynodSpec . Formula SynodSpec equals $\exists \text{chosen}, \text{allInput} : \text{ISpec}$.⁷ To prove such a formula, we must find Skolem functions with which to instantiate the bound variables chosen and allInput , and then prove that DiskSynodSpec implies ISpec , when chosen and allInput are defined to equal those Skolem functions. The choice of Skolem functions is called a *refinement mapping*. However, we cannot define such a refinement mapping because chosen and allInput record history that is not present in the actual state of the algorithm. Instead, we add chosen and allInput to the algorithm specification as *history variables*. Formally, we define a specification HDiskSynodSpec such that

$$\text{DiskSynodSpec} \equiv \exists \text{chosen}, \text{allInput} : \text{HDiskSynodSpec}$$

We then prove that HDiskSynodSpec implies ISpec , from which we infer by simple logic that DiskSynodSpec implies SynodSpec .

The initial predicate HInit of HDiskSynodSpec is the conjunction of the initial predicate Init of DiskSynodSpec with formulas that specify the initial values of chosen and allInput . Its next-state action HNext is the conjunction of the next-state action Next of DiskSynodSpec with formulas that specify the values of chosen' and $\text{allInput}'$ as functions of the (unprimed and primed) values of the other variables. A general theorem of TLA asserts that, if the variable x does not occur in I , N , or the tuple \mathbf{y} of variables, then

$$I \wedge \square[N]_{\mathbf{y}} \equiv \exists x : (I \wedge (x = f(\mathbf{y}))) \wedge \square[N \wedge (x' = g(x, \mathbf{y}, \mathbf{y}'))]_{\langle x, \mathbf{y} \rangle}$$

⁷Actually, $\exists \text{chosen}, \text{allInput} : \text{ISpec}$ is not a legal TLA⁺ formula; we should instead write $\exists \text{chosen}, \text{allInput} : \text{IS}(\text{chosen}, \text{allInput})! \text{ISpec}$.

for any f and g . This result implies that the specification obtained from $HDiskSynodSpec$ by hiding (existentially quantifying) $chosen$ and $allInput$ is equivalent to $DiskSynodSpec$.

We define $HDiskSynodSpec$ in a module $HDiskSynod$ that extends the $DiskSynod$ module and declares $chosen$ and $allInput$ as variables.

MODULE $HDiskSynod$
EXTENDS $DiskSynod$ VARIABLES $allInput, chosen$
<p>The initial values of $chosen$ and $allInput$ are the same as in the initial predicate of $ISpec$.</p> $ \begin{aligned} HInit &\triangleq \wedge Init \\ &\wedge chosen = NotAnInput \\ &\wedge allInput = \{input[p] : p \in Proc\} \end{aligned} $ <p>The action $HNext$ ensures that $chosen$ equals the first $output$ value that is different from $NotAnInput$, and that $allInput$ always equals the set of all $input$ values that have appeared thus far.</p> $ \begin{aligned} HNext &\triangleq \\ &\wedge Next \\ &\wedge chosen' = \text{LET } hasOutput(p) \triangleq output'[p] \neq NotAnInput \\ &\quad \text{IN IF } \vee chosen \neq NotAnInput \\ &\quad \quad \vee \forall p \in Proc : \neg hasOutput(p) \\ &\quad \quad \text{THEN } chosen \\ &\quad \quad \text{ELSE } output'[\text{CHOOSE } p \in Proc : hasOutput(p)] \\ &\wedge allInput' = allInput \cup \{input'[p] : p \in Proc\} \end{aligned} $ <p>The module then defines $HDiskSynodSpec$ in the usual way, and asserts that it implies $ISpec$, with $chosen$ and $allInput$ replaced by the variables of the same name declared in the current module. (Again, the details of how this is expressed in TLA^+ are not important.)</p> $ HDiskSynodSpec \triangleq HInit \wedge \square[HNext]_{\langle vars, chosen, allInput \rangle} $ <p>THEOREM $HDiskSynodSpec \Rightarrow IS(chosen, allInput)!ISpec$</p>

We now outline the proof of this theorem. Let $ivars$ be the tuple of all variables of $ISpec$:

$ivars \triangleq \langle input, output, chosen, allInput \rangle$

To prove that *HDiskSynodSpec* implies *ISpec* we must prove⁸

THEOREM *R1* $HInit \Rightarrow IInit$

THEOREM *R2* $HInit \wedge \square[HNEXT]_{\langle vars, chosen, allInput \rangle} \Rightarrow \square[INEXT]_{ivars}$

The proof of *R1* is trivial. To prove *R2*, standard TLA reasoning shows that it suffices to find a state predicate *HInv* for which we can prove:

THEOREM *R2a* $HInit \wedge \square[HNEXT]_{\langle vars, chosen, allInput \rangle} \Rightarrow \square HInv$

THEOREM *R2b* $HInv \wedge HInv' \wedge HNEXT \Rightarrow INEXT \vee (\text{UNCHANGED } ivars)$

A predicate *HInv* satisfying *R2a* is said to be an invariant of the specification $HInit \wedge \square[HNEXT]_{\langle vars, chosen, allInput \rangle}$. To prove *R2a*, we make *HInv* strong enough to satisfy:

THEOREM *I1* $Hinit \Rightarrow HInv$

THEOREM *I2* $HInv \wedge HNEXT \Rightarrow HInv'$

A predicate *HInv* satisfying *I2* is said to be an invariant of the action *HNEXT*. A standard TLA theorem asserts that *I1* and *I2* imply *R2a*.

There are two general approaches to defining *HInv*. In both, we write *HInv* as a conjunction $HI_1 \wedge \dots \wedge HI_k$. In the bottom-up method, we define the HI_i in increasing order of i , so that each conjunction $HI_1 \wedge \dots \wedge HI_k$ is an invariant of *HNEXT*. We stop when we obtain an invariant strong enough to prove *R2b*. In the top-down method, we start by defining HI_k so that *R2b* is satisfied with HI_k substituted for *HInv*. We then define the HI_i in decreasing order of i so that $HI_i \wedge \dots \wedge HI_k \wedge HNEXT \Rightarrow HI'_{i+1}$, stopping when we obtain an invariant of *HNEXT*. In practice, one uses a combination of the two methods—with a lot of backtracking. Here, we present the invariant in a bottom-up fashion.

If the set of disks is empty, then *IsMajority*(D) is false for all subsets D of *Disk*. (This follows from the assumption about *IsMajority* by substituting D for both S and T .) Hence, *HDiskSynodSpec* implies that the system remains forever in its initial state, trivially satisfying *ISpec*. It therefore suffices to consider only the case when *Disk* is nonempty:

ASSUME $Disk \neq \{\}$

⁸The symbols *IInit* and *INEXT* are not defined in the current context; to be rigorous, we should define them to equal $IS(chosen, allInput)!IInit$ and $IS(chosen, allInput)!INEXT$, respectively.

The standard starting point for a TLA proof is a simple “type invariant”, which we call $HInv1$, asserting that all variables have the correct type:

$$\begin{aligned}
HInv1 &\triangleq \\
&\wedge \textit{input} \in [\textit{Proc} \rightarrow \textit{Inputs}] \\
&\wedge \textit{output} \in [\textit{Proc} \rightarrow \textit{Inputs} \cup \{\textit{NotAnInput}\}] \\
&\wedge \textit{disk} \in [\textit{Disk} \rightarrow [\textit{Proc} \rightarrow \textit{DiskBlock}]] \\
&\wedge \textit{phase} \in [\textit{Proc} \rightarrow 0 \dots 3] \\
&\wedge \textit{dblock} \in [\textit{Proc} \rightarrow \textit{DiskBlock}] \\
&\wedge \textit{output} \in [\textit{Proc} \rightarrow \textit{Inputs} \cup \{\textit{NotAnInput}\}] \\
&\wedge \textit{disksWritten} \in [\textit{Proc} \rightarrow \text{SUBSET } \textit{Disk}] \\
&\wedge \textit{blocksRead} \in [\textit{Proc} \rightarrow [\textit{Disk} \rightarrow \\
&\hspace{10em} \text{SUBSET } [\textit{block} : \textit{DiskBlock}, \textit{proc} : \textit{Proc}]]] \\
&\wedge \textit{allInput} \in \text{SUBSET } \textit{Inputs} \\
&\wedge \textit{chosen} \in \textit{Inputs} \cup \{\textit{NotAnInput}\} \\
&\wedge \textit{input} \in [\textit{Proc} \rightarrow \textit{Inputs}]
\end{aligned}$$

Our first lemma asserts that $HInv1$ is an invariant of $HNext$:

$$\text{LEMMA } I2a \quad HInv1 \wedge HNext \Rightarrow HInv1'$$

The proofs of Theorem $R2b$ and of most lemmas appear in Section A.4 below.

Before going any further, we define some useful state functions. First, we let $MajoritySet$ be the set of all subsets of the set of disks containing a majority of them; we let $blocksOf(p)$ be the set of all copies of p 's disk blocks in the system—that is, $dblock[p]$, p 's blocks on disk, and all blocks of p read by some processor; and we let $allBlocks$ be the set of all copies of all disk blocks of all processors.

$$\begin{aligned}
MajoritySet &\triangleq \{D \in \text{SUBSET } \textit{Disk} : IsMajority(D)\} \\
blocksOf(p) &\triangleq \\
&\text{LET } rdBy(q, d) \triangleq \{br \in \textit{blocksRead}[q][d] : br.proc = p\} \\
&\text{IN } \{dblock[p]\} \cup \{disk[d][p] : d \in \textit{Disk}\} \\
&\quad \cup \{br.block : br \in \text{UNION } \{rdBy(q, d) : q \in \textit{Proc}, d \in \textit{Disk}\}\} \\
allBlocks &\triangleq \text{UNION } \{blocksOf(p) : p \in \textit{Proc}\}
\end{aligned}$$

The next conjunct of $HInv$ describes some simple relations between the values of the different variables.

$$\begin{aligned}
HInv2 &\triangleq \\
&\wedge \forall p \in Proc : \\
&\quad \forall bk \in blocksOf(p) : \wedge bk.mbal \in Ballot(p) \cup \{0\} \\
&\quad \quad \wedge bk.bal \in Ballot(p) \cup \{0\} \\
&\quad \quad \wedge (bk.bal = 0) \equiv (bk.inp = NotAnInput) \\
&\quad \quad \wedge bk.mbal \geq bk.bal \\
&\wedge \forall p \in Proc, d \in Disk : \\
&\quad \wedge (d \in disksWritten[p]) \Rightarrow \wedge phase[p] \in \{1, 2\} \\
&\quad \quad \wedge disk[d][p] = dblock[p] \\
&\quad \wedge (phase[p] \in 1, 2) \Rightarrow \wedge (blocksRead[p][d] \neq \{\}) \Rightarrow \\
&\quad \quad \quad (d \in disksWritten[p]) \\
&\quad \quad \wedge \neg hasRead(p, d, p) \\
&\wedge \forall p \in Proc : \\
&\quad \wedge (phase[p] = 0) \Rightarrow \wedge dblock[p] = InitDB \\
&\quad \quad \wedge disksWritten[p] = \{\} \\
&\quad \quad \wedge \forall d \in Disk : \forall br \in blocksRead[p][d] : \\
&\quad \quad \quad \wedge br.proc = p \\
&\quad \quad \quad \wedge br.block = disk[d][p] \\
&\quad \wedge (phase[p] \neq 0) \Rightarrow \wedge dblock[p].mbal \in Ballot(p) \\
&\quad \quad \wedge dblock[p].bal \in Ballot(p) \cup \{0\} \\
&\quad \quad \wedge \forall d \in Disk : \forall br \in blocksRead[p][d] : \\
&\quad \quad \quad br.block.mbal < dblock[p].mbal \\
&\quad \wedge (phase[p] \in \{2, 3\}) \Rightarrow (dblock[p].bal = dblock[p].mbal) \\
&\quad \wedge output[p] = \text{IF } phase[p] = 3 \text{ THEN } dblock[p].inp \text{ ELSE } NotAnInput \\
&\wedge chosen \in allInput \cup \{NotAnInput\} \\
&\wedge \forall p \in Proc : \wedge input[p] \in allInput \\
&\quad \quad \wedge (chosen = NotAnInput) \Rightarrow (output[p] = NotAnInput)
\end{aligned}$$

The invariance of $HInv1 \wedge HInv2$ follows from Lemma *I2a* and:

LEMMA *I2b* $HInv1 \wedge HInv2 \wedge HNext \Rightarrow HInv2'$

The next conjunct of $HInv$ expresses the observation that if processors p and q have each read the other's block from disk d during their current phases, then at least one of them has read the other's current block.

$$\begin{aligned}
HInv3 &\triangleq \forall p, q \in Proc, d \in Disk : \\
&\quad \wedge phase[p] \in \{1, 2\} \\
&\quad \wedge phase[q] \in \{1, 2\} \\
&\quad \wedge hasRead(p, d, q) \\
&\quad \wedge hasRead(q, d, p)
\end{aligned}$$

$$\begin{aligned} \Rightarrow & \vee [block \mapsto dblock[q], proc \mapsto q] \in blocksRead[p][d] \\ & \vee [block \mapsto dblock[p], proc \mapsto p] \in blocksRead[q][d] \end{aligned}$$

LEMMA I2c $HInv1 \wedge HInv3 \wedge HNext \Rightarrow HInv3'$

The next conjunct of the invariant expresses relations among the *mbal* and *bal* values of a processor and of its disk blocks. Its first conjunct asserts that, when p is not recovering from a failure, its *mbal* value is at least as large as the *bal* field of any of its blocks, and at least as large as the *mbal* field of its block on some disk in any majority set. Its second conjunct asserts that, in phase 1, its *mbal* value is actually greater than the *bal* field of any of its blocks. Its third conjunct asserts that, in phase 2, its *bal* value is the *mbal* field of all its blocks on some majority set of disks. The fourth conjunct asserts that the *bal* field of any of its blocks is at most as large as the *mbal* field of all its disk blocks on some majority set of disks.

$$\begin{aligned} HInv4 & \triangleq \\ \forall p \in Proc : & \\ \wedge (phase[p] \neq 0) \Rightarrow & \\ \wedge \forall bk \in blocksOf(p) : dblock[p].mbal \geq bk.bal & \\ \wedge \forall D \in MajoritySet : & \\ \exists d \in D : \wedge dblock[p].mbal \geq disk[d][p].mbal & \\ \wedge dblock[p].bal \geq disk[d][p].bal & \\ \wedge (phase[p] = 1) \Rightarrow (\forall bk \in blocksOf(p) : dblock[p].mbal > bk.bal) & \\ \wedge (phase[p] \in \{2, 3\}) \Rightarrow & \\ (\exists D \in MajoritySet : \forall d \in D : disk[d][p].mbal = dblock[p].bal) & \\ \wedge \forall bk \in blocksOf(p) : & \\ \exists D \in MajoritySet : \forall d \in D : disk[d][p].mbal \geq bk.bal & \end{aligned}$$

LEMMA I2d $HInv1 \wedge HInv2 \wedge HInv2' \wedge HInv4 \wedge HNext \Rightarrow HInv4'$

Before going further, we define $maxBalInp(b, v)$ to assert that every block in *allBlocks* with *bal* field at least b has *inp* field v .

$$maxBalInp(b, v) \triangleq \forall bk \in allBlocks : (bk.bal \geq b) \Rightarrow (bk.inp = v)$$

We now come to a conjunct of *HInv* that provides some high-level insight into why the algorithm is correct. It asserts that, if a processor p is in phase 2, then either its *bal* and *inp* values satisfy $maxBalInp$, or else p must eventually abort its current ballot. Processor p will eventually abort its ballot if there is some processor q and majority set D such that p has not read q 's block on any disk in D , and all of those blocks have *mbal* values greater than $dblock[p].bal$. (Since p must read at least one of those disks, it must eventually read one of those blocks and abort.)

$$\begin{aligned}
HIInv5 &\triangleq \\
&\forall p \in Proc : \\
&\quad (phase[p] = 2) \Rightarrow \vee maxBallInp(dblock[p].bal, dblock[p].inp) \\
&\quad \vee \exists D \in MajoritySet, q \in Proc : \\
&\quad \quad \forall d \in D : \wedge disk[d][q].mbal > dblock[p].bal \\
&\quad \quad \wedge \neg hasRead(p, d, q)
\end{aligned}$$

LEMMA I2e

$$HIInv1 \wedge HIInv2 \wedge HIInv2' \wedge HIInv3 \wedge HIInv4 \wedge HIInv5 \wedge HNext \Rightarrow HIInv5'$$

Before defining our final conjunct, we define a predicate $valueChosen(v)$ that is true if v is the only possible value that can be chosen as an output. It asserts that there is some ballot number b such that $maxBallInp(b, v)$ is true. This condition is satisfied if there is no block bk in $allBlocks$ with $bk.bal \geq b$. So, $valueChosen(v)$ must require that some processor p has written blocks with bal field at least b to a majority set D of the disks. (By $maxBallInp(b, v)$, those blocks must have inp field v). We also ensure that, once $valueChosen(v)$ becomes true, it can never be made false. This requires the additional condition that no processor q that is currently executing phase 1 with $mbal$ value at least b can fail to see those blocks that p has written. So, $valueChosen(v)$ also asserts that, for every disk d in D , if q has already read $disk[d][p]$, then it has read a block with bal field at least b .

$$\begin{aligned}
valueChosen(v) &\triangleq \\
&\exists b \in \text{UNION } \{Ballot(p) : p \in Proc\} : \\
&\quad \wedge maxBallInp(b, v) \\
&\quad \wedge \exists p \in Proc, D \in MajoritySet : \\
&\quad \quad \forall d \in D : \wedge disk[d][p].bal \geq b \\
&\quad \quad \wedge \forall q \in Proc : \\
&\quad \quad \quad \wedge phase[q] = 1 \\
&\quad \quad \quad \wedge dblock[q].mbal \geq b \\
&\quad \quad \quad \wedge hasRead(q, d, p) \\
&\quad \quad \Rightarrow (\exists br \in blocksRead[q][d] : br.bal \geq b)
\end{aligned}$$

It's obvious that, if $valueChosen(v) = valueChosen(w)$, then $v = w$.

The final conjunct of $HIInv$ asserts that, once an output has been chosen, $valueChosen(chosen)$ holds, and each processor's output equals either $chosen$ or $NotAnInput$.

$$\begin{aligned}
HIInv6 &\triangleq \wedge (chosen \neq NotAnInput) \Rightarrow valueChosen(chosen) \\
&\quad \wedge \forall p \in Proc : output[p] \in \{chosen, NotAnInput\}
\end{aligned}$$

LEMMA I2f $HIInv1 \wedge HIInv2 \wedge HIInv2' \wedge HIInv3 \wedge HIInv6 \wedge HNext \Rightarrow HIInv6'$

We define $HInv$ to be the conjunction of $HInv1$ – $HInv6$.

$$HInv \triangleq HInv1 \wedge HInv2 \wedge HInv3 \wedge HInv4 \wedge HInv5 \wedge HInv6$$

Theorem $I2$ then follows easily from Lemmas $I2a$ – $I2f$.

A.4 Proofs

We now sketch the proofs of most of the lemmas from Section A.3 and of Theorem $R2b$. We give hierarchically structured proofs [9]. A structured proof consists of a sequence of statements and their proofs; each of those proofs is either a structured proof or an ordinary paragraph-style proof. The j^{th} step in the current level- i proof is numbered $\langle i \rangle j$. Within a paragraph-style proof, $\langle i \rangle j$ denotes the most recent statement with that number. The proof statement “ $\langle i \rangle j$. Q.E.D.” denotes the current goal—that is, the level $i - 1$ statement being proved by this step. A proof statement

ASSUME: A
 PROVE: P

asserts that the assumption A implies P . If P is the current goal, then this is abbreviated as

CASE: A

An assumption $\text{CONSTANT } c \in S$ asserts that c is a new constant parameter that we assume is in S . We prove $\forall c \in S : P(c)$ by proving

ASSUME: $\text{CONSTANT } c \in S$
 PROVE: $P(c)$

The assumption $\text{CONSTANT } c \in S \text{ s.t. } A(c)$ also assumes that c also satisfies $A(c)$. A proof statement

$\langle i \rangle j$ CHOOSE $c \in S \text{ s.t. } P(c)$

asserts the existence of a value c in S satisfying $P(c)$, and defines c to be such a value. To prove this statement, we must demonstrate the existence of c .

We recommend that proofs be read hierarchically, from the top level down. To read the proof of a long level- i step, you should first read the level- $(i + 1)$ statements that form its proof, together with the proof of the final “Q.E.D.” step (which is usually a short paragraph), and then read the proof of each level- $(i + 1)$ step, in any desired order.

We also use a hierarchical scheme for naming subformulas of a formula. If F is the name of a formula that is a conjunction, then $F.i$ is the name of its i^{th} conjunct. A similar scheme is used for a disjunction, except using letters instead of numbers, so $F.c$ is the name of the third disjunct of F . If F is the name of the formula $P \Rightarrow Q$, then $F.L$ is the name of P and $F.R$ is the name of Q . If F is the name of the formula $\exists x : P(x)$ or $\forall x : P(x)$, then $F(e)$ is the name of the formula $P(e)$, for any expression e . This is generalized in the obvious way for abbreviated quantifications like $\exists x, y : P(x, y)$. For example, $HInv5(n).R.b(E, m)(dd).2$ is the formula $\neg hasRead(n, dd, m)$.

We now give the proofs. We omit the proofs of Lemmas *I2a* and *I2b*, which require a simple but tedious case analysis for the different disjuncts of *Next*. In the informal paragraph-style proofs, we use *HInv1* implicitly in many places by tacitly assuming that variables have values of the right type. For example, we deduce $phase'[p] = 2$ from

$$phase' = [phase \text{ EXCEPT } ![p] = 2]$$

without mentioning that this follows only if $phase$ is a function whose domain contains p , which is implied by *HInv1*.

A.4.1 Lemma I2c

We prove Lemma *I2c* by proving:

ASSUME: 1. $HInv1 \wedge HInv3 \wedge HNext$
 2. CONSTANTS $p, q \in Proc, d \in Disk$
 3. $HInv3(p, q, d).L'$

PROVE: $HInv3(p, q, d).R'$

(1)1. CASE: $\neg HInv3(p, q, d).L$

(2)1. CASE: $Phase1or2Read(p, d, q)$

PROOF: Action $Phase1or2Read(p, d, q)$ adds the record

$[block \mapsto dblock[q], proc \mapsto q]$

to $blocksRead[p][d]$, making $HInv3(p, q, d).R.a'$ true.

(2)2. CASE: $Phase1or2Read(q, d, p)$

PROOF: Action $Phase1or2Read(q, d, p)$ adds the record

$[block \mapsto dblock[p], proc \mapsto p]$

to $blocksRead[q][d]$, making $HInv3(p, q, d).R.b'$ true.

(2)3. CASE: $EndPhase0(p)$

PROOF: This implies $\neg hasRead(p, d, q)$, so $HInv3(p, q, d).L'$ is false, making $HInv3(p, q, d).R'$ true.

(2)4. CASE: $EndPhase0(q)$

PROOF: This implies $\neg hasRead(q, d, p)$, so $HInv3(p, q, d).L'$ is false, making $HInv3(p, q, d)'$ true.

$\langle 2 \rangle 5$. Q.E.D.

PROOF: By assumption 3 and the level $\langle 1 \rangle$ case assumption, one of the four conjuncts of $HInv3(p, q, d).L$ is changed from false to true. Steps $\langle 2 \rangle 1$ – $\langle 2 \rangle 4$ covers the four subactions of $Next$ that can make one of those conjuncts false.

$\langle 1 \rangle 2$. CASE: $HInv3(p, q, d).L$

PROOF: By $HInv3$ (which holds by assumption 1), the case assumption implies $HInv3(p, q, d).R$. The only subactions of $HNext$ that make $HInv3(p, q, d).R$ false are ones that remove elements from $blocksRead[p][d]$ or $blocksRead[q][d]$ or that change $dblock[p]$ or $dblock[q]$. The only such subactions are ones with an $InitializePhase(p)$ or $InitializePhase(q)$ conjunct, which make $HInv3(p, q, d).R'$ false, contrary to assumption 3.

$\langle 1 \rangle 3$. Q.E.D.

PROOF: Immediate from steps $\langle 1 \rangle 1$ and $\langle 1 \rangle 2$.

A.4.2 Lemma BksOf

We now state and prove a simple result that will be used below.

LEMMA *BksOf*

$Next \wedge HInv1 \Rightarrow$

$(\forall p \in Proc :$

$blocksOf(p)' \subseteq (blocksOf(p) \setminus \{dblock[p]\}) \cup \{dblock'[p]\})$

It is proved as follows.

ASSUME: $p \in Proc$

PROVE: $blocksOf(p)' \subseteq (blocksOf(p) \setminus \{dblock[p]\}) \cup \{dblock'[p]\}$

PROOF: The only way an $HNext$ step creates a new block for p , rather than copying an existing one, is by changing $dblock[p]$.

A.4.3 Lemma I2d

ASSUME: 1. $HInv1 \wedge HInv2 \wedge HInv2' \wedge HInv4 \wedge HNext$

2. CONSTANT $p \in Proc$

PROVE: $HInv4(p)'$

$\langle 1 \rangle 1$. $HInv4(p).1'$

$\langle 2 \rangle 1$. CASE: $(phase[p] = 0) \wedge (phase'[p] \neq 0)$

$\langle 3 \rangle 1$. $EndPhase0(p)$

PROOF: By the level ⟨2⟩ case assumption, since $EndPhase0(p)$ is the only subaction of $HNext$ that changes $phase[p]$ from zero to a nonzero value.

⟨3⟩2. ASSUME: CONSTANT $bk \in blocksOf(p)'$ s.t. $bk \neq dblock'[p]$

PROVE: $dblock'[p].mbal \geq bk.bal$

⟨4⟩1. $bk \in blocksOf(p)$

PROOF: Lemma $BksOf$ and the level ⟨3⟩ assumption.

⟨4⟩2. CHOOSE $D1 \in MajoritySet$ s.t.

$\forall d \in D : disk[d][p].mbal \geq bk.bal$

PROOF: $HInv4.4$ and ⟨4⟩1 imply the existence of $D1$.

⟨4⟩3. $\forall D \in MajoritySet : \exists d \in D : disk[d][p].mbal \geq bk.bal$

PROOF: Using ⟨4⟩2, for any majority set D we can choose d to be a disk in $D1 \cap D$, which is nonempty because any two majority sets have an element in common.

⟨4⟩4. $\exists d \in Disk : \exists rb \in readBlock[p][d] : rb.block.mbal \geq bk.bal$

PROOF: By $HInv2.3(p).1.R.3$, which holds by assumption 1 and case assumption ⟨2⟩, $hasRead(p, d, p)$ implies that $blocksRead[p][d]$ consists of a single element whose $block$ field equals $disk[d][p]$, for any disk d . Step ⟨3⟩1 implies that $hasRead(p, d, p)$ is true for all d in some majority set D of disks. The level ⟨3⟩ goal then follows from ⟨4⟩3.

⟨4⟩5. Q.E.D.

PROOF: ⟨4⟩4 and ⟨3⟩1 imply $dblock'[p].mbal > bk.bal$.

⟨3⟩3. $HInv4(p).1.R.2'$

⟨4⟩1. $\exists D \in MajoritySet :$

$\forall d \in D : \wedge dblock'[p].mbal > disk[d][p].mbal$
 $\wedge dblock'[p].bal \geq disk[d][p].bal$

PROOF: By ⟨3⟩1, $dblock'[p].mbal > br.mbal$ and $dblock'[p].bal \geq br.bal$ for all $br \in allBlocksRead(p)$. By ⟨3⟩1, the level ⟨2⟩ case assumption, and $HInv2.3(p).1$, $allBlocksRead(p)$ contains all blocks $disk[d][p]$ for d in some majority set D of disks.

⟨4⟩2. Q.E.D.

PROOF: $HInv4(p).1.R.2'$ follows from ⟨4⟩1 and ⟨3⟩1, which implies that $disk$ is unchanged.

⟨3⟩4. Q.E.D.

PROOF: By ⟨3⟩2 and ⟨3⟩3, since ⟨3⟩2 implies $HInv4(p).1.R.1(bk)'$ except for the case $bk = dblock'[p]$; and $HInv4(p).1.R.1(bk)'$ follows from ⟨3⟩1 and $HInv2.1(p).4$ in that case.

⟨2⟩2. CASE: $(phase[p] \neq 0) \wedge (phase'[p] \neq 0)$

⟨3⟩1. $\wedge dblock'[p].mbal \geq dblock[p].mbal$
 $\wedge dblock'[p].bal \geq dblock[p].bal$

PROOF: The only subactions of *Next* that change $dblock[p]$ are

StartBallot(p), *EndPhase1or2(p)*, *EndPhase0(p)*, and *Fail(p)*

Of these, only *Fail(p)* can decrease $dblock[p].mbal$ or $dblock[p].bal$.

(*HInv2.1(p).4* implies that $EndPhase1or2(p) \wedge phase[p] = 1$ cannot decrease $dblock[p].bal$.) The level ⟨2⟩ case assumption implies $\neg Fail(p)$.

⟨3⟩2. *HInv4(p).1.R.1'*

PROOF: If $bk \in blocksOf(p)$, then $HInv4(p).1.R.1(bk)'$ follows from ⟨3⟩1 and *HInv4(p).1.R.1* (which holds by assumption 1 and the level ⟨2⟩ case assumption). If $bk = dblock'[p]$, then $HInv4(p).1.R.1(bk)'$ follows from *HInv2.1(p)(bk).4'*. We then obtain *HInv4(p).1.R.1'* from Lemma *BksOf*.

⟨3⟩3. *HInv4(p).1.R.2'*

PROOF: *HNext* implies that $disk'[p][d]$ equals $disk[p][d]$ or $dblock[p]$, so *HInv4(p).1.R.2'* follows from ⟨3⟩1 and *HInv4(p).1.R.1*, which holds by assumption 1 and the level ⟨2⟩ case assumption.

⟨3⟩4. Q.E.D.

PROOF: By ⟨3⟩2 and ⟨3⟩3.

⟨2⟩3. Q.E.D.

PROOF: By ⟨2⟩1 and ⟨2⟩2, since $HInv4(p).1'$ is trivially true if $phase'[p]$ equals 0.

⟨1⟩2. *HInv4(p).2'*

⟨2⟩1. CASE: $(phase[p] \neq 1) \wedge (phase'[p] = 1)$

⟨3⟩1. CASE: $phase[p] = 0$

⟨4⟩1. *EndPhase0(p)*

PROOF: By the levels ⟨2⟩ and ⟨3⟩ case assumptions.

⟨4⟩2. $\forall bk \in blocksOf(p) :$

$\exists D \in MajoritySet : \forall d \in D : disk[d][p].mbal \geq bk.bal$

PROOF: By *HInv4(p).4*.

⟨4⟩3. $\forall bk \in blocksOf(p) :$

$\forall D \in MajoritySet : \exists d \in D : disk[d][p].mbal \geq bk.bal$

PROOF: By ⟨4⟩2, since any two majority sets have a disk in common.

⟨4⟩4. $\forall bk \in blocksOf(p) :$

$\exists br \in allBlocksRead(p) : br.mbal \geq bk.bal$

PROOF: By ⟨4⟩1, ⟨4⟩3, and *HInv2.3(p).1.R.3* (which holds by assumption 1 and the level ⟨3⟩ case assumption).

⟨4⟩5. Q.E.D.

PROOF: $\langle 4 \rangle 1$ implies
 $\forall br \in allBlocksRead(p) : dblock'[p].mbal > br.mbal$
Therefore, $HInv4(p).2.R(bk)'$ follows from $\langle 4 \rangle 4$ and $\langle 4 \rangle 1$ if $bk \in blocksOf(p)$. Step $\langle 4 \rangle 1$ also implies
 $\exists bk \in blocksOf(p) : dblock'[p].bal = bk.bal$
so $HInv4(p).2.R(bk)'$ follows from $\langle 4 \rangle 4$ if $bk = dblock'[p]$. By Lemma *BksOf*, this proves $HInv4(p).2.R'$.

$\langle 3 \rangle 2$. CASE: $phase[p] \in \{2, 3\}$
 $\langle 4 \rangle 1$. $\forall bk \in blocksOf(p) : dblock[p].mbal \geq bk.bal$
PROOF: $HInv4(p).1$ and the level $\langle 3 \rangle$ case assumption (which imply $HInv4(p).1.R.1$).
 $\langle 4 \rangle 2$. $\wedge dblock'[p].mbal > dblock[p].mbal$
 $\wedge dblock'[p].bal = dblock[p].bal$
PROOF: By *HNext* and the level $\langle 2 \rangle$ and $\langle 3 \rangle$ case assumptions, which imply *StartBallot*(p).
 $\langle 4 \rangle 3$. Q.E.D.
PROOF: $\langle 4 \rangle 1$ and $\langle 4 \rangle 2$ imply $HInv4(p).2.R(bk)'$ for $bk = dblock'[p]$ and $bk \in blocksOf(p)$. Lemma *BksOf* then implies $HInv4(p).2.R'$.

$\langle 3 \rangle 3$. Q.E.D.
PROOF: The level $\langle 2 \rangle$ case assumption implies that $\langle 3 \rangle 1$ and $\langle 3 \rangle 2$ cover all possibilities.

$\langle 2 \rangle 2$. CASE: $(phase[p] = 1) \wedge (phase'[p] = 1)$
PROOF: By *HNext*, this implies $dblock'[p] = dblock[p]$, so Lemma *BksOf* implies that $HInv4(p).2'$ follows from $HInv4(p).2$.

$\langle 2 \rangle 3$. Q.E.D.
PROOF: Since $HInv4(p).2'$ is trivially true if $phase'[p] \neq 1$, the cases of $\langle 2 \rangle 1$ and $\langle 2 \rangle 2$ are exhaustive.

$\langle 1 \rangle 3$. $HInv4(p).3'$
 $\langle 2 \rangle 1$. CASE: $(phase[p] \neq 2) \wedge (phase'[p] = 2)$
 $\langle 3 \rangle 1$. $EndPhase1or2(p) \wedge (phase[p] = 1)$
PROOF: By *HNext* and the level $\langle 2 \rangle$ case assumption.
 $\langle 3 \rangle 2$. $\exists D \in MajoritySet : \forall d \in D : disk[d][p].mbal = dblock[p].mbal$
PROOF: By $\langle 3 \rangle 1$ and $HInv2.2(p).1$.
 $\langle 3 \rangle 3$. Q.E.D.
PROOF: $\langle 3 \rangle 1$ implies $dblock'[p].bal = dblock[p].mbal$ and $disk' = disk$, which by $\langle 3 \rangle 2$ implies $HInv4(p).3'$

$\langle 2 \rangle 2$. CASE: $(phase[p] \in \{2, 3\}) \wedge (phase'[p] \in \{2, 3\})$
 $\langle 3 \rangle 1$. $dblock'[p].bal = dblock[p].bal$
PROOF: By *HNext* and the level $\langle 2 \rangle$ case assumption.

⟨3⟩2. $\forall d \in \text{Disk}$:

$$\text{Phase1or2Write}(p, d) \Rightarrow (\text{disk}'[d][p].\text{mbal} = \text{dblock}[p].\text{bal})$$

PROOF: By the level ⟨2⟩ case assumption and $\text{HInv2.3}(p).3$.

⟨3⟩3. Q.E.D.

PROOF: $\text{HInv4}(p).3'$ follows from $\text{HInv4}(p).3$, ⟨3⟩1, and ⟨3⟩1, since $\text{HNext} \wedge \neg \text{Phase1or2Write}(p, d)$ implies $\text{disk}'[d][p] = \text{disk}[d][p]$, for any disk d .

⟨2⟩3. Q.E.D.

PROOF: $\text{HInv4}(p).3'$ follows from ⟨2⟩1 and ⟨2⟩2 because it is trivially true if $\text{phase}'[p] \notin \{2, 3\}$, and $\text{HNext} \wedge (\text{phase}'[p] = 3)$ implies $\text{phase}[p] \in \{2, 3\}$,

⟨1⟩4. $\text{HInv4}(p).4'$

⟨2⟩1. CASE: $\text{EndPhase1or2}(p) \wedge (\text{phase}[p] = 1)$

⟨3⟩1. $\exists D \in \text{MajoritySet} : \forall d \in D : \text{disk}'[d][p].\text{mbal} \geq \text{dblock}'[p].\text{bal}$

PROOF: By ⟨1⟩3 and the level ⟨2⟩ case assumption, which implies $\text{phase}'[p] = 2$.

⟨3⟩2. $\text{disk}' = \text{disk}$

PROOF: By the level ⟨2⟩ case assumption.

⟨3⟩3. Q.E.D.

PROOF: If $bk \neq \text{dblock}'[p]$, then $\text{HInv4}(p).4(bk)'$ follows from ⟨3⟩2 and $\text{HInv4}(p).4(bk)$. If $bk = \text{dblock}'[p]$, then it follows from ⟨3⟩1.

⟨2⟩2. CASE: $\text{Fail}(p)$

PROOF: $\text{Fail}(p)$ implies $\text{disk}' = \text{disk}$, so $\text{HInv4}(p).4(bk)'$ follows from $\text{HInv4}(p).4(bk)$ if $bk \neq \text{dblock}'[p]$. For $bk = \text{dblock}[p]$, $\text{HInv4}(p).4(bk)'$ is trivial because $\text{Fail}(p)$ implies $\text{dblock}[p].\text{bal} = 0$.

⟨2⟩3. CASE: $\exists d \in \text{Disk} : \text{Phase1or2Write}(p, d)$

PROOF: In this case, we have

$$\exists d \in \text{Disk} : [\text{disk}' = \text{disk} \text{ EXCEPT } ![d][p] = \text{dblock}[p]]$$

and $\text{phase}'[p] \neq 0$. From this, $\text{HInv4}(p).4$, and $\text{HInv4}(p).1.R.1'$ (which holds by ⟨1⟩1), we deduce $\text{HInv4}(p).4(bk)'$ for $bk \in \text{blocksOf}(p) \cap \text{blocksOf}(p)'$. For $bk = \text{dblock}'[p]$, we obtain $\text{HInv4}(p).4(bk)'$ from $\text{HInv2.1}(p)(bk).4'$. By Lemma BksOf , this proves $\text{HInv4}(p).4'$.

⟨2⟩4. Q.E.D.

PROOF: By ⟨2⟩1, ⟨2⟩2, ⟨2⟩3, since they consider all the subactions of HNext that change $b\text{Val}$ or $\text{disk}[d][p]$, for some disk d .

⟨1⟩5. Q.E.D.

PROOF: By steps ⟨1⟩1–⟨1⟩4.

A.4.4 Lemma I2e

Simple logic shows that, to prove Lemma I2e, it suffices to prove:

ASSUME: 1. $HInv1 \wedge HInv2 \wedge HInv2' \wedge HInv3 \wedge HInv4 \wedge HInv5 \wedge HNext$
 2. CONSTANT $p \in Proc$
 3. $phase'[p] = 2$
 4. $\neg HInv5(p).R.a'$

PROVE: $HInv5(p).R.b'$

$\langle 1 \rangle 1$. CASE: $(phase[p] \neq 2)$

$\langle 2 \rangle 1$. $EndPhase1or2(p) \wedge (phase[p] = 1)$

PROOF: By $HNext$, assumption 3, and the level $\langle 1 \rangle$ case assumption.

$\langle 2 \rangle 2$. CHOOSE $bk \in allBlocks$ s.t.

$(bk.bal \geq dblocks'[p].bal) \wedge (bk \neq dblocks')$

PROOF: Assumption 4 and the definition of $maxBalInp$ imply that there exists $bk \in allBlocks'$ such that $bk.bal \geq dblocks'[p].bal$ and $bk.inp \neq dblocks'[p].inp$. By Lemma $BksOf$ and the definition of $allBlocks$, $\langle 2 \rangle 1$ implies $bk \in allBlocks$.

$\langle 2 \rangle 3$. CHOOSE $q \in Proc \setminus \{p\}$ s.t. $bk \in blocksOf(q)$

PROOF: By $\langle 2 \rangle 2$ and the definition of $allBlocks$, there is some processor q such that $bk \in blocksOf(q)$. $\langle 2 \rangle 1$ and $\langle 2 \rangle 2$ imply $bk.bal \geq dblock[p].mbal$, so $\langle 2 \rangle 1$ and $HInv4(p).2$ imply $q \neq p$.

$\langle 2 \rangle 4$. $\exists D \in MajoritySet : \forall d \in D : disk[d][q].mbal \geq dblock'[p].bal$

PROOF: By $\langle 2 \rangle 3$, $HInv4(q).4$, and $\langle 2 \rangle 2$.

$\langle 2 \rangle 5$. $\exists D \in MajoritySet : \forall d \in D : disk[d][q].mbal > dblock'[p].bal$

PROOF: By $\langle 2 \rangle 3$ (which implies $p \neq q$) and $\langle 2 \rangle 4$, since $\langle 2 \rangle 1$ (which implies $dblock'[p].bal > 0$), $HInv2.1$, and the assumption that different processors have distinct ballot numbers imply $disk[d][q].mbal \neq dblock'[p].bal$.

$\langle 2 \rangle 6$. Q.E.D.

PROOF: $\langle 2 \rangle 1$ implies $\neg hasRead(p, d, q)'$, for all disks d . Hence, $\langle 2 \rangle 5$ implies $HInv5(p).R.b'$.

$\langle 1 \rangle 2$. CASE: $(phase[p] = 2) \wedge HInv5(p).R.a$

$\langle 2 \rangle 1$. CHOOSE $q \in Proc \setminus \{p\}$ s.t. $\wedge EndPhase1or2(q) \wedge (phase[q] = 1)$
 $\wedge dblock'[q].bal > dblock[p].bal$
 $\wedge dblock'[q].inp \neq dblock[p].inp$

PROOF: $HNext$, Assumption 3, and the level $\langle 1 \rangle$ case assumption imply that $dblock'[p] = dblock[p]$. Assumption 4, the level $\langle 1 \rangle$ case assump-

tion, the definition of $maxBalInp$, and Lemma $BlksOf$ imply

$$\begin{aligned} & \wedge dblock'[q].bal \geq dblock[p].bal \\ & \wedge dblock'[q].inp \neq dblock[p].inp \\ & \wedge (dblock'[q].bal \neq dblock[q].bal) \vee (dblock'[q].bal \geq dblock[q].bal) \end{aligned}$$

for some processor $q \neq p$. By $HNEXT$, this implies $EndPhase1or2(q) \wedge (phase[q] = 1)$. By $HInv2.1$ and the assumption that different processors have different ballot numbers, this also implies $dblock'[q].bal \neq dblock[p].bal$.

$\langle 2 \rangle 2$. CHOOSE $D \in MajoritySet$ S.T.

$$\begin{aligned} \forall d \in D : & \wedge disk[d][q].mbal > dblock[p].bal \\ & \wedge hasRead(q, d, p) \end{aligned}$$

PROOF: By $HInv2.2(q, d).1$ and $\langle 2 \rangle 1$, there is a majority set D such that $hasRead(q, d, p)$ and $disk[d][q].mbal = dblock'[q].bal$, for all $d \in D$. The result then follows from $\langle 2 \rangle 1$.

$\langle 2 \rangle 3$. $\forall d \in D : [block \mapsto dblock[p], proc \mapsto p] \notin blocksRead[q][d]$

PROOF: By the level $\langle 1 \rangle$ case assumption, $\langle 2 \rangle 1$, and the definitions of $maxBalInp$ and $EndPhase1or2$, if $dblock[p]$ were in $allBlocksRead(q)$, then $dblock'[q].inp$ would equal $dblock[p].inp$, contradicting $\langle 2 \rangle 1$.

$\langle 2 \rangle 4$. $\forall d \in D : \neg \exists br \in blocksRead[p][d] : br.block.mbar \geq dblock[p].bal$

PROOF: By the level $\langle 1 \rangle$ case assumption (which implies $phase[p] = 2$), $HInv2.3(p).2.R.3$, and $HInv2.3(p).3$.

$\langle 2 \rangle 5$. $\forall d \in D : \neg hasRead(p, d, q)$

PROOF: By $HInv3$, $\langle 2 \rangle 2$ (which implies $hasRead(q, d, p)$ for $d \in D$), the level $\langle 1 \rangle$ case assumption (which implies $phase[p] = 2$), and $\langle 2 \rangle 1$ (which implies $phase[q] = 1$), $hasRead(p, d, q)$ implies

$$dblock[q] \in allBlocksRead(p)$$

which is impossible by $\langle 2 \rangle 2$ and $\langle 2 \rangle 4$.

$\langle 2 \rangle 6$. Q.E.D.

PROOF: Since $\langle 2 \rangle 1$ implies that $disk$, $dblock[p].bal$ and $hasRead(p, d, q)$ are unchanged, for all $d \in Disk$, $\langle 2 \rangle 2$ and $\langle 2 \rangle 5$ imply $HInv5(p).R.b'$.

$\langle 1 \rangle 3$. CASE: $(phase[p] = 2) \wedge HInv5(p).R.b$

$\langle 2 \rangle 1$. CHOOSE $D \in MajoritySet$, $q \in Proc$ S.T.

$$(q \neq p) \wedge HInv5(p).R.b(D, q)$$

PROOF: The level $\langle 1 \rangle$ case assumption implies the existence of D and q satisfying $HInv5(p).R.b(D, q)$. Since any two majority sets have a disk in common, $HInv4(p).3$ then implies $q \neq p$.

$\langle 2 \rangle 2$. CASE: $\exists d \in D : Phase1or2Write(q, d)$

$\langle 3 \rangle 1$. $\exists d \in D : disk' = [disk \text{ EXCEPT } ![d][q] = dblock[q]]$

PROOF: By the level $\langle 2 \rangle$ case assumption.

$\langle 3 \rangle 2$. $dblock[q].mbal > dblock[p].bal$.

PROOF: By $\langle 2 \rangle 1$ and $HInv4(q).1.R(2)$.

$\langle 3 \rangle 3$. Q.E.D.

PROOF: $\langle 3 \rangle 1$, $\langle 3 \rangle 2$, and $\langle 2 \rangle 1$ imply $HInv5(p).R.b(D, q)'$.

$\langle 2 \rangle 3$. CASE: $\exists d \in D : Phase1or2Read(p, d, q)$

PROOF: In this case, $HInv5(p).R.b$ (the level $\langle 1 \rangle$ case assumption) implies $phase'[p] = 1$ (because the ballot must abort), contradicting assumption 3.

$\langle 2 \rangle 4$. Q.E.D.

PROOF: Assumption 3, the level $\langle 1 \rangle$ case assumption, and $HNext$ imply that $dblock[p]$ is unchanged; and $HNext$ implies that, for any $d \in D$:

$\wedge (disk'[d][q] \neq disk[d][q]) \Rightarrow Phase1or2Write(q, d)$

$\wedge hasRead(p, d, q)' \wedge \neg hasRead(p, d, q) \Rightarrow Phase1or2Read(p, d, q)$

Hence, $\langle 2 \rangle 2$ and $\langle 2 \rangle 3$ cover the only cases in which $HInv5(p).R.b(D, q)$ can be made false. In all other cases, $HInv5(p).R.b'$ follows from $\langle 2 \rangle 1$.

$\langle 1 \rangle 4$. Q.E.D.

PROOF: By $HInv5(p)$, the cases in steps $\langle 1 \rangle 1$, $\langle 1 \rangle 2$, and $\langle 1 \rangle 3$ are exhaustive.

A.4.5 Lemma I2f

The proof of Lemma I2f uses:

LEMMA VC $\forall v \in Inputs : HInv1 \wedge HInv4 \wedge HNext \wedge valueChosen(v) \Rightarrow valueChosen(v)'$

We prove Lemma VC by proving:

ASSUME: 1. CONSTANT $b \in \text{UNION } \{Ballot(p) : p \in Proc\}$
 2. CONSTANTS $v \in Inputs, p \in Proc, D \in MajoritySet$
 3. $maxBalInp(b, v)$
 4. $valueChosen(v)(b).2(p, D)$

PROVE: $maxBalInp(b, v)' \wedge valueChosen(v)(b).2(p, D)'$

$\langle 1 \rangle 1$. $maxBalInp(b, v)'$

$\langle 2 \rangle 1$. CASE: $\exists q \in Proc : EndPhase1or2(q) \wedge (phase[q] = 1)$

$\langle 3 \rangle 1$. CHOOSE $q \in Proc$ s.t. $EndPhase1or2(q) \wedge (phase[q] = 1)$

PROOF: q exists by the level $\langle 2 \rangle$ case assumption.

$\langle 3 \rangle 2$. CASE: $dblock[q].mbal \geq b$

$\langle 4 \rangle 1$. $\exists d \in D : hasRead(q, d, p)$

PROOF: By $\langle 3 \rangle 1$ (which implies $hasRead(q, d, p)$ for all d in some majority set), since any two majority sets have a disk in common.

$\langle 4 \rangle 2$. $\exists d \in D : \exists br \in blocksRead[q][d] : br.block.bal \geq b$

PROOF: By $\langle 4 \rangle 1$, case assumption $\langle 3 \rangle$, and assumption 4.

$\langle 4 \rangle 3$. $dblock'[q].inp = v$
PROOF: By $\langle 4 \rangle 2$, $maxBalInp(b, v)$ (assumption 3), $\langle 3 \rangle 1$, and the definition of $EndPhase1or2$.

$\langle 4 \rangle 4$. Q.E.D.
PROOF: $\langle 4 \rangle 3$, $maxBalInp(b, v)$ (assumption 3), $\langle 3 \rangle 1$ (which implies $blocksOf(r)' = blocksOf(r)$ for $r \neq q$), and Lemma $BlksOf$ imply $maxBalInp(b, v)'$.

$\langle 3 \rangle 3$. CASE: $dblock[q].mbal < b$
PROOF: By $\langle 3 \rangle 1$, this implies $dblock'[q].bal < b$, so $maxBalInp(b, v)$ (assumption 3), $\langle 3 \rangle 1$, and Lemma $BlksOf$ imply $maxBalInp(b, v)'$.

$\langle 3 \rangle 4$. Q.E.D.
PROOF: By $\langle 3 \rangle 2$ and $\langle 3 \rangle 3$.

$\langle 2 \rangle 2$. Q.E.D.
PROOF: By $\langle 2 \rangle 1$, since $HNext \wedge (allBlocks' \neq allBlocks)$ implies
 $\exists q \in Proc : \vee EndPhase1or2(q) \wedge (phase[q] = 1)$
 $\vee Fail(q)$
and $maxBalInp(b, v) \wedge Fail(q)$ obviously implies $maxBalInp(b, v)'$.

$\langle 1 \rangle 2$. $valueChosen(v)(b).2(p, D)'$
 $\langle 2 \rangle 1$. ASSUME: CONSTANT $d \in D$
PROVE: $disk'[d][p].bal \geq b$
 $\langle 3 \rangle 1$. CASE: $Phase1or2Write(p, d)$
 $\langle 4 \rangle 1$. $\exists dd \in D : dblock[p].bal \geq disk[dd][p].bal$
PROOF: By $HInv4(p).1.R.2$ (which holds because the level $\langle 3 \rangle$ case assumption implies $phase[p] \neq 0$) and assumption 2.
 $\langle 4 \rangle 2$. $dblock[p].bal \geq b$
PROOF: By $\langle 4 \rangle 1$ and assumption 4, which implies $disk[dd][p].bal \geq b$ for all $dd \in D$.
 $\langle 4 \rangle 3$. Q.E.D.
PROOF: By the level $\langle 3 \rangle$ case assumption, $disk'[d][p] = dblock[p]$, so $\langle 4 \rangle 2$ implies $disk'[d][p].bal \geq b$.

$\langle 3 \rangle 2$. CASE: $disk'[d][p] = disk[d][p]$
PROOF: In this case, assumption 4 and the level $\langle 2 \rangle$ assumption imply $disk'[d][p] \geq b$.

$\langle 3 \rangle 3$. Q.E.D.
PROOF: By $\langle 3 \rangle 1$ and $\langle 3 \rangle 2$, since:
 $HNext \wedge (disk'[d][p] \neq disk[d][p]) \Rightarrow Phase1or2Write(p, d)$

$\langle 2 \rangle 2$. ASSUME: 1. CONSTANTS $q \in Proc$, $d \in D$
2. $phase'[q] = 1$
3. $dblock'[q].mbal \geq b$

4. $hasRead(q, d, p)'$

PROVE: $\exists br \in blocksRead'[q][d] : br.block.bal \geq b$

$\langle 3 \rangle 1$. $phase[q] = 1$

PROOF: By the level $\langle 2 \rangle$ assumptions 2 and 4, since:

$HNext \wedge (phase'[q] \neq phase[q]) \Rightarrow InitializePhase(q)$

and $InitializePhase(q)$ implies $\neg hasRead(q, d, p)'$.

$\langle 3 \rangle 2$. $dblock'[q].mbal = dblock[q].mbal$

PROOF: By the level $\langle 2 \rangle$ assumption 4, since:

$HNext \wedge (dblock'[q] \neq dblock[q]) \Rightarrow InitializePhase(q)$

and $InitializePhase(q)$ implies $\neg hasRead(q, d, p)'$.

$\langle 3 \rangle 3$. CASE: $Phase1or2Read(q, d, p)$

PROOF: Assumption 4 implies and the level $\langle 2 \rangle$ assumption 2 imply $disk[d][p].bal \geq b$. The case assumption and the level $\langle 2 \rangle$ assumption 4 imply

$[block \mapsto disk[d][p], proc \mapsto p] \in blocksRead'[q][d]$

proving the level $\langle 2 \rangle$ goal.

$\langle 3 \rangle 4$. CASE: $\neg Phase1or2Read(q, d, p)$

$\langle 4 \rangle 1$. $hasRead(q, d, p)$

PROOF: By the level $\langle 3 \rangle$ case assumption and the level $\langle 2 \rangle$ assumption 4, since:

$HNext \wedge \neg hasRead(q, d, p) \wedge hasRead(q, d, p)'$

$\Rightarrow Phase1or2Read(q, d, p)$

$\langle 4 \rangle 2$. $\exists br \in blocksRead[q][d] : br.block.bal \geq b$

PROOF: $\langle 3 \rangle 1$, $\langle 3 \rangle 2$ and the level $\langle 2 \rangle$ assumption 3 (which imply $dblock[q].mbal \geq b$), and assumption 4.

$\langle 4 \rangle 3$. Q.E.D.

PROOF: By $\langle 4 \rangle 2$ and the level $\langle 2 \rangle$ assumption 4, since:

$HNext \wedge hasRead(q, d, p)' \Rightarrow$

$(blocksRead[q][d] \subseteq blocksRead[q][d]')$

$\langle 3 \rangle 5$. Q.E.D.

PROOF: By $\langle 3 \rangle 3$ and $\langle 3 \rangle 4$.

$\langle 2 \rangle 3$. Q.E.D.

PROOF: $\langle 2 \rangle 1$ and $\langle 2 \rangle 2$ imply $valueChosen(v)(b).2(p, D)'$.

$\langle 1 \rangle 3$. Q.E.D.

PROOF: By $\langle 1 \rangle 1$ and $\langle 1 \rangle 2$.

We now prove Lemma *I2f* by proving:

ASSUME: $HInv1 \wedge HInv2 \wedge HInv2' \wedge HInv3 \wedge HInv5 \wedge HInv6 \wedge HNext$

PROVE: $HInv6'$

$\langle 1 \rangle 1$. ASSUME: $chosen' \neq NotAnInput$

PROVE: $valueChosen(chosen)'$

⟨2⟩1. CASE: $chosen = NotAnInput$

⟨3⟩1. CHOOSE $p \in Proc$ s.t. $EndPhase1or2(p) \wedge (phase[p] = 2)$

PROOF: $HInv2.5$, $HInv2.5'$, and the level ⟨1⟩ and ⟨2⟩ assumptions imply the existence of a $p \in Proc$ such that:

$(output[p] = NotAnInput) \wedge (output'[p] \neq NotAnOutput)$

By $HNext$, this implies $EndPhase1or2(p) \wedge (phase[p] = 2)$.

⟨3⟩2. $maxBalInp(dblock[p].bal, dblock[p].inp)$

PROOF: ⟨3⟩1 implies

$\exists D \in MajoritySet : \forall d \in D, q \in Proc : hasRead(p, d, q)$

Since any two majority sets have a disk in common, this implies $\neg HInv5(p).R.b$. Hence, $HInv5$ and ⟨3⟩1 (which implies $phase[p] = 2$) imply $HInv5(p).R.a$.

⟨3⟩3. $maxBalInp(dblock[p].bal, chosen)'$

PROOF: ⟨3⟩1 implies

$(chosen' = dblock[p].inp) \wedge (dblock'[p].bal = dblock[p].bal)$

which by ⟨3⟩2 implies $maxBalInp(dblock'[p].bal, chosen')$. Lemma $BksOf$ and ⟨3⟩1 imply $maxBalInp(b, v)' = maxBalInp(b, v)$ for any constants b and v .

⟨3⟩4. CHOOSE $D \in MajoritySet$ s.t.

$\forall d \in D, q \in Proc : hasRead(p, d, q) \wedge (disk[d][p] = dblock[p])$

PROOF: D exists by ⟨3⟩1 and $HInv2.2(p, d).1$.

⟨3⟩5. ASSUME: CONSTANTS $q \in Proc, d \in D$ s.t.

$\wedge phase[q] = 1$

$\wedge dblock[q].mbal \geq dblock[p].bal$

$\wedge hasRead(q, d, p)$

PROVE: $[block \mapsto dblock[p], proc \mapsto p] \in blocksRead[q][d]$

PROOF: ⟨3⟩1 and $HInv2.3(p).3$ imply $dblock[p].bal = dblock[p].mbal$; $HInv2.3(p).2.R.3$ and the assumption $dblock[q].mbal \geq dblock[p].bal$ then imply

$[block \mapsto dblock[q], proc \mapsto q] \notin blocksRead[p][d]$

The result now follows from $HInv3$ and ⟨3⟩4.

⟨3⟩6. $\forall q \in Proc, d \in D :$

$\wedge phase'[q] = 1$

$\wedge dblock'[q].mbal \geq dblock[p].bal$

$\wedge hasRead(q, d, p)'$

$\Rightarrow (\exists br \in blocksRead'[q][d] : br.block.bal = dblock[p].bal)$

PROOF: By ⟨3⟩5, since ⟨3⟩1 implies that, if $p \neq q$, then $phase[q]$, $dblock[q]$, $hasRead(q, d, p)$, and $blocksRead$ are unchanged, for any disk d ; and that $phase'[q] = 1$ implies $p \neq q$.

- ⟨3⟩7. Q.E.D.
 PROOF: ⟨3⟩3 implies $valueChosen(chosen)'(dblock[p].bal).1$; ⟨3⟩6 implies $valueChosen(chosen)'(dblock[p].bal).2(p, D)$.
- ⟨2⟩2. CASE: $chosen \neq NotAnInput$
 ⟨3⟩1. $chosen' = chosen$
 PROOF: By $HNext$ and the level ⟨2⟩ case assumption.
- ⟨3⟩2. Q.E.D.
 PROOF: The level ⟨2⟩ case assumption and $HInv6$ imply $valueChosen(chosen)$
 By Lemma VC and ⟨3⟩1, this implies $valueChosen(chosen)'$.
- ⟨2⟩3. Q.E.D.
 PROOF: Immediate from ⟨2⟩1 and ⟨2⟩2.
- ⟨1⟩2. ASSUME: CONSTANT $p \in Proc$ s.t. $output'[p] \neq NotAnInput$
 PROVE: $output'[p] = chosen'$
- ⟨2⟩1. CASE: $chosen = NotAnInput$
 ⟨3⟩1. $\forall q \in Proc : output[q] = NotAnInput$
 PROOF: By $HInv2.5$ and the level ⟨2⟩ case assumption.
- ⟨3⟩2. Q.E.D.
 PROOF: ⟨3⟩1, the level ⟨2⟩ case assumption, and $HNext$ imply that if $output'[p] \neq NotAnInput$, then $chosen' = output'[p]$.
- ⟨2⟩2. CASE: $chosen \neq NotAnInput$
 ⟨3⟩1. $valueChosen(chosen)$
 PROOF: By the level ⟨2⟩ case assumption and $HInv6.1$.
- ⟨3⟩2. $valueChosen(chosen)'$
 PROOF: By ⟨1⟩1, since the level ⟨2⟩ case assumption and $HNext$ imply $chosen' \neq NotAnInput$.
- ⟨3⟩3. $chosen' = chosen$
 PROOF: By ⟨3⟩1, ⟨3⟩2, and Lemma VC , since $valueChosen(v)$ and $valueChosen(w)$ imply $v = w$.
- ⟨3⟩4. CASE: $output[p] = NotAnInput$
 ⟨4⟩1. $EndPhase1or2(p) \wedge (phase[p] = 2)$
 PROOF: By the level ⟨1⟩ assumption, the level ⟨3⟩ case assumption, and $HNext$.
- ⟨4⟩2. $\exists D \in MajoritySet : \forall q \in Proc : hasRead(p, d, q)$
 PROOF: By ⟨4⟩1
- ⟨4⟩3. $maxBalInp(dblock[p].bal, dblock[p].inp)$
 PROOF: By $HInv5(p)$ and ⟨4⟩1, since ⟨4⟩2 implies $\neg HInv5(p).R.b$ (because any two majority sets have a disk in common).

⟨4⟩4. $\exists bk \in allBlocks, b \in \text{UNION } \{Ballot(p) : p \in Proc\} :$
 $\wedge maxBalInp(b, chosen)$
 $\wedge bk.bal \geq b$

PROOF: By ⟨3⟩1 and the definition of *valueChosen*.

⟨4⟩5. $dblock[p].inp = chosen$

PROOF: By ⟨4⟩3, ⟨4⟩4, and the definition of *maxBalInp*.

⟨4⟩6. Q.E.D.

PROOF: ⟨3⟩3, ⟨4⟩1 (which implies $output'[p] = dblock[p].inp$), and
⟨4⟩5 imply $output'[p] = chosen'$.

⟨3⟩5. CASE: $output[p] \neq NotAnInput$

PROOF: In this case, *HInv2.3(p).4*, the level ⟨1⟩ assumption, and
HNext imply $output'[p] = output[p]$; and *HInv6.2* and ⟨3⟩3 imply
 $output'[p] = chosen'$.

⟨3⟩6. Q.E.D.

PROOF: By ⟨3⟩4 and ⟨3⟩5

⟨2⟩3. Q.E.D.

PROOF: By ⟨2⟩1 and ⟨2⟩2

⟨1⟩3. Q.E.D.

PROOF: *HInv6'* follows immediately from ⟨1⟩1 and ⟨1⟩2.

A.4.6 Theorem R2b

We now prove Theorem *R2b*. First, we define *IFail(p)* and *IChoose(p)*
to be the actions *Fail(p)* and *Choose(p)* from submodule *Inner* of module
SynodSpec (with *chosen* and *allInput* being the variables declared in the
current context).

LET: $IFail(p) \triangleq IS(chosen, allInput)!Fail(p)$

$IChoose(p) \triangleq IS(chosen, allInput)!IChoose$

ASSUME: $HInv \wedge HInv' \wedge HNext$

PROVE: $(\exists p \in Proc : IFail(p) \vee IChoose(p)) \vee (\text{UNCHANGED } ivars)$

⟨1⟩1. CASE: $\exists p \in Proc : Fail(p)$

PROOF: *Fail(p)* implies *IFail(p).1* and the existence of $ip \in Inputs$ such
that *IFail(p).2(ip).1*; and *HNext* then implies *IFail(p).2(ip).2*. We deduce
IFail(p).3 from *Fail(p).4*, *HNext* and *HInv2.5*.

⟨1⟩2. CASE: $\exists p \in Proc : (phase[p] = 2) \wedge EndPhase1or2$

⟨2⟩1. CHOOSE $p \in Proc$ S.T. $(phase[p] = 2) \wedge EndPhase1or2$

PROOF: p exists by the level ⟨1⟩ case assumption.

⟨2⟩2. CASE: $chosen = NotAnInput$

⟨3⟩1. $\forall q \in Proc : output[q] = NotAnInput$

PROOF: By the level $\langle 2 \rangle$ case assumption and *HInv2.5*.

$\langle 3 \rangle 2$. $\forall q \in Proc \setminus \{p\} : output[q] = NotAnInput$
PROOF: $\langle 3 \rangle 1$ and $\langle 2 \rangle 1$.

$\langle 3 \rangle 3$. $chosen' = output'[p]$
PROOF: By $\langle 3 \rangle 2$, $\langle 2 \rangle 1$ (which implies $output'[p] \neq NotAnInput$), and *HNext*.

$\langle 3 \rangle 4$. Q.E.D.
PROOF: $\langle 3 \rangle 1$ implies *IChoose*(p).1; $\langle 2 \rangle 1$, $\langle 3 \rangle 3$, the level $\langle 2 \rangle$ case assumption, and *HNext* imply *IChoose*(p).2; and $\langle 2 \rangle 1$ and *HInv2.5* imply *IChoose*(p).3.

$\langle 2 \rangle 3$. CASE: $chosen \neq NotAnInput$
 $\langle 3 \rangle 1$. $chosen' = chosen$
PROOF: By *HNext* and the level $\langle 2 \rangle$ case assumption.

$\langle 3 \rangle 2$. $output'[p] = chosen$
PROOF: By *HInv6'.2*, $\langle 2 \rangle 1$ (which implies $output'[p] \neq NotAnInput$), and $\langle 3 \rangle 1$.

$\langle 3 \rangle 3$. Q.E.D.
PROOF: $\langle 2 \rangle 1$ implies *IChoose*(p).1; $\langle 3 \rangle 1$, $\langle 3 \rangle 2$ and the level $\langle 2 \rangle$ case assumption imply *IChoose*(p).2; and $\langle 2 \rangle 1$, *HNext*, and *HInv2.5* imply *IChoose*(p).3.

$\langle 2 \rangle 4$. Q.E.D.
PROOF: By $\langle 2 \rangle 2$ and $\langle 2 \rangle 3$.

$\langle 1 \rangle 3$. Q.E.D.
PROOF: By $\langle 1 \rangle 1$ and $\langle 1 \rangle 2$, since
 $HInv2.5 \wedge HNext \wedge (ivars' \neq ivars) \Rightarrow$
 $(input' \neq input) \vee (output' \neq output)$
and
 $HNext \wedge ((input' \neq input) \vee (output' \neq output)) \Rightarrow$
 $\exists p \in Proc : Fail(p) \vee ((phase[p] = 2) \wedge EndPhase1or2)$