

WAIT-FREE REGISTERS

Ellis Michael

CONGRATS! YOU'RE NOW PAXOS EXPERTS!

DRAWBACKS OF PAXOS

DRAWBACKS OF PAXOS

- Leader is a single bottleneck, processes $O(n)$ messages on every request.

DRAWBACKS OF PAXOS

- Leader is a single bottleneck, processes $O(n)$ messages on every request.
- FLP means that liveness not guaranteed.

DRAWBACKS OF PAXOS

- Leader is a single bottleneck, processes $O(n)$ messages on every request.
- FLP means that liveness not guaranteed.
- More practically, Paxos can have bad availability during failure scenarios (e.g., if a leader fails, it takes time to elect a new one).

ALTERNATIVES

ALTERNATIVES

- Allow randomness (see Ben-Or lecture).

ALTERNATIVES

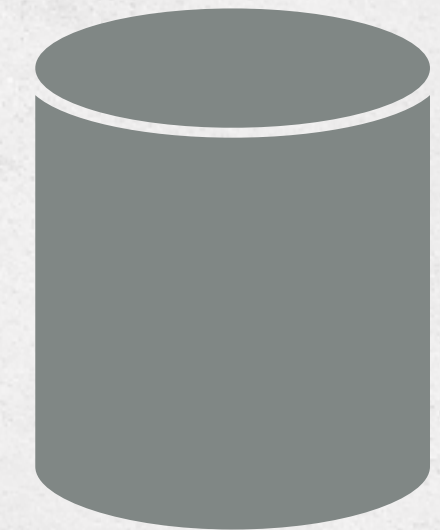
- Allow randomness (see Ben-Or lecture).
- Weaken the safety guarantees and accept weaker consistency (at your own peril).

ALTERNATIVES

- Allow randomness (see Ben-Or lecture).
- Weaken the safety guarantees and accept weaker consistency (at your own peril).
- Constrain the problem.

REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



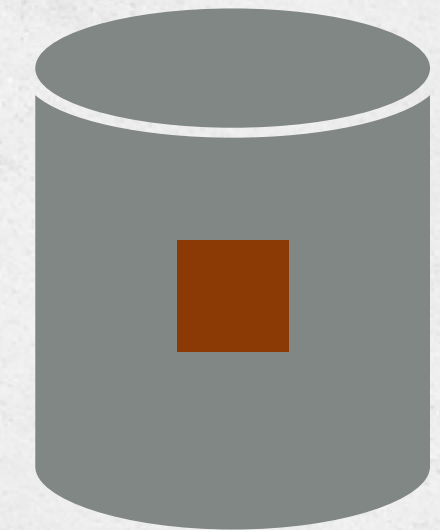
REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



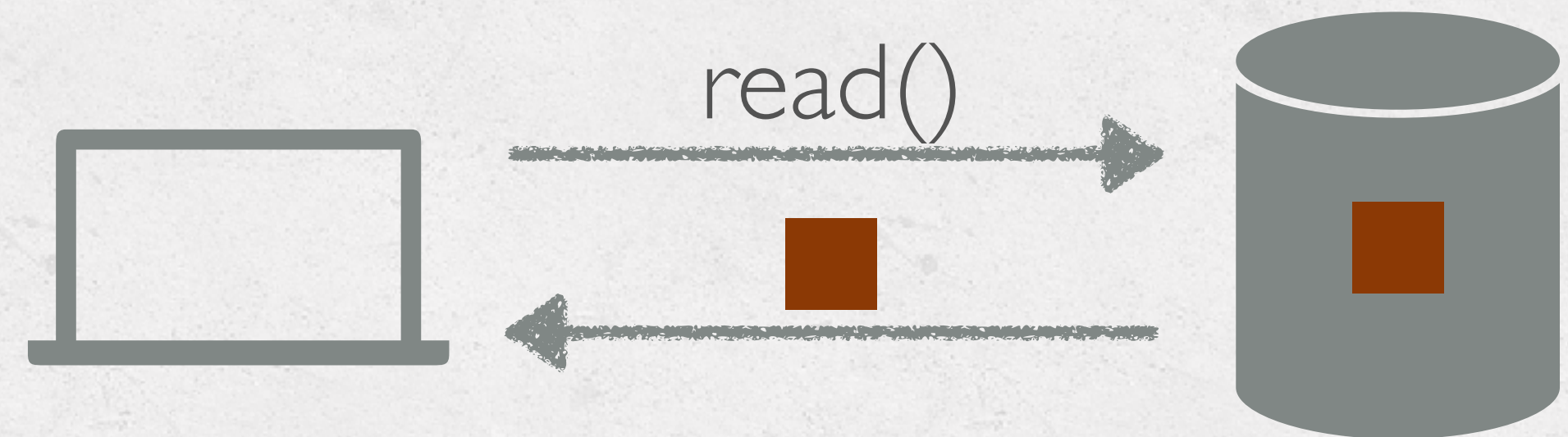
REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



REGISTERS

- Hold a single value. Want multiple values? Use multiple registers.
- Allows reads and writes only. Does not allow appends or other read-modify-write operations.
- Recall **safe**, **regular**, and **atomic/linearizable** semantics. We want *linearizability*.



WHY NO APPENDS?

Simple way to implement consensus:

WHY NO APPENDS?

Simple way to implement consensus:

- All processes append their input value.

WHY NO APPENDS?

Simple way to implement consensus:

- All processes append their input value.
- All processes read the value.

WHY NO APPENDS?

Simple way to implement consensus:

- All processes append their input value.
- All processes read the value.
- They all decide the first value that was appended.

WHY NO APPENDS?

Simple way to implement consensus:

- All processes can wait-free implement an appendable register, you can solve consensus (safety and liveness), which is impossible.
- All processes can wait-free implement an appendable register, you can solve consensus (safety and liveness), which is impossible.
- They all decide the first value that was appended.

IMPLEMENTING A REGISTER

- We will use the **client/server** model, where servers are replicas storing the value and clients send **reads** and **writes**.
- We want linearizability of reads and writes.
- As usual, we want to tolerate up to f server *crash failures*. Clients can also fail by crashing.

NON-BLOCKING ALGORITHMS

- **Lock-free** algorithms guarantee system-wide progress.
- **Wait-free** algorithms guarantee per-client progress. That is, no matter what steps other processes take, a correct client's operations are always completed in a finite number of steps.

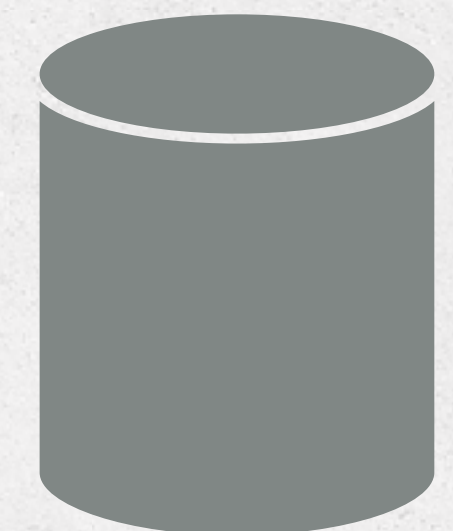
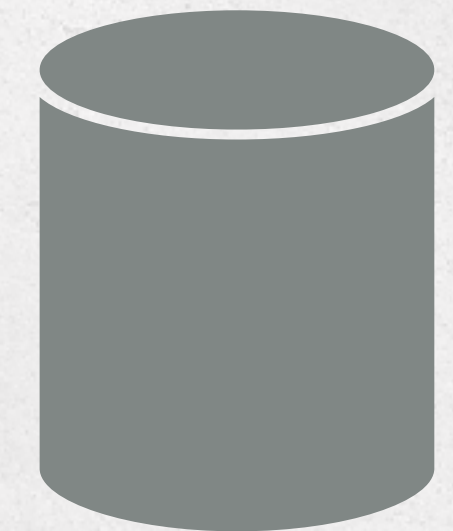
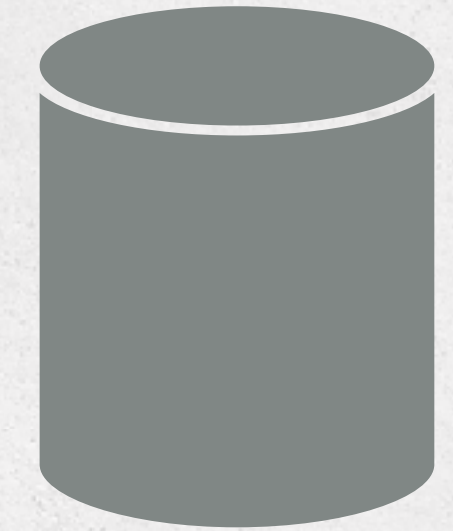
NON-BLOCKING ALGORITHMS

- **Lock-free** algorithms guarantee system-wide progress.
- **Wait-free** algorithms guarantee per-client progress. That is, no matter what steps other processes take, a correct client's operations are always completed in a finite number of steps.

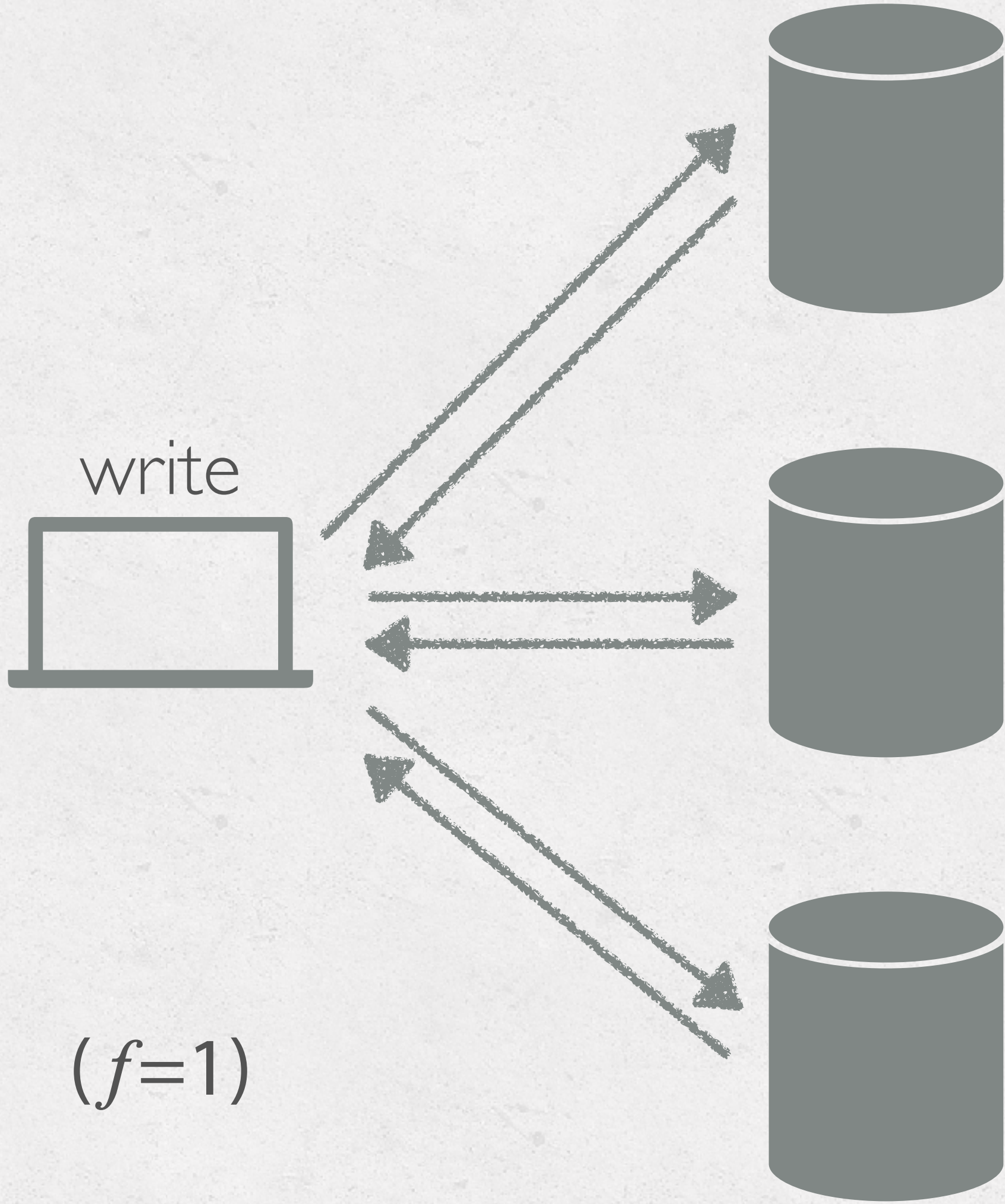
HOW MANY SERVERS DO WE NEED?



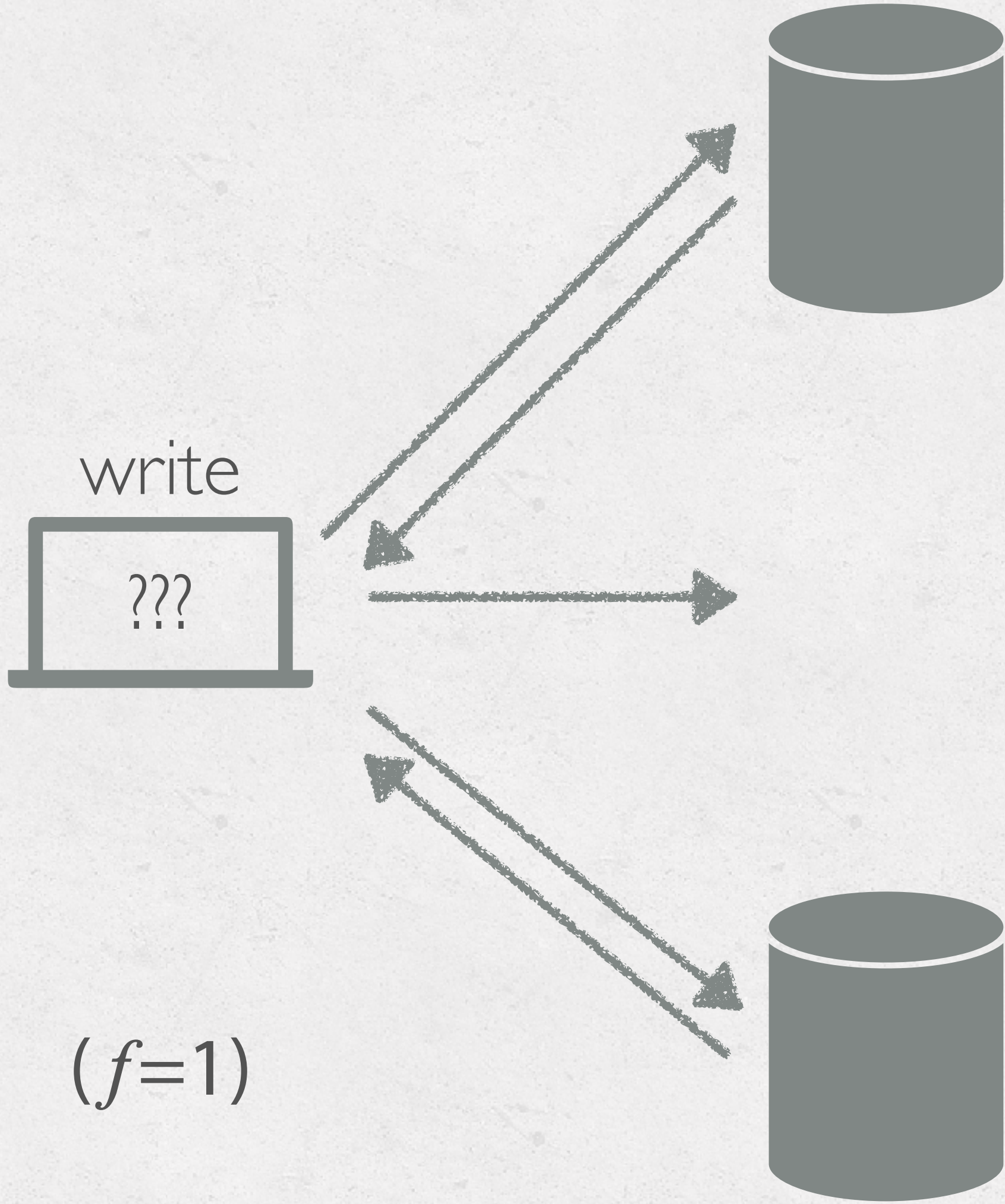
$(f=1)$



HOW MANY SERVERS DO WE NEED?

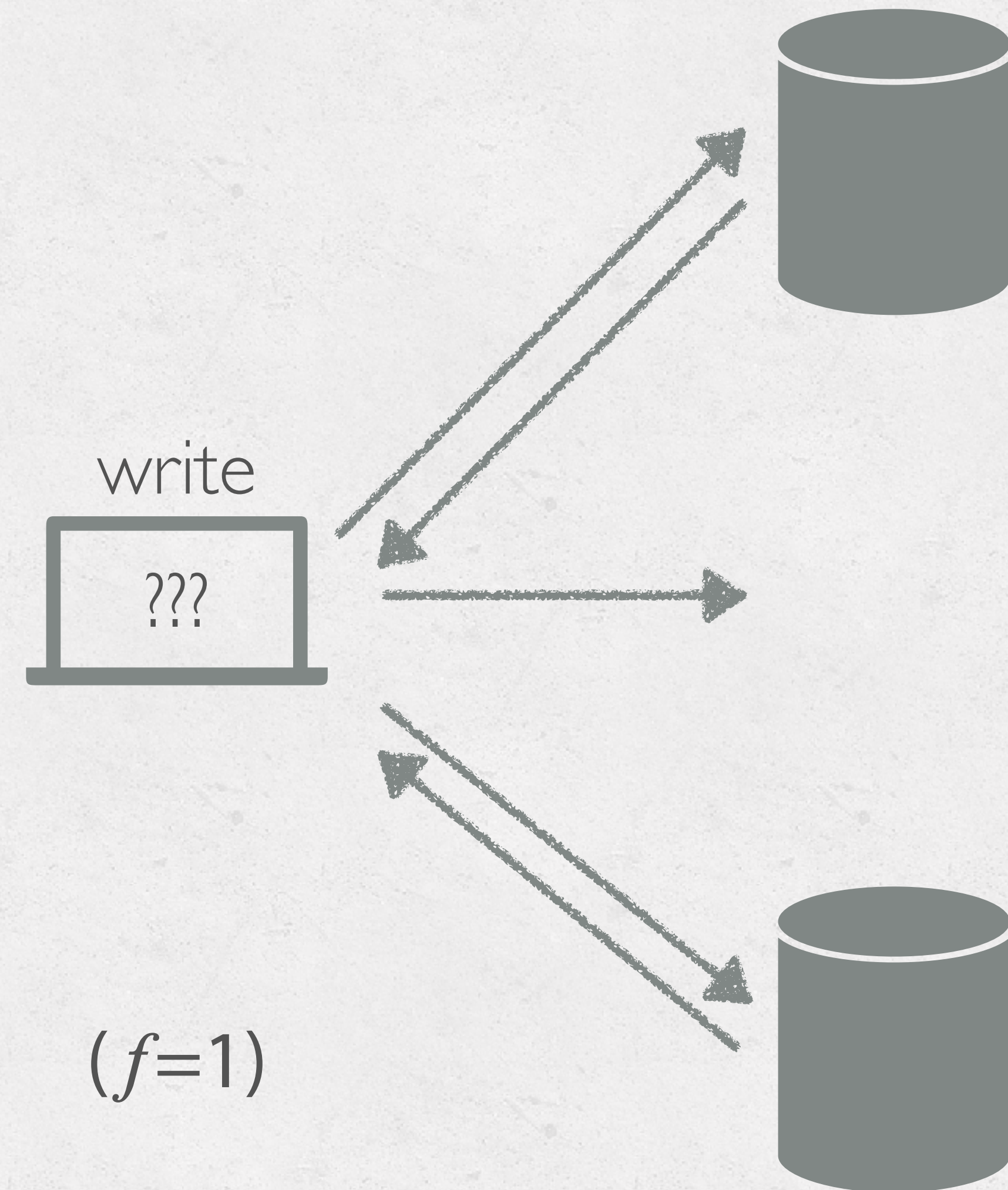


HOW MANY SERVERS DO WE NEED?



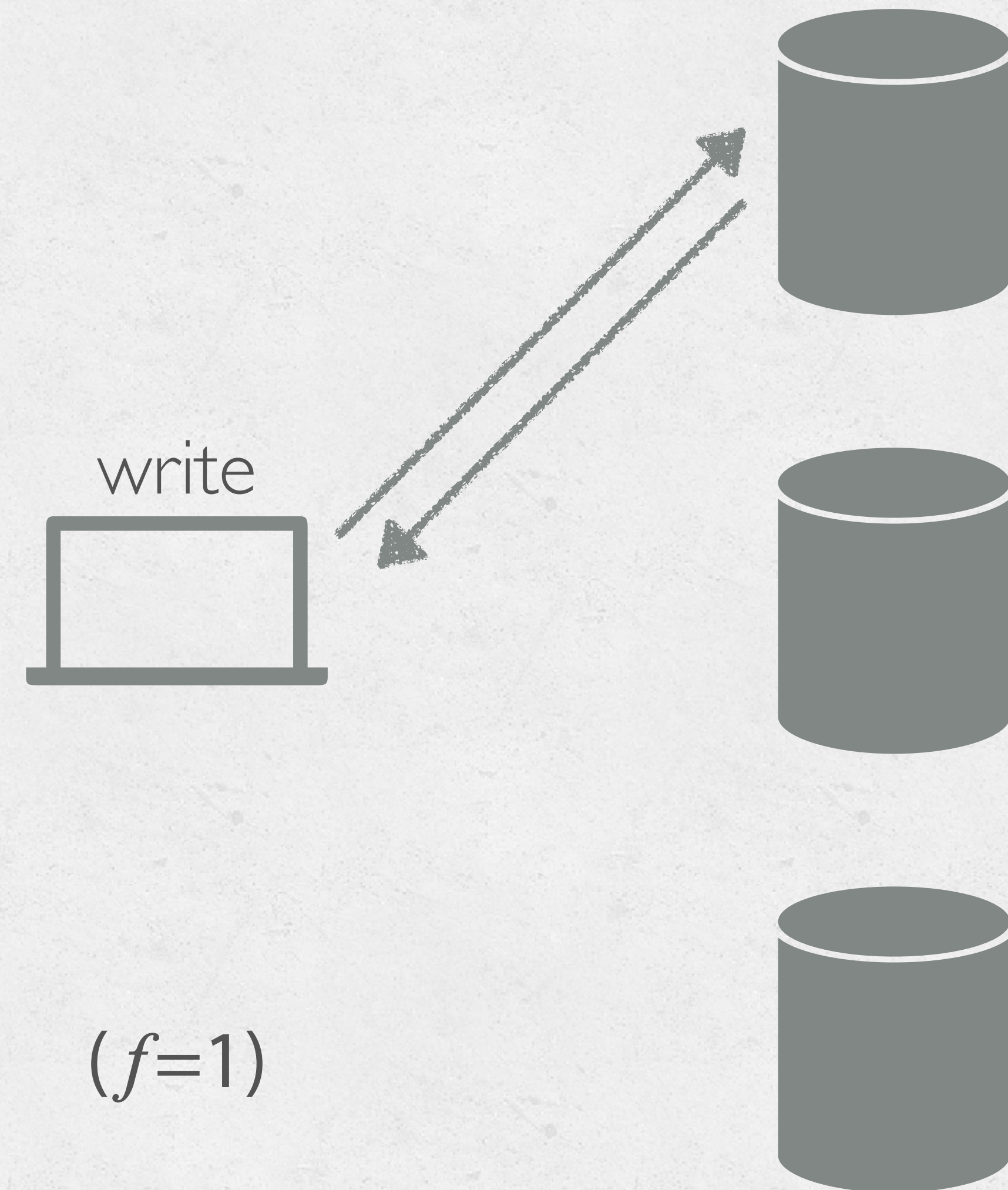
HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.



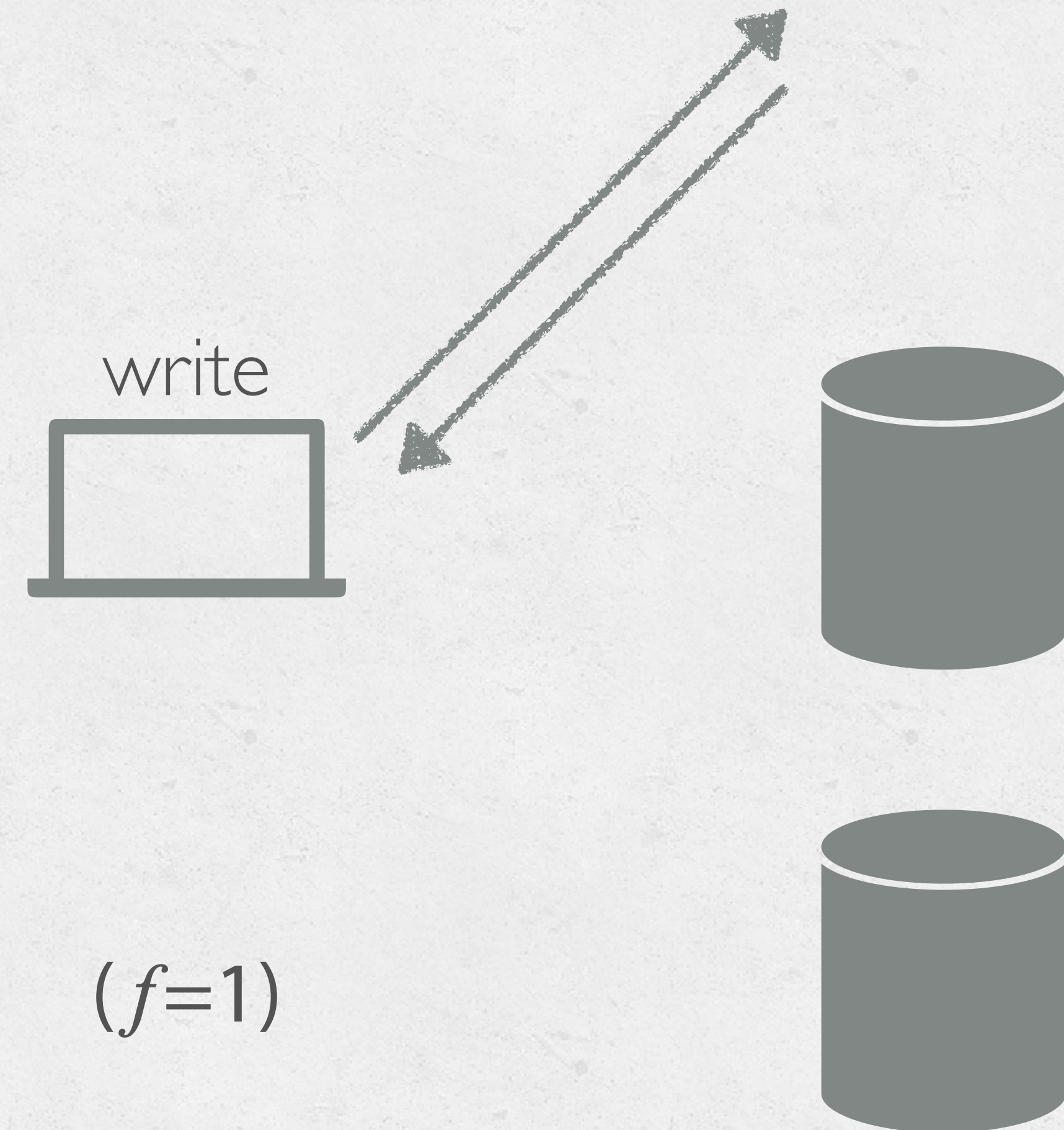
HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.



HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.



HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.

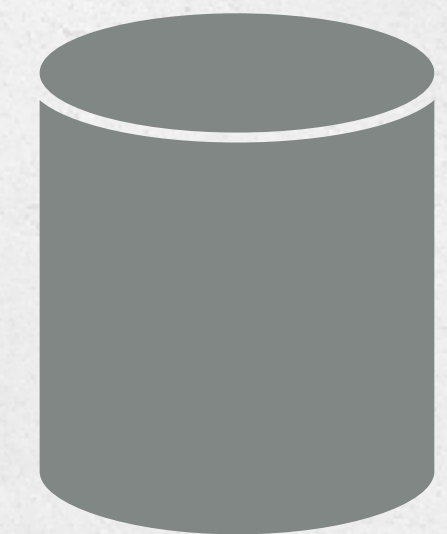
read



write



$(f=1)$

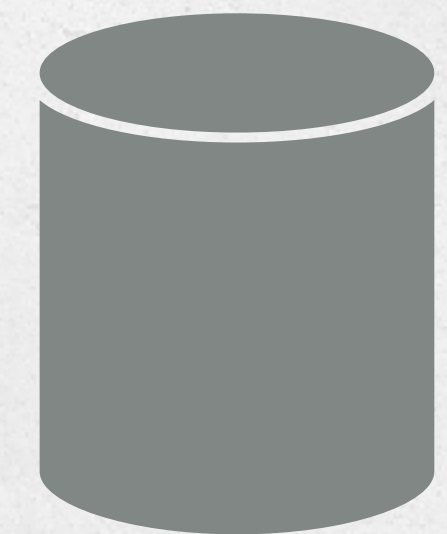


HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.
- We need to send writes to $> f$ replicas, otherwise they could get lost forever.



$(f=1)$

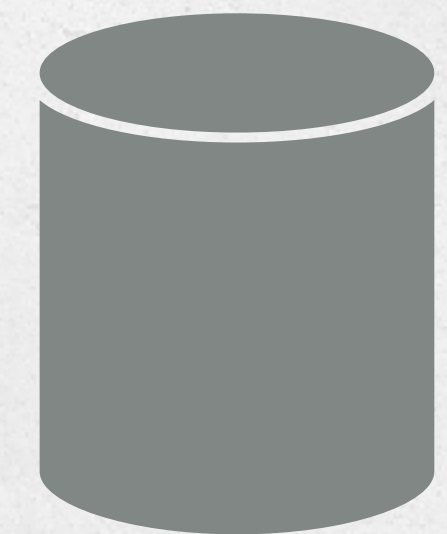


HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.
- We need to send writes to $> f$ replicas, otherwise they could get lost forever.
- So we need *at least* $2f + 1$ servers. And, in fact, we will use $2f + 1$.



$(f=1)$

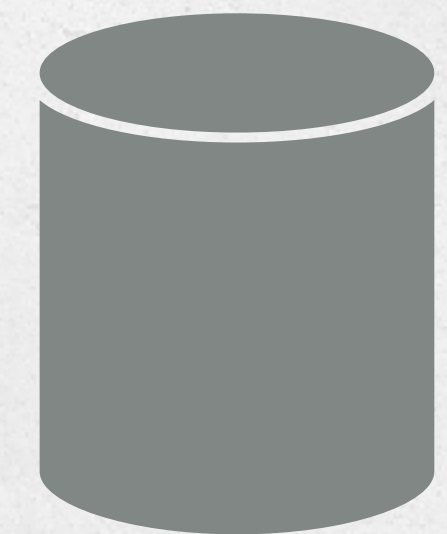


HOW MANY SERVERS DO WE NEED?

- If we want to make progress even when f servers crash, we can wait for at most $n - f$ responses.
- We need to send writes to $> f$ replicas, otherwise they could get lost forever.
- So we need *at least* $2f+1$ servers. And, in fact, we will use $2f+1$.
- Read quorum size plus write quorum size should be greater than n (i.e., they should overlap). We'll use simple majorities.



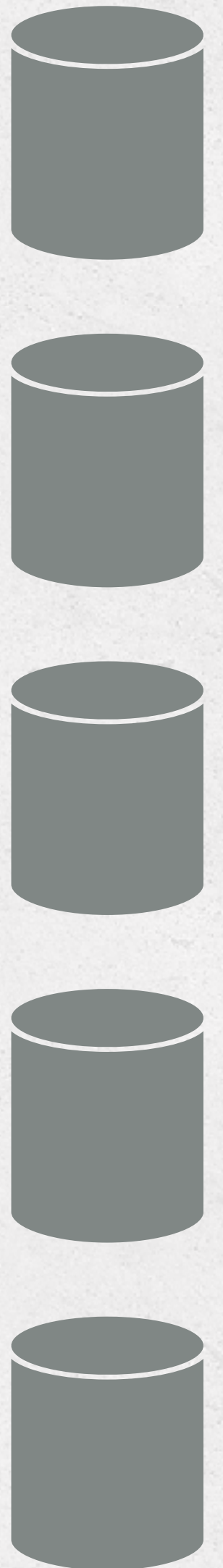
$(f=1)$



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

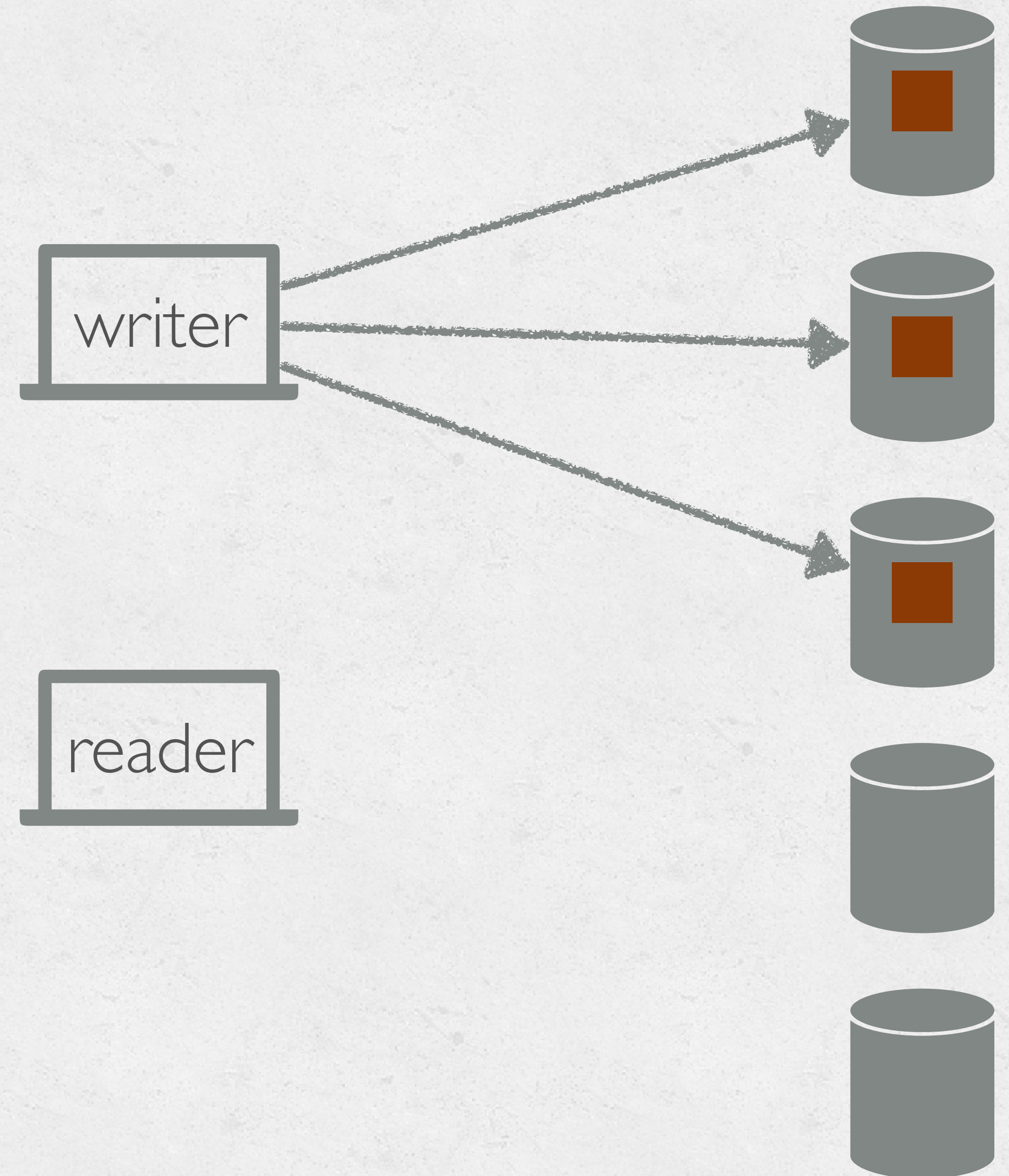
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

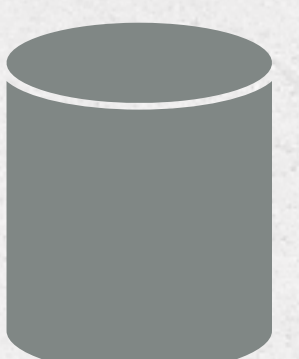
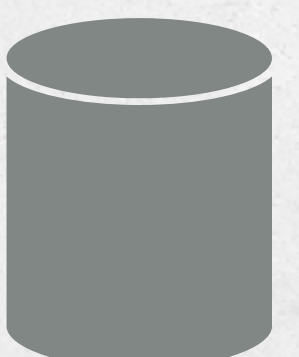
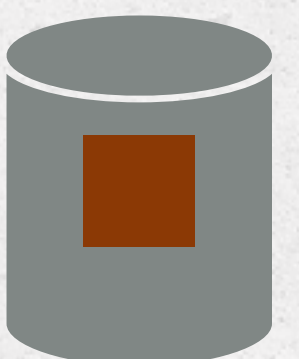
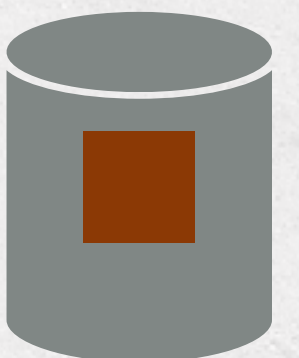
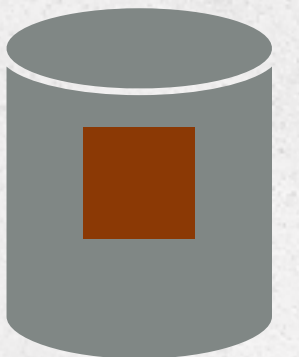
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

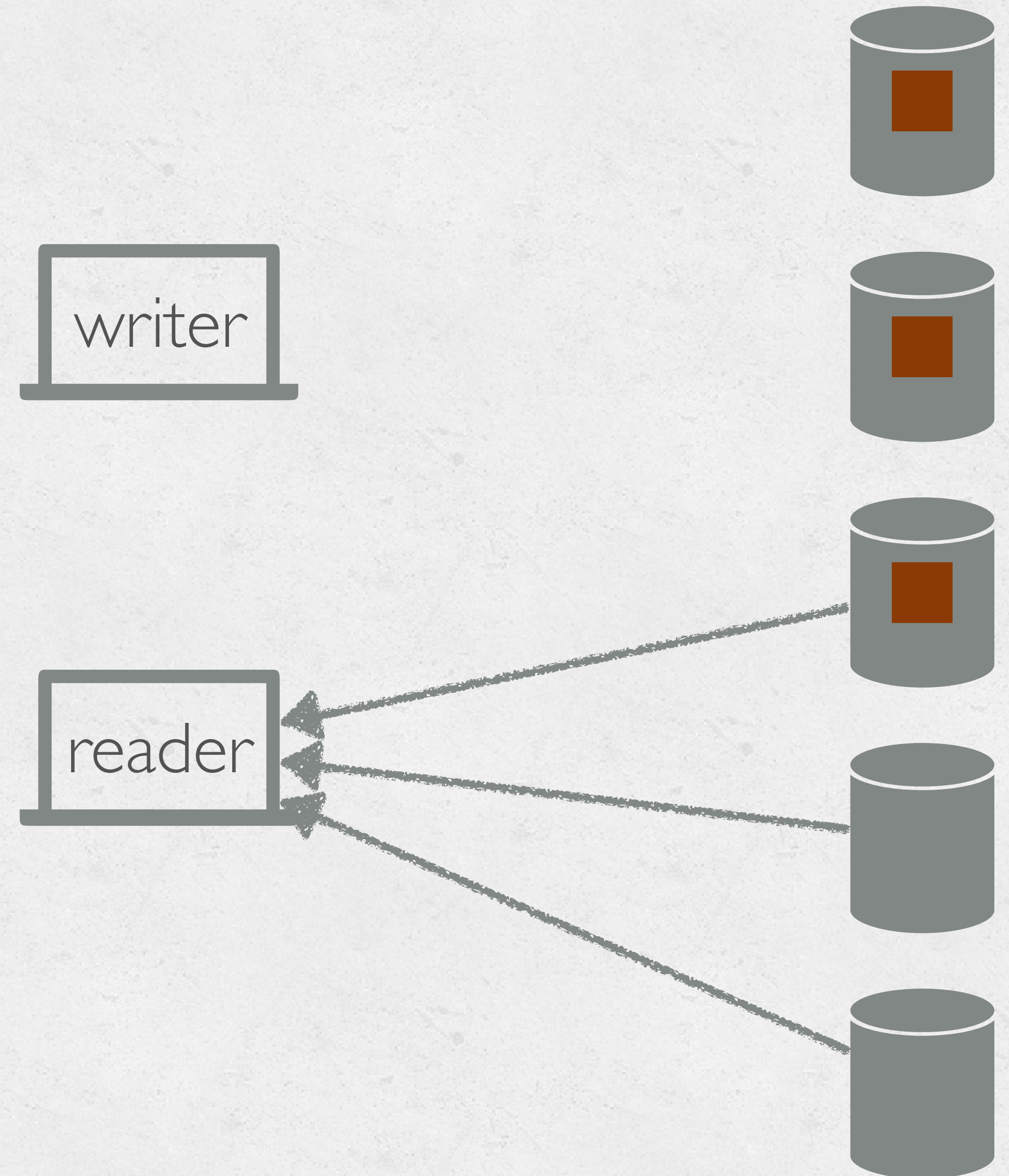
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

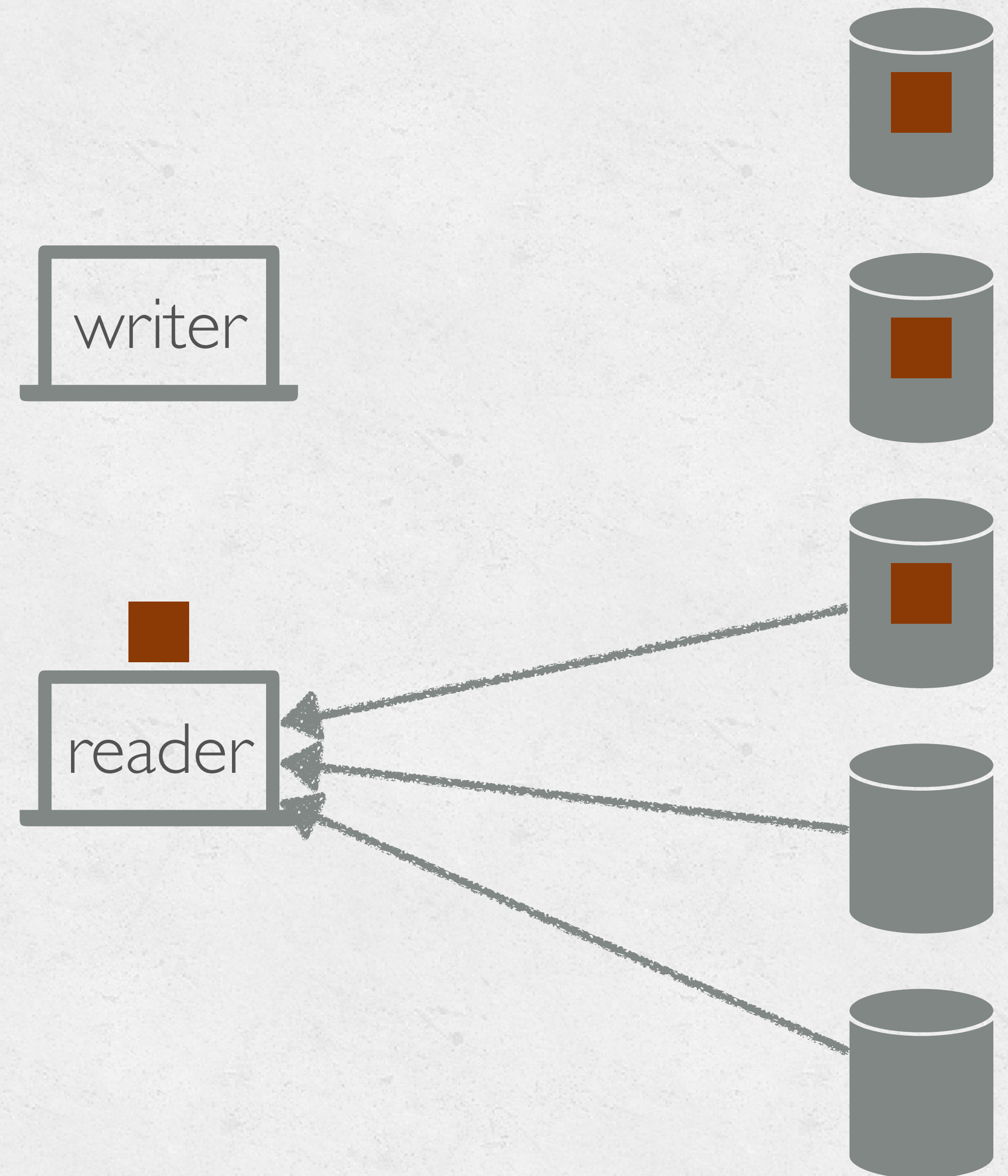
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

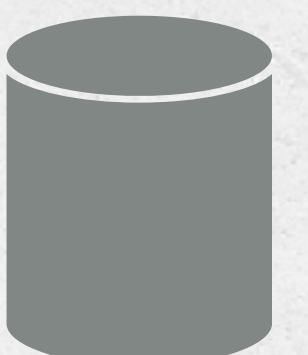
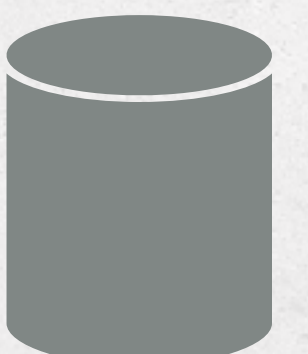
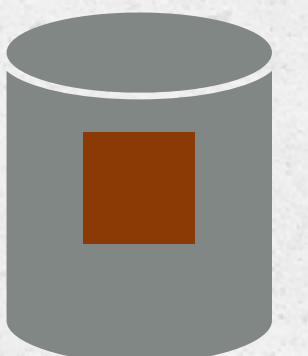
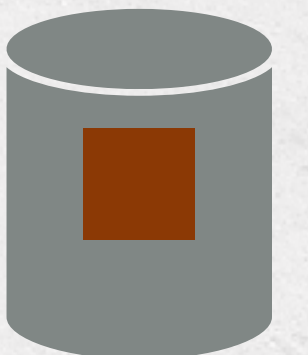
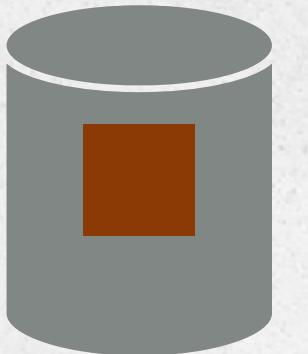
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

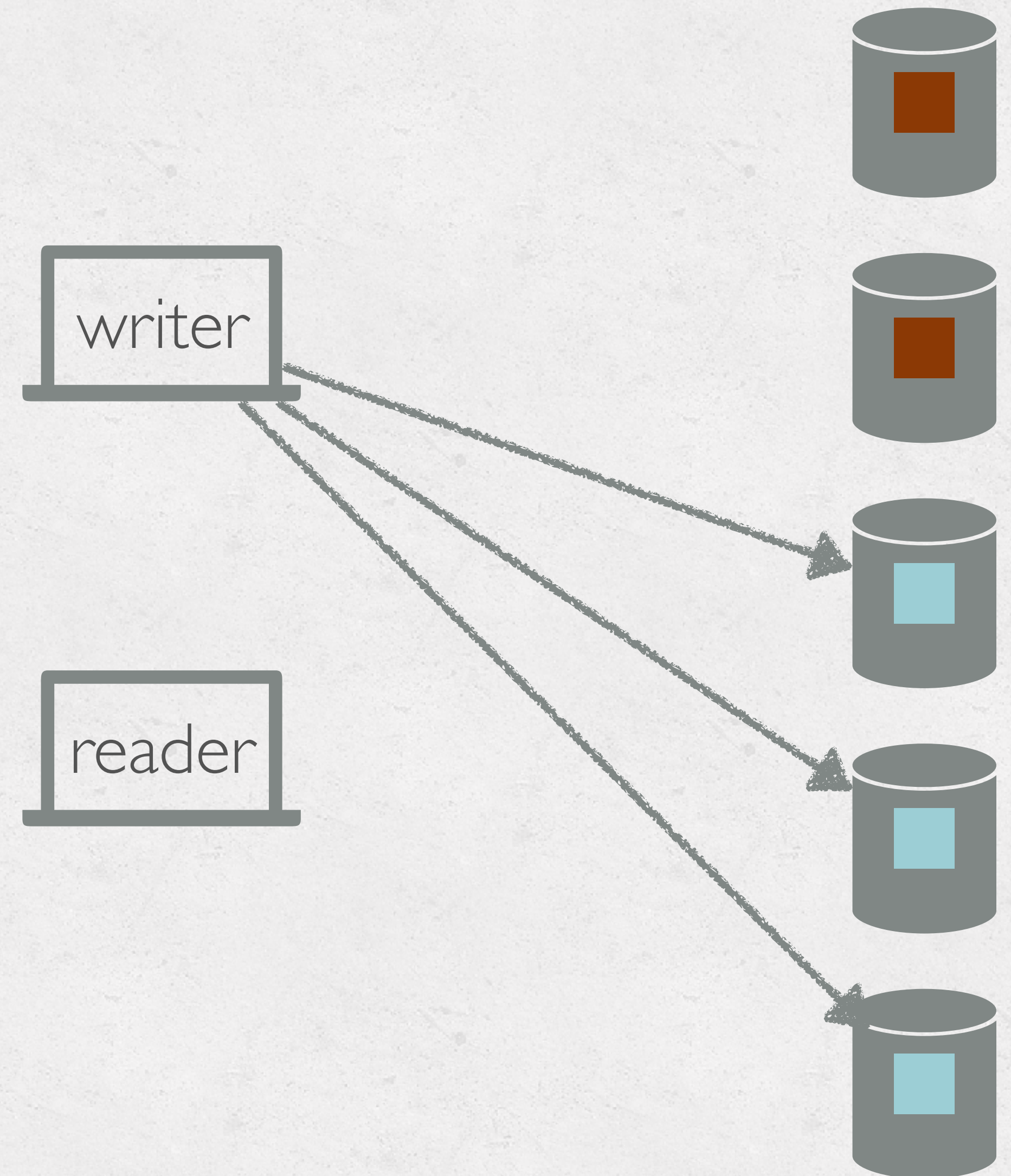
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

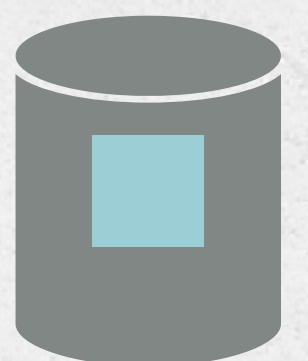
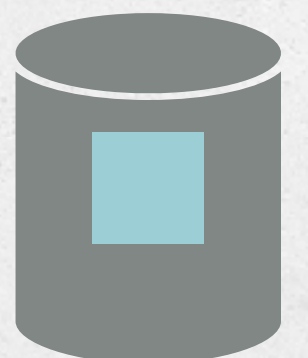
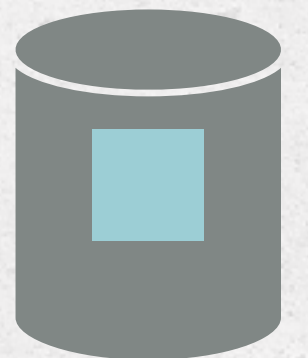
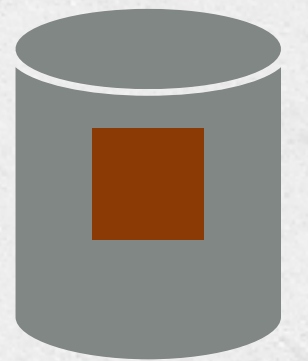
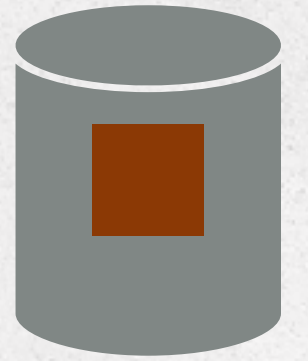
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

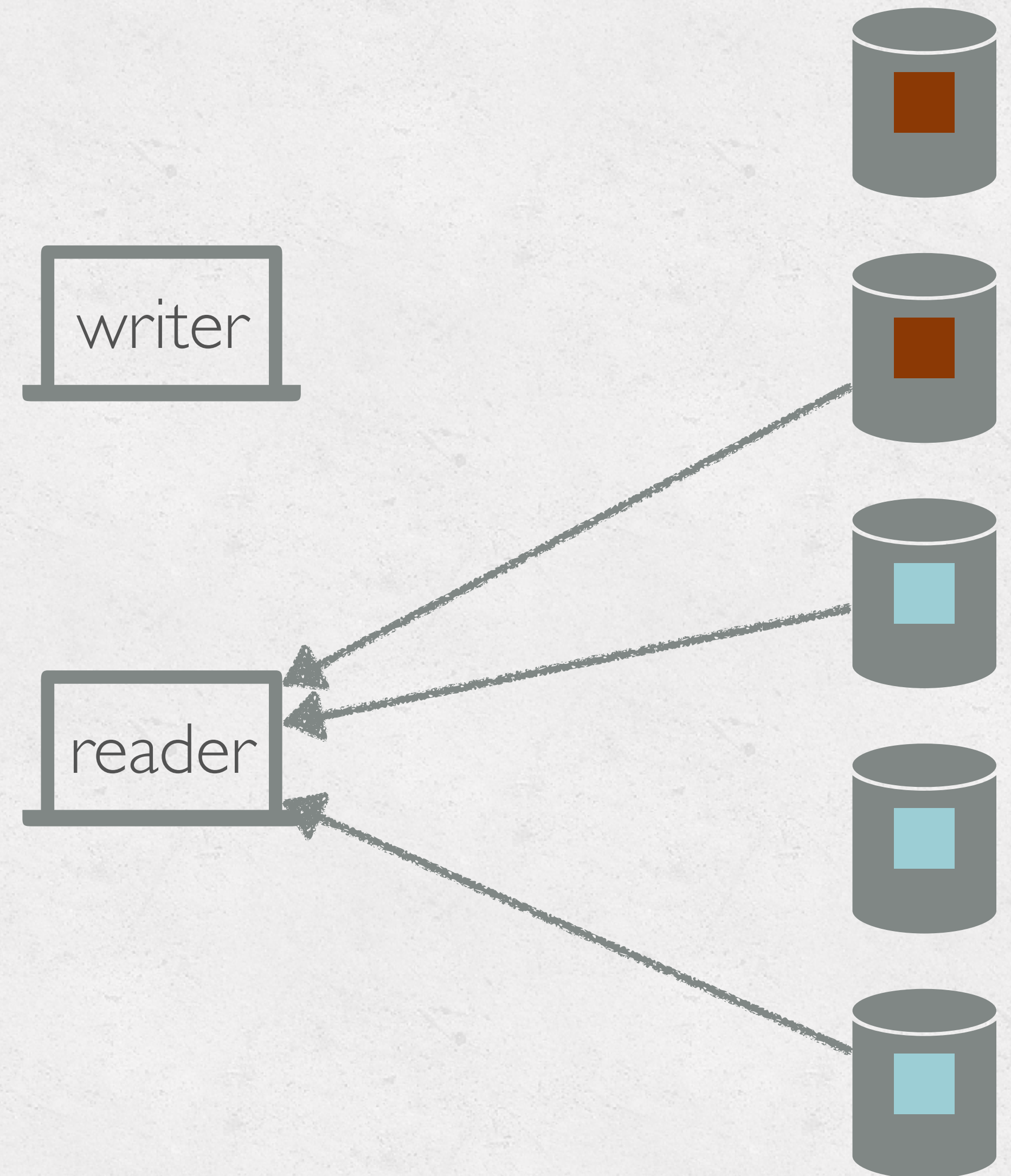
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

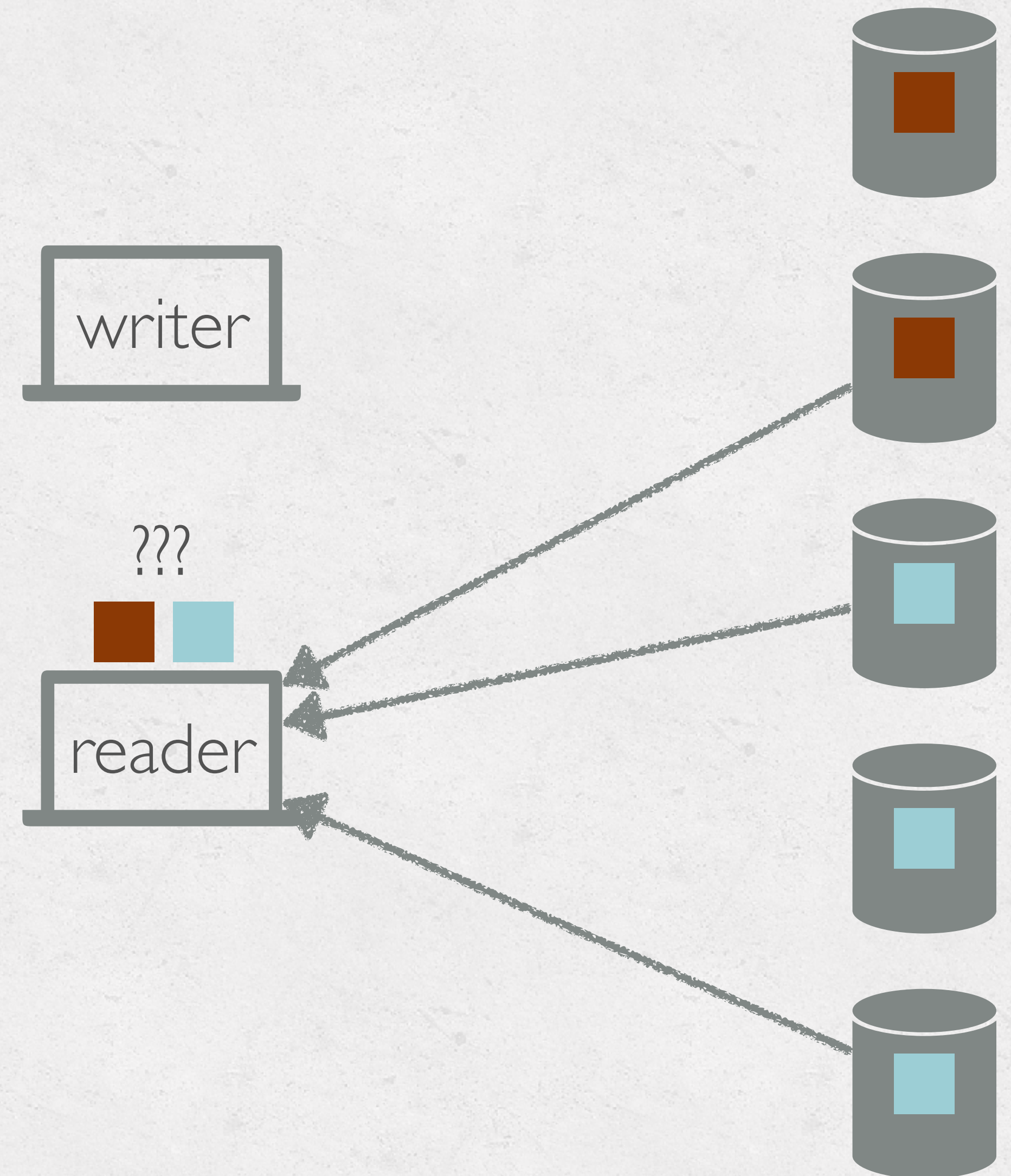
Does this work?



FIRST STEP: SINGLE READER, SINGLE WRITER (SWSR)

- Writer sends value to a majority.
- Reader reads value from a majority.
- Since majorities intersect, reader reads writer's value.

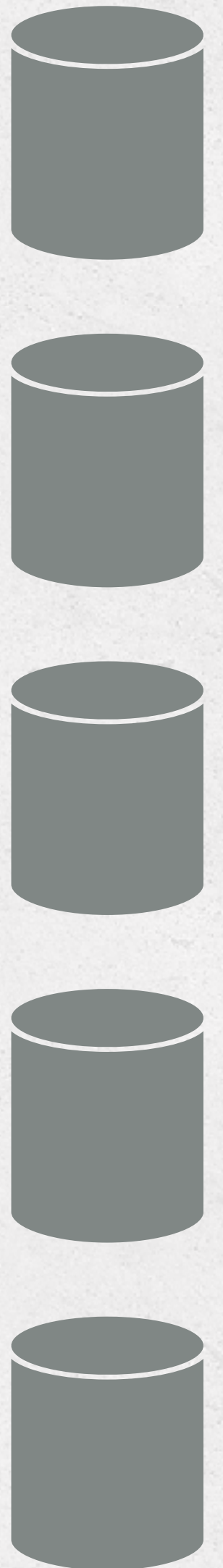
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

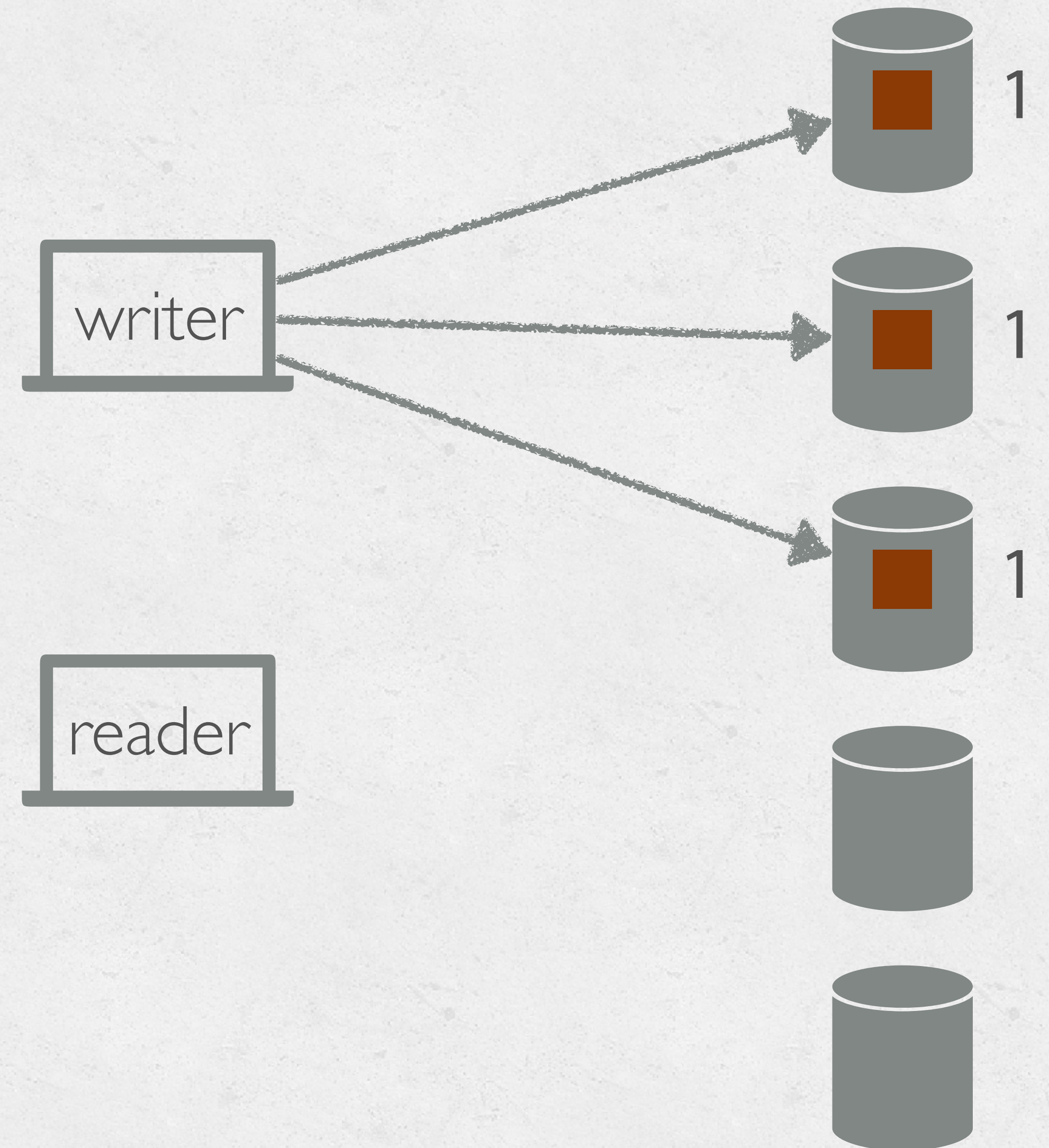
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

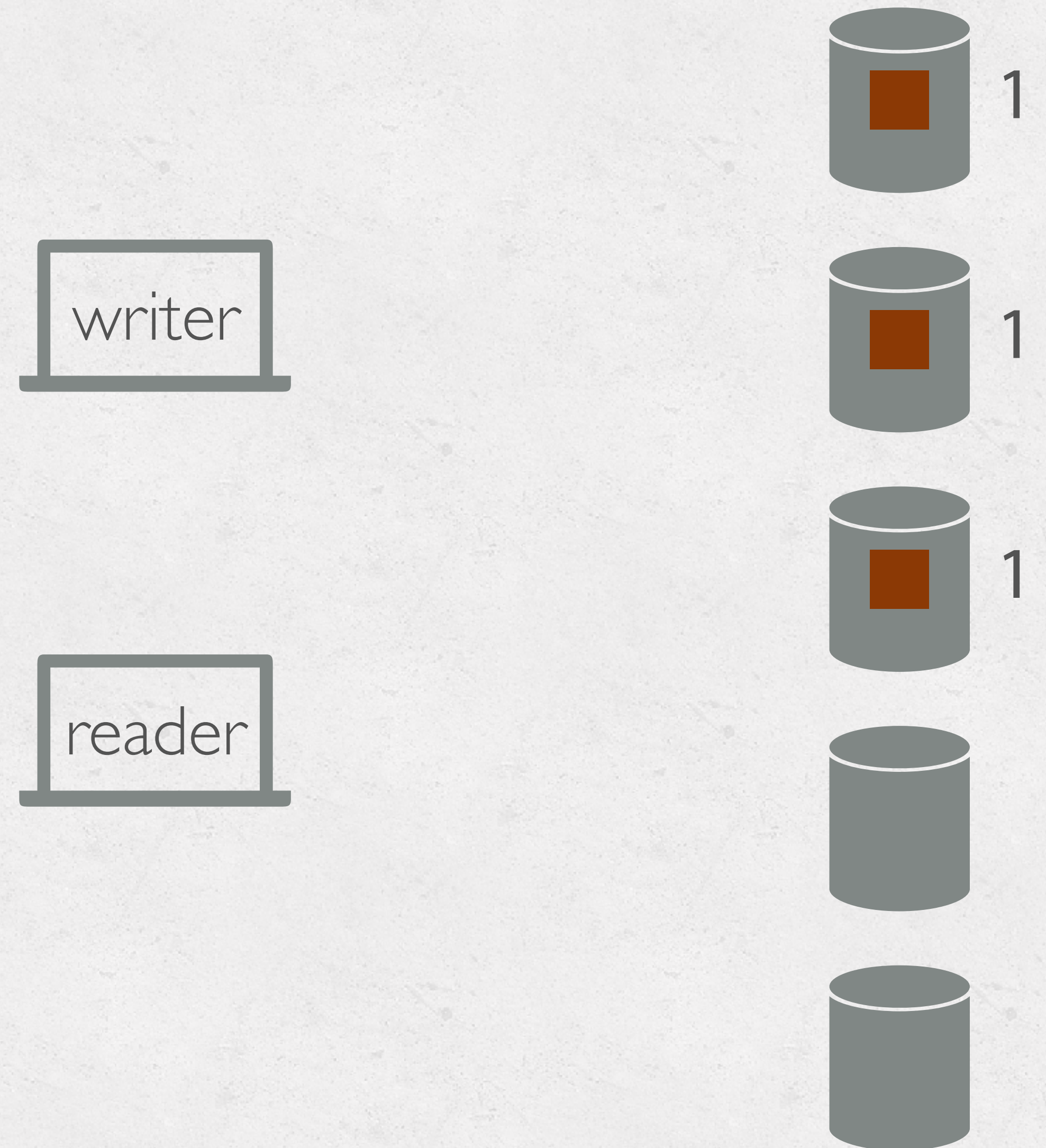
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

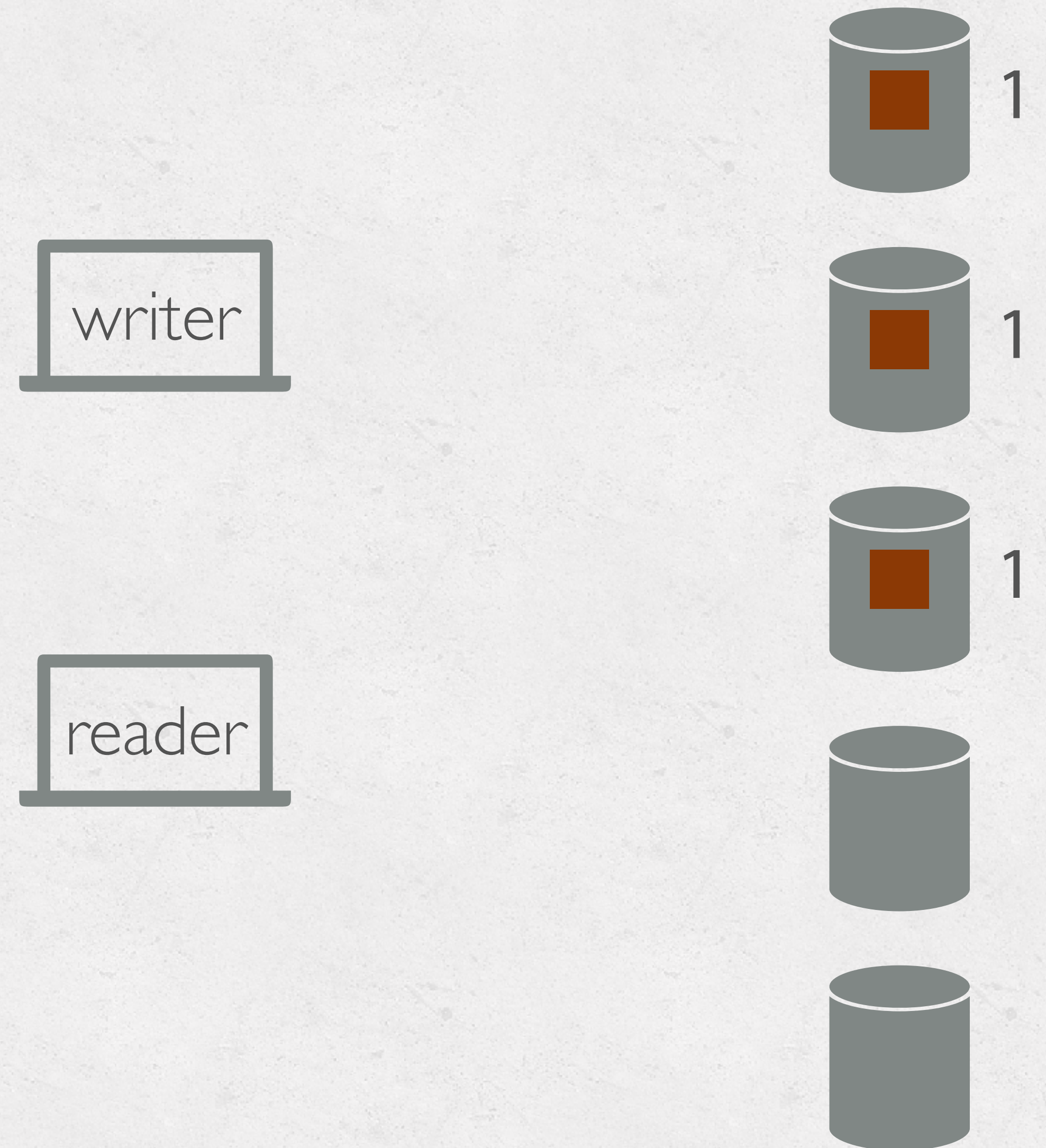
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

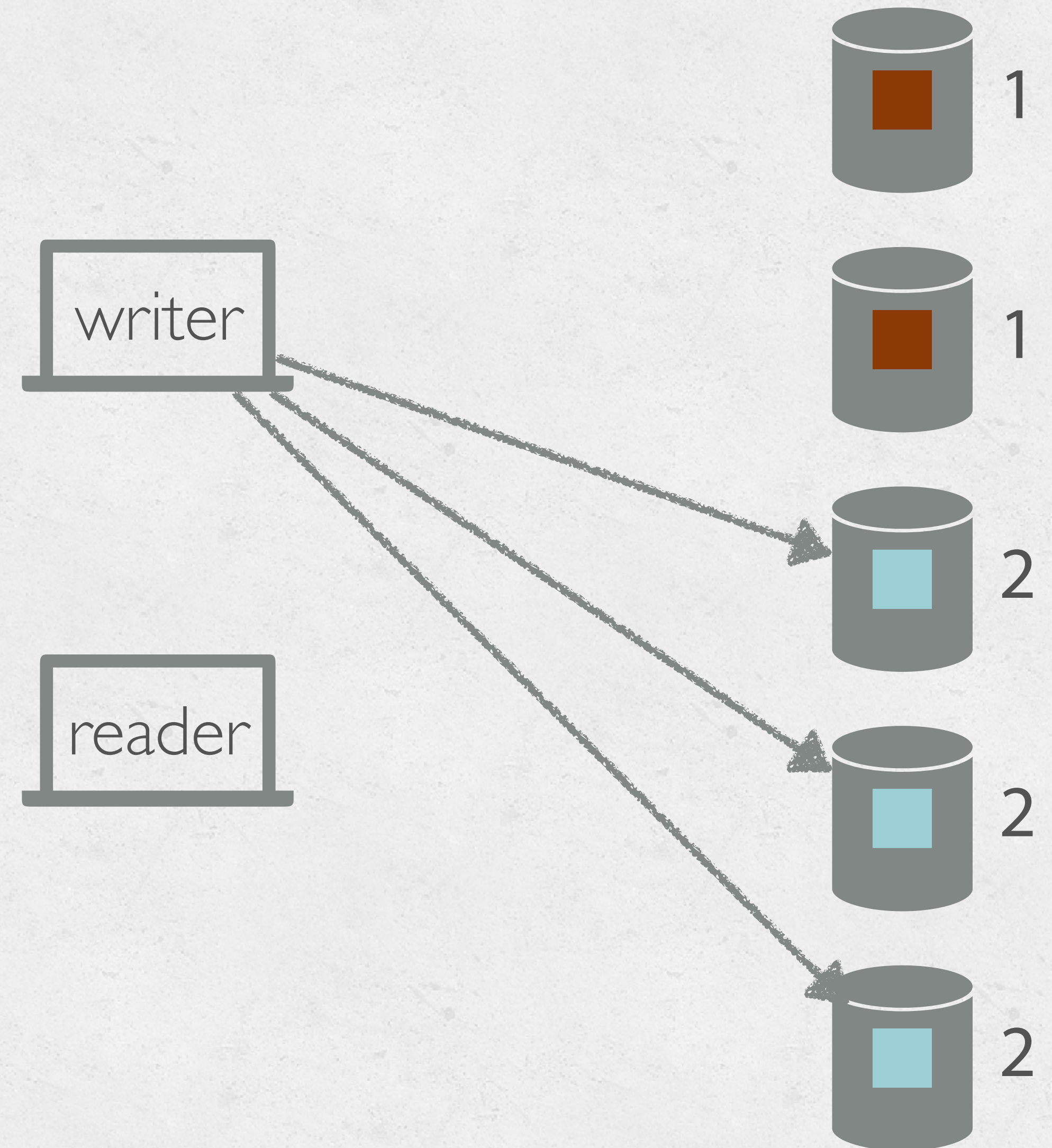
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

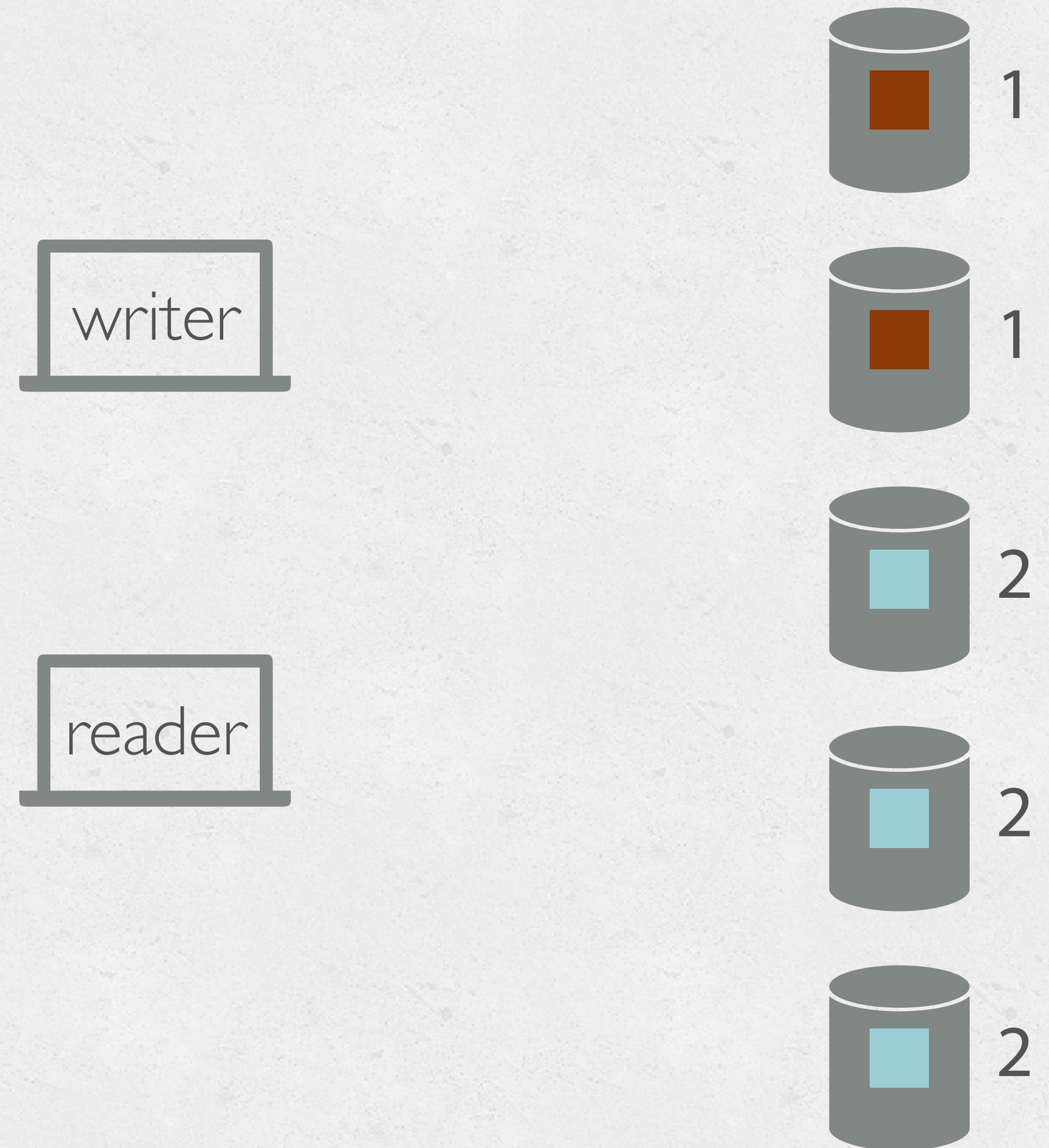
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

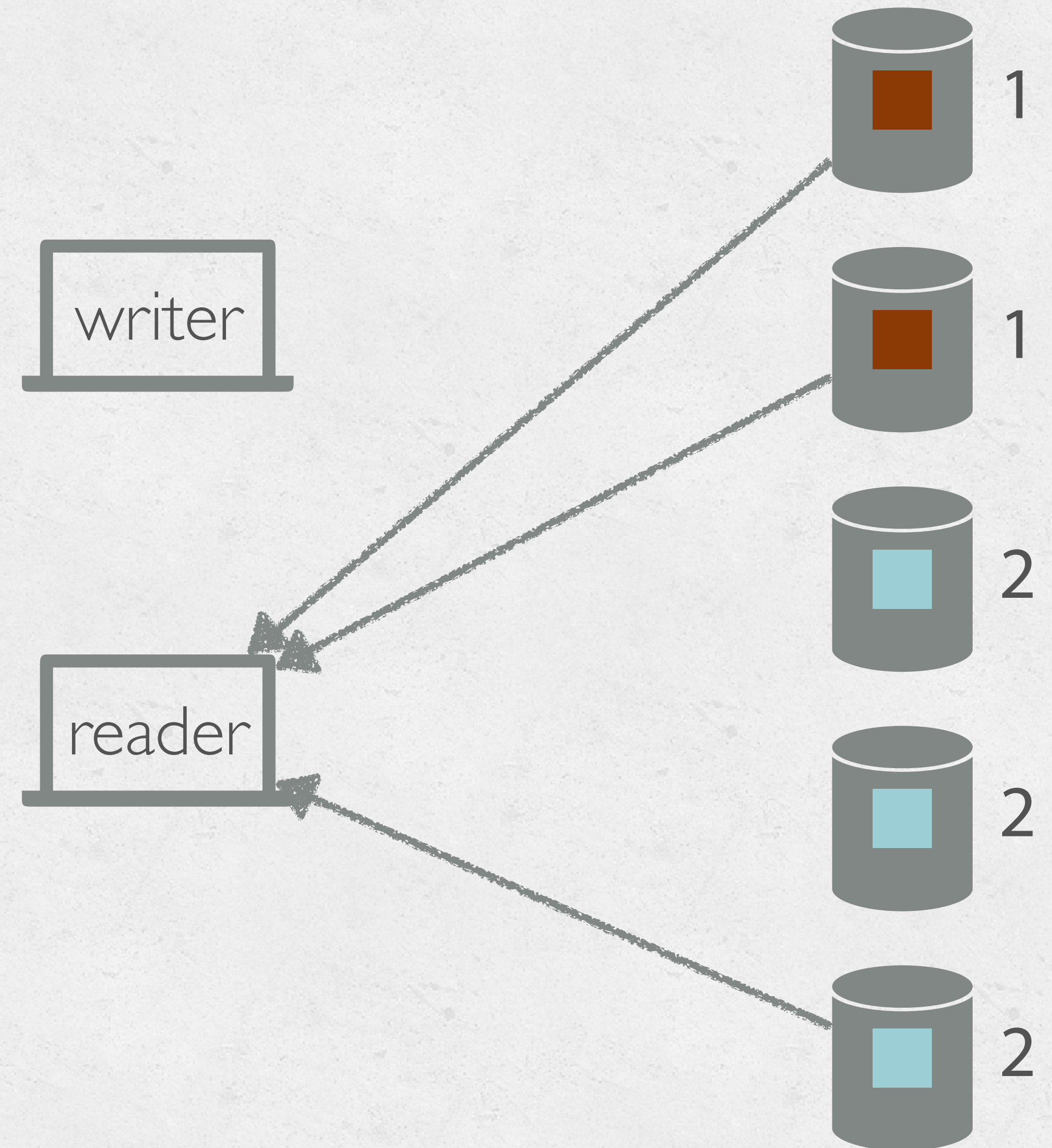
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

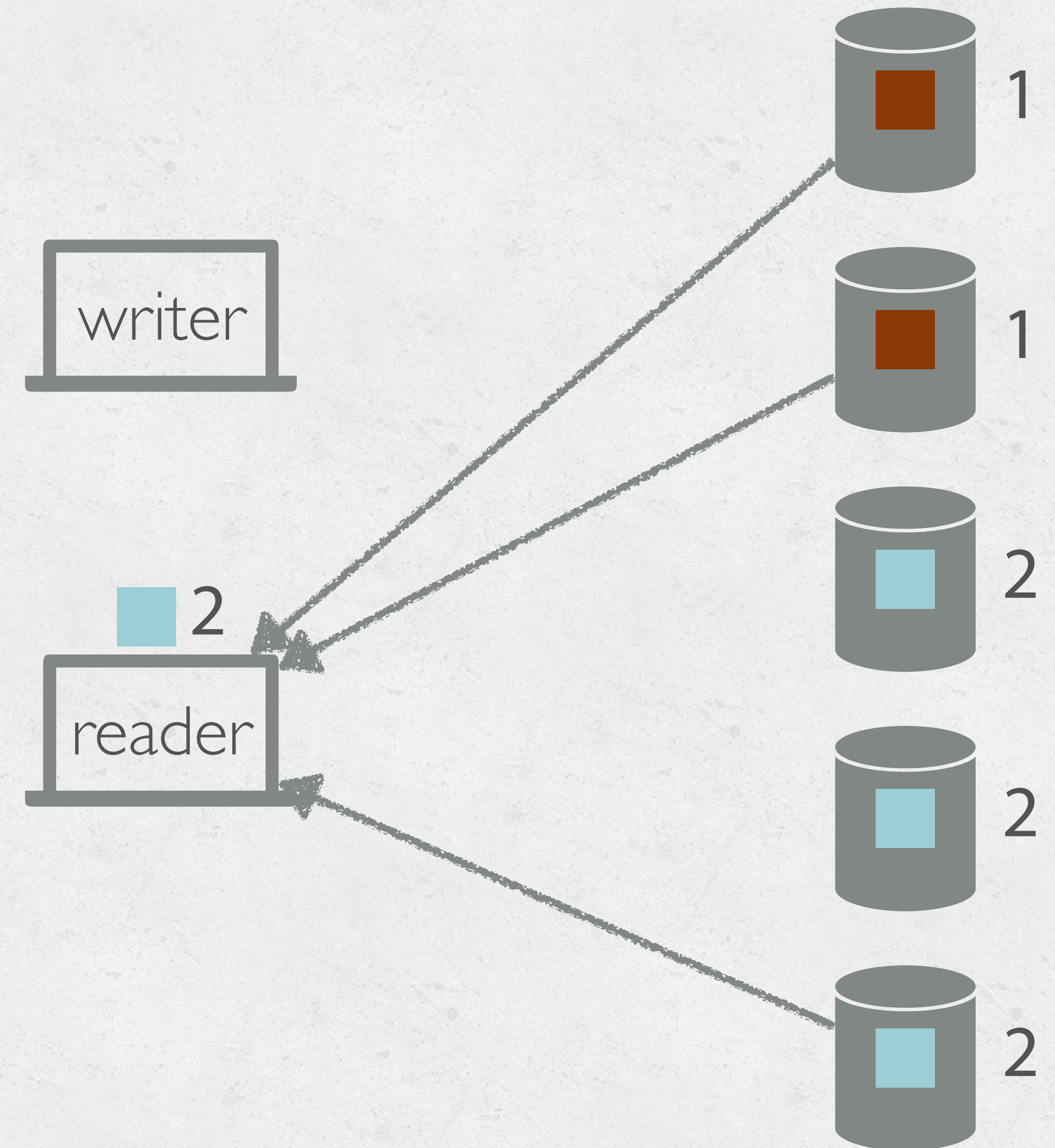
Does this work?



SWSR II

- Writer sends a **timestamped** value to a majority.
- Reader reads value from a majority, takes the one with the **highest timestamp**.
- Since majorities intersect, reader reads writer's value.

Does this work?



SWSR III

Servers' algorithm:

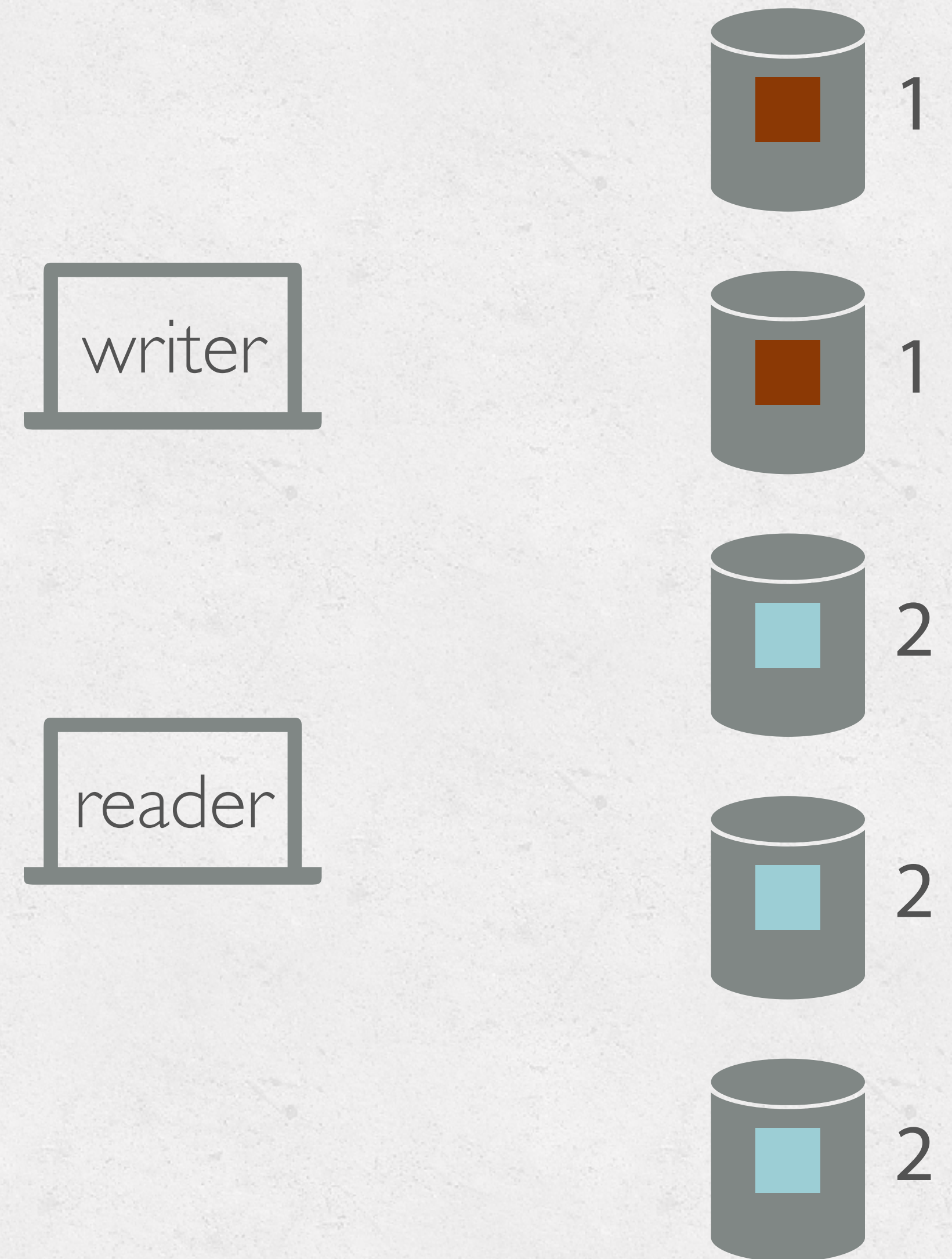
- Upon receiving a write, update local timestamp and value if write's timestamp is greater; send ack.
- Respond to reads with local timestamp and value.

Writer's algorithm:

- When writing, increment local timestamp, send timestamp and value to all.
- Wait for acks from a majority.

Reader's algorithm:

- Read from a majority, take value with highest timestamp.
- Maintain local value, return local value if servers' timestamps smaller.



SWSR III

Servers' algorithm:

- Upon receiving a request, write the value if
- Respond to requests with responses (i.e., ignore responses from old requests)

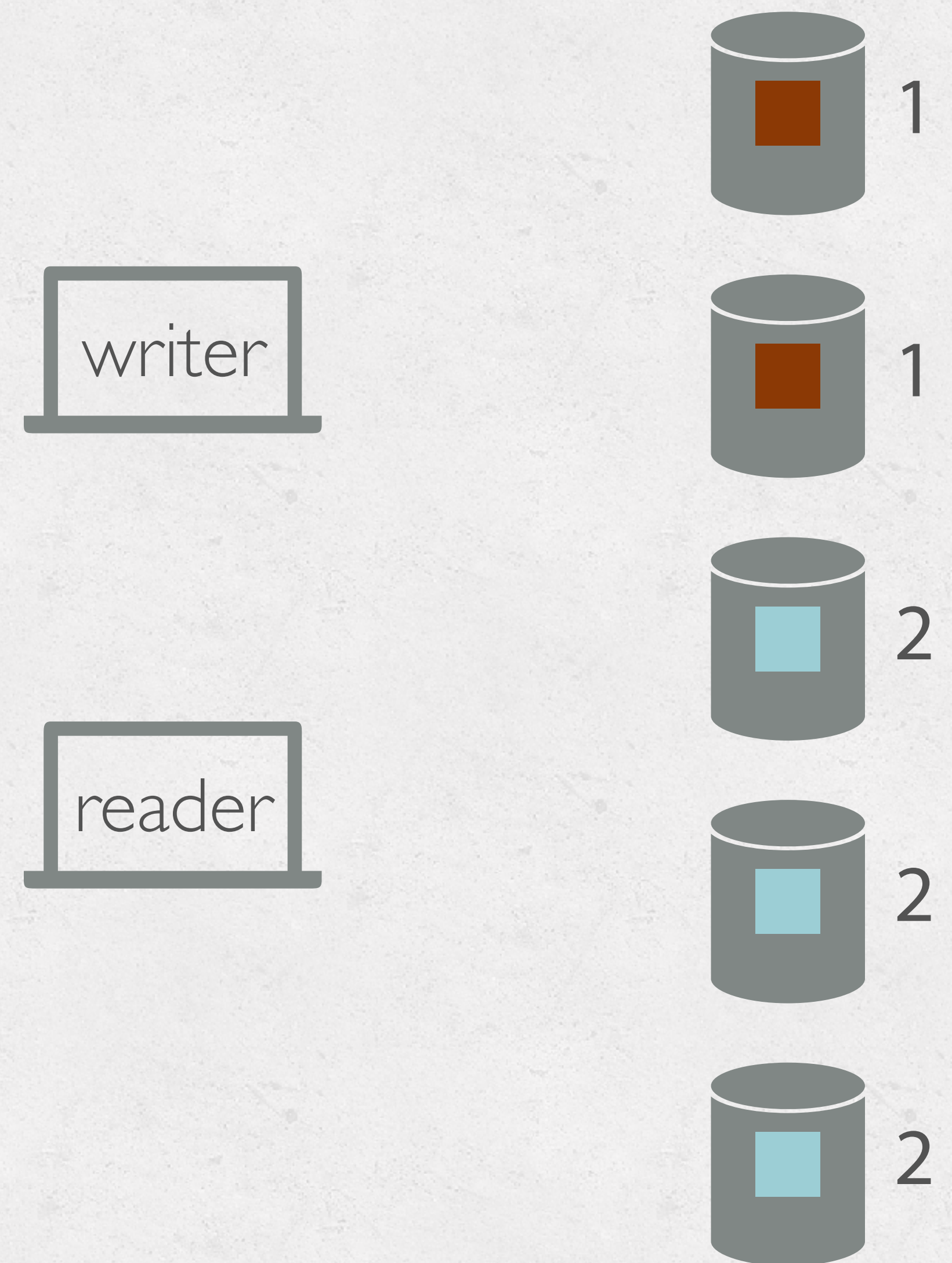
Writer's

- When writing, increment local timestamp, send timestamp and value to all.
- Wait for acks from a majority.

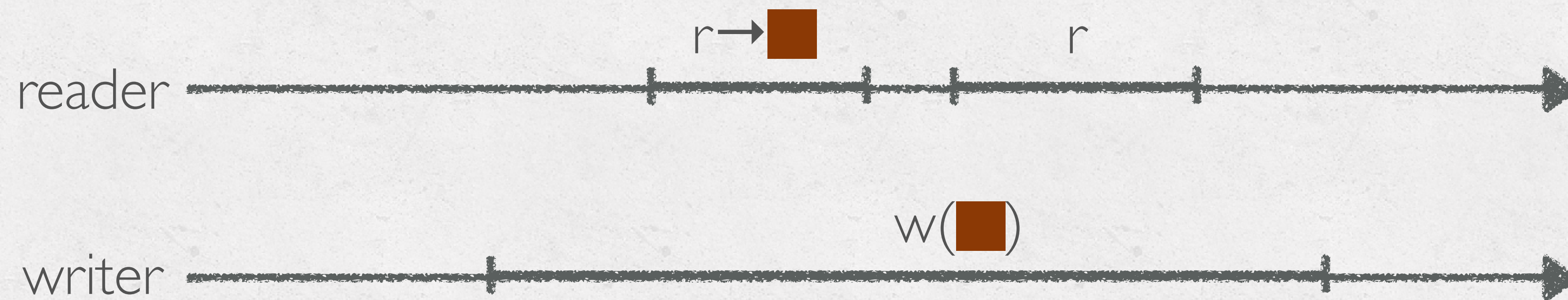
Reader's algorithm:

- Read from a majority, take value with highest timestamp.
- Maintain local value, return local value if servers' timestamps smaller.

Assume clients can associate requests with responses (i.e., ignore responses from old requests)

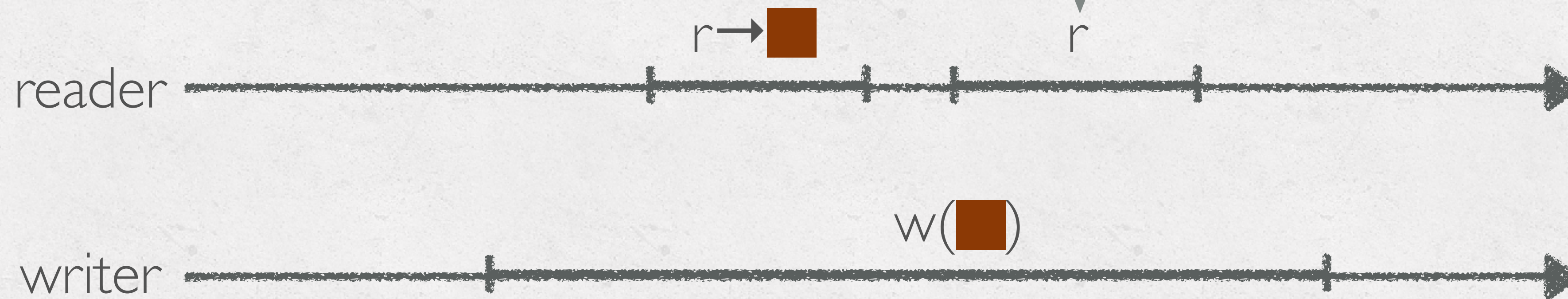


SRSW: WHAT ABOUT MULTIPLE READS?



SRSW: WHAT ABOUT MULTIPLE READS?

Guaranteed to return the red value, stored in the reader's cache.



If there's only one writer and one reader, why do we need the servers at all? Couldn't the writer just send its value to the reader directly?

MULTIPLE READERS, SINGLE WRITER (MRSW)

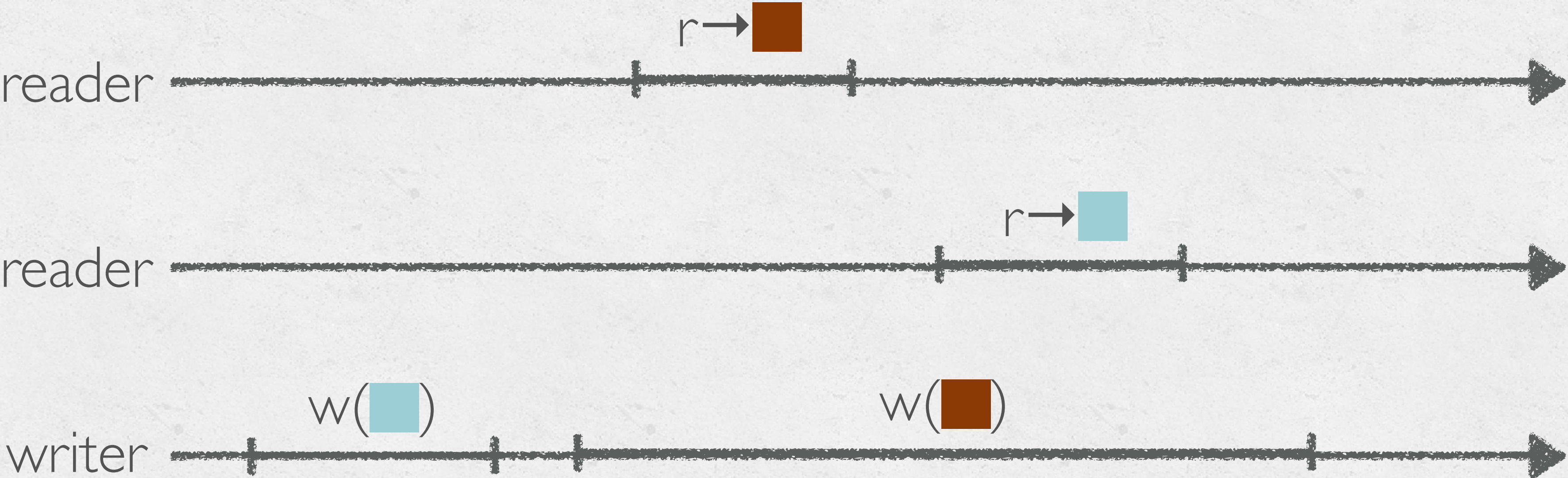
Does this previous solution just work?

MULTIPLE READERS, SINGLE WRITER (MRSW)

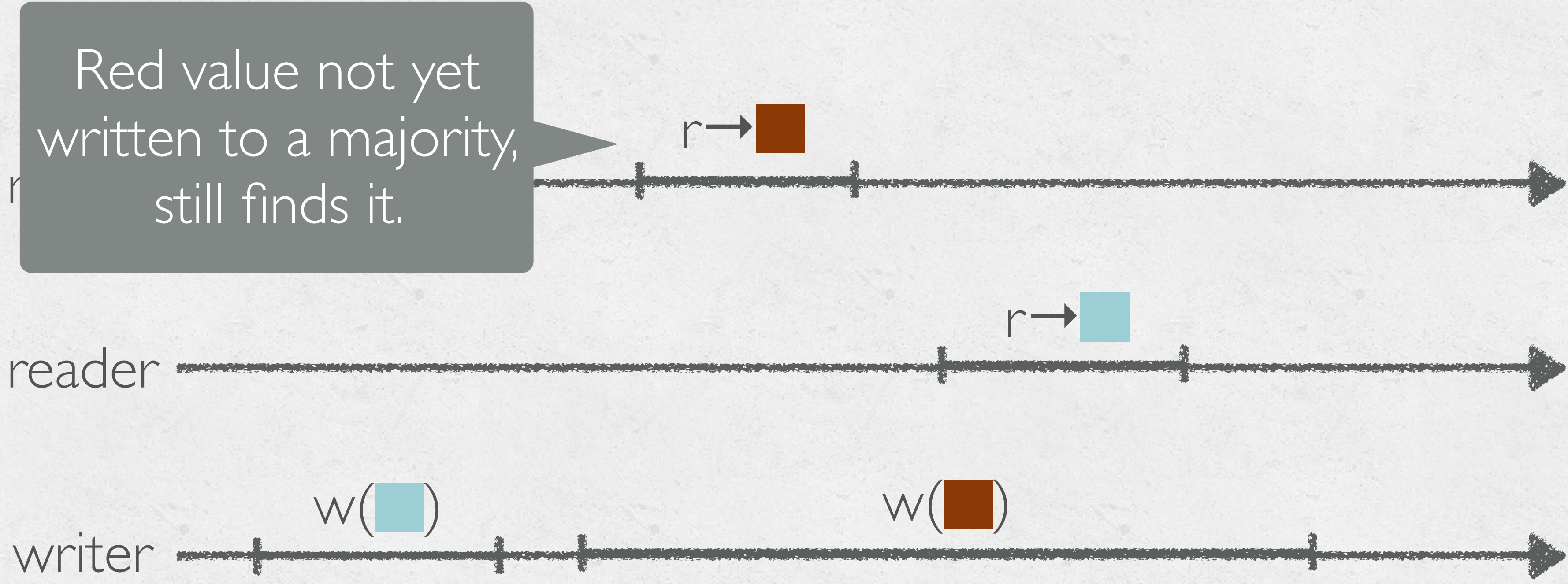
Does this previous solution just work?

What happens if there are multiple reads by **different processes** overlapping the same write?

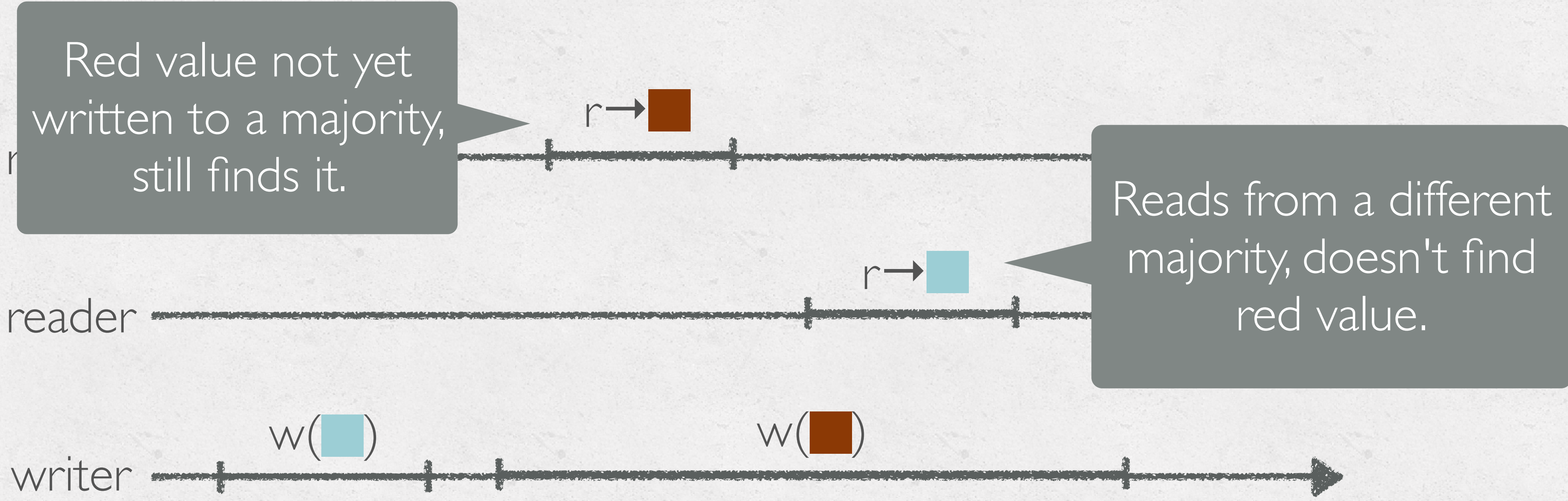
MRSW: REVENGE OF THE READS



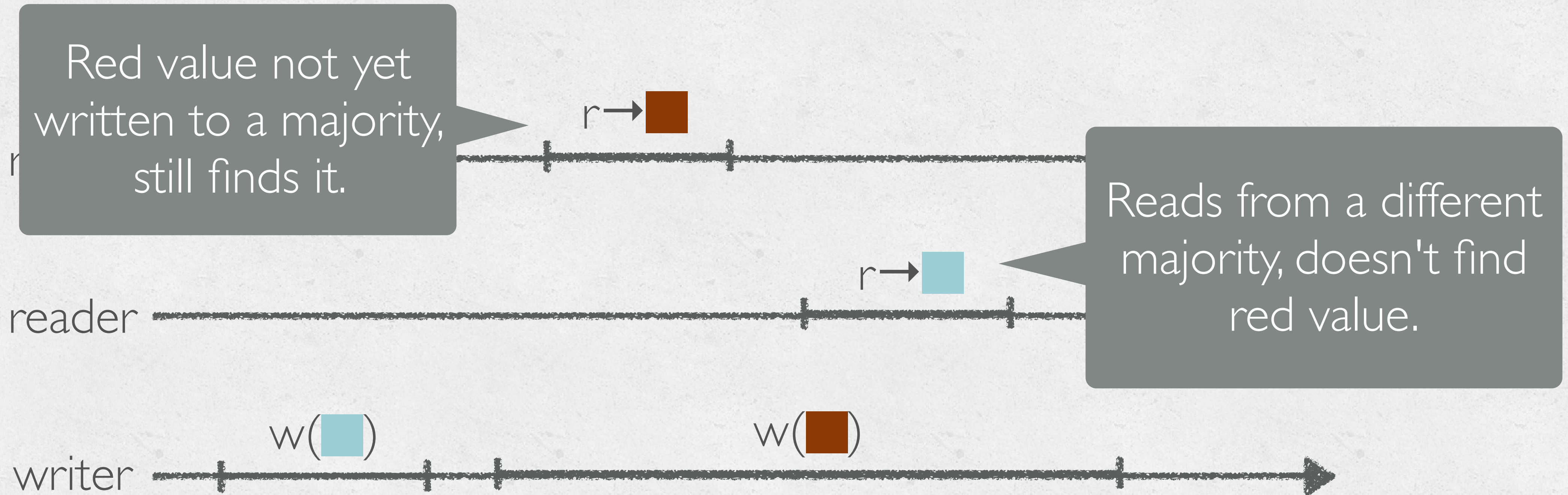
MRSW: REVENGE OF THE READS



MRSW: REVENGE OF THE READS



MRSW: REVENGE OF THE READS

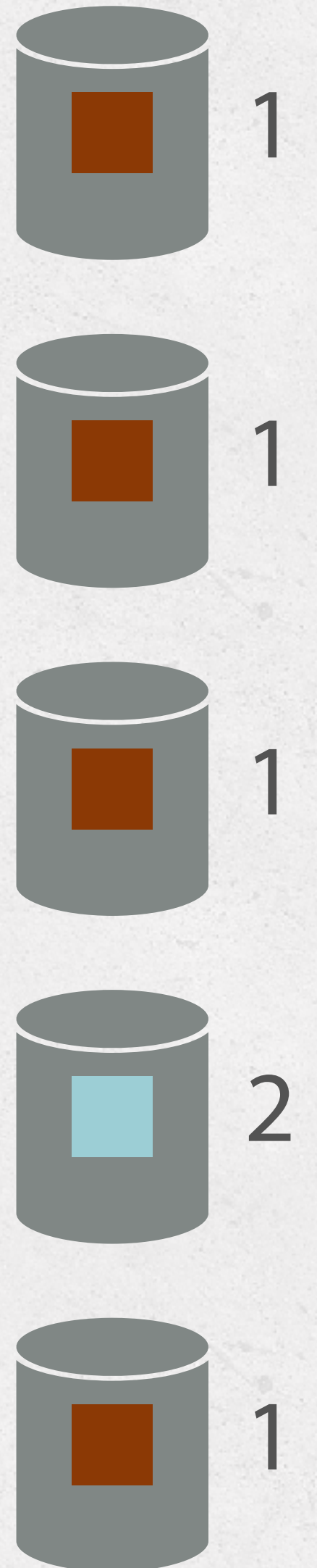


Not linearizable!

MRSW II

Suppose a write is ongoing (or the writer died).

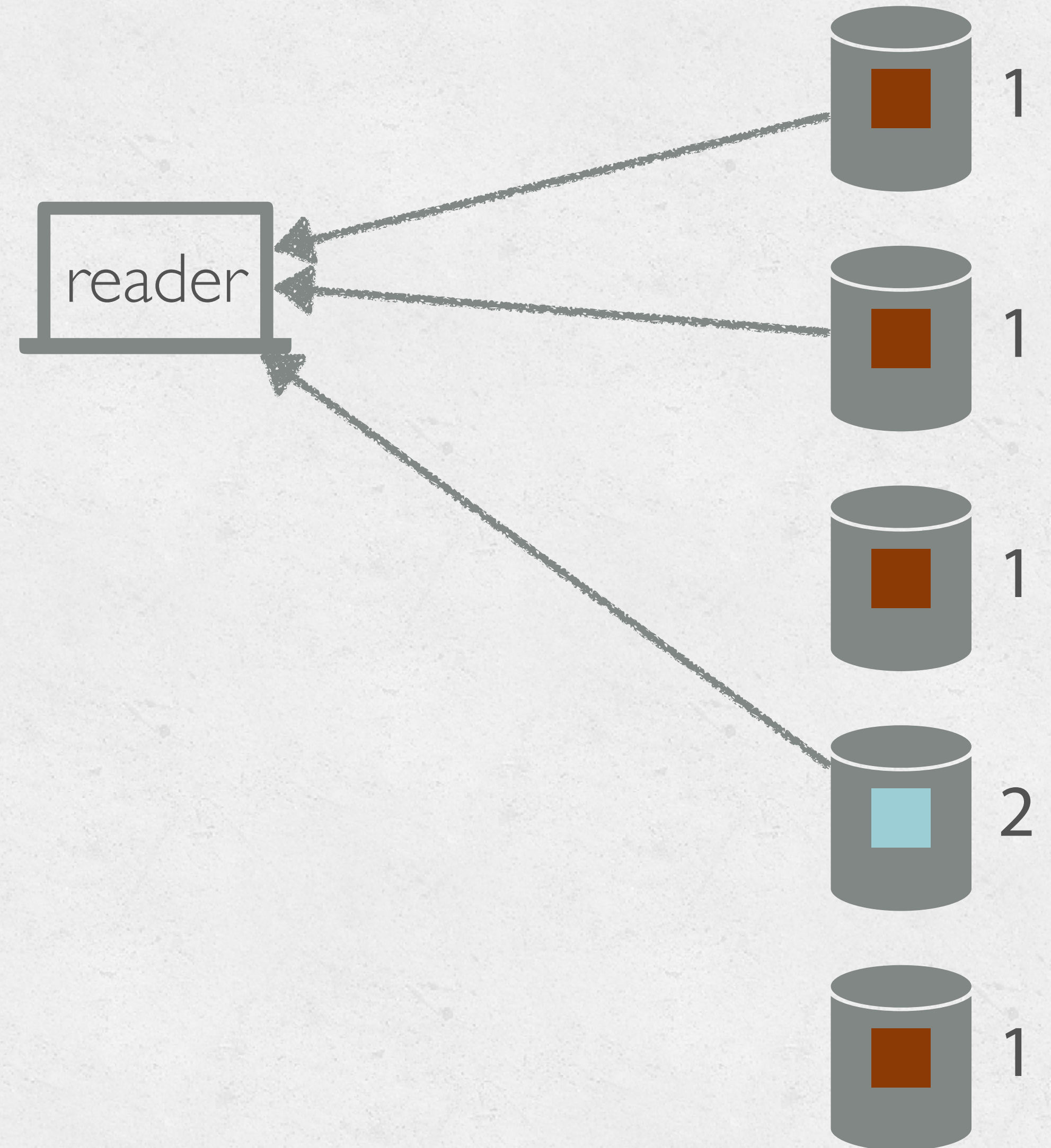
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW II

Suppose a write is ongoing (or the writer died).

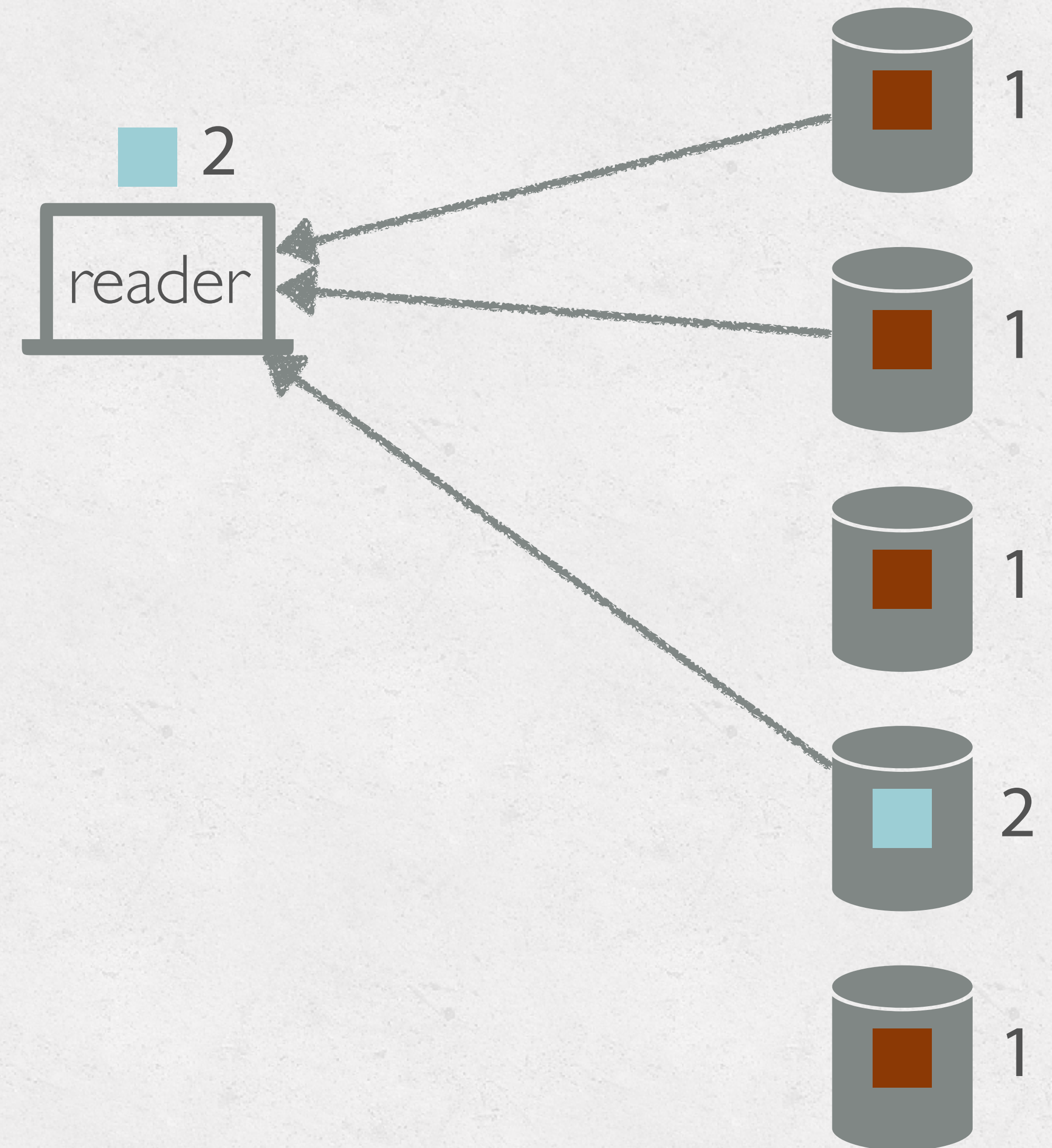
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW II

Suppose a write is ongoing (or the writer died).

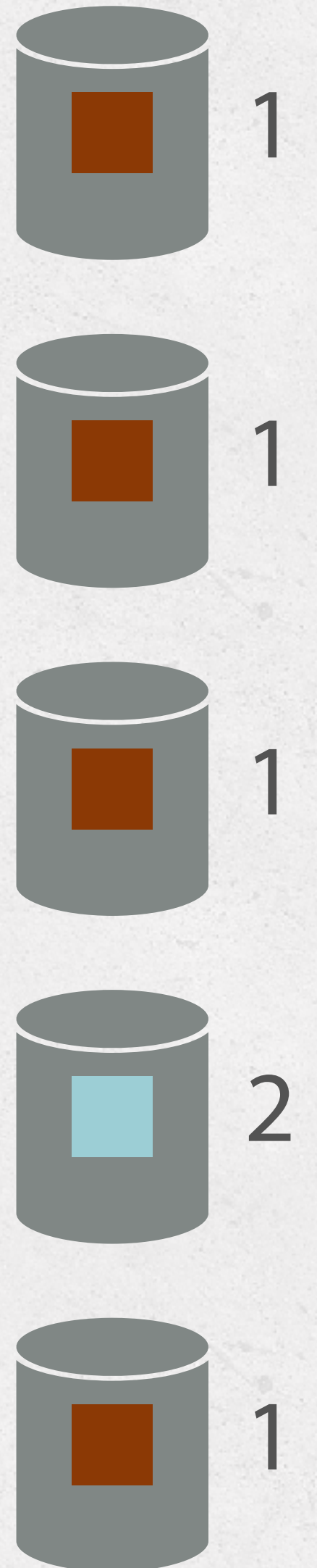
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW II

Suppose a write is ongoing (or the writer died).

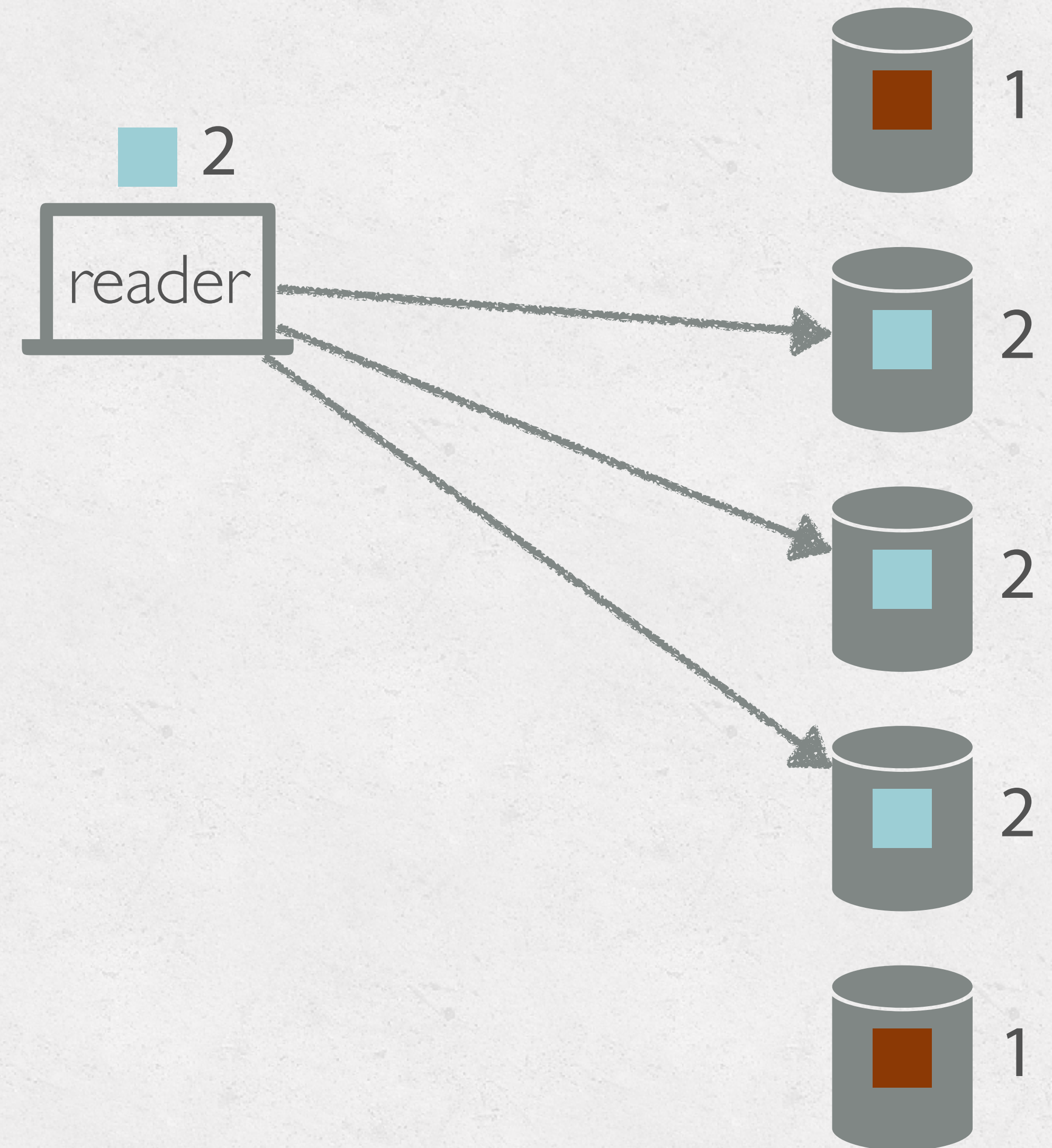
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW II

Suppose a write is ongoing (or the writer died).

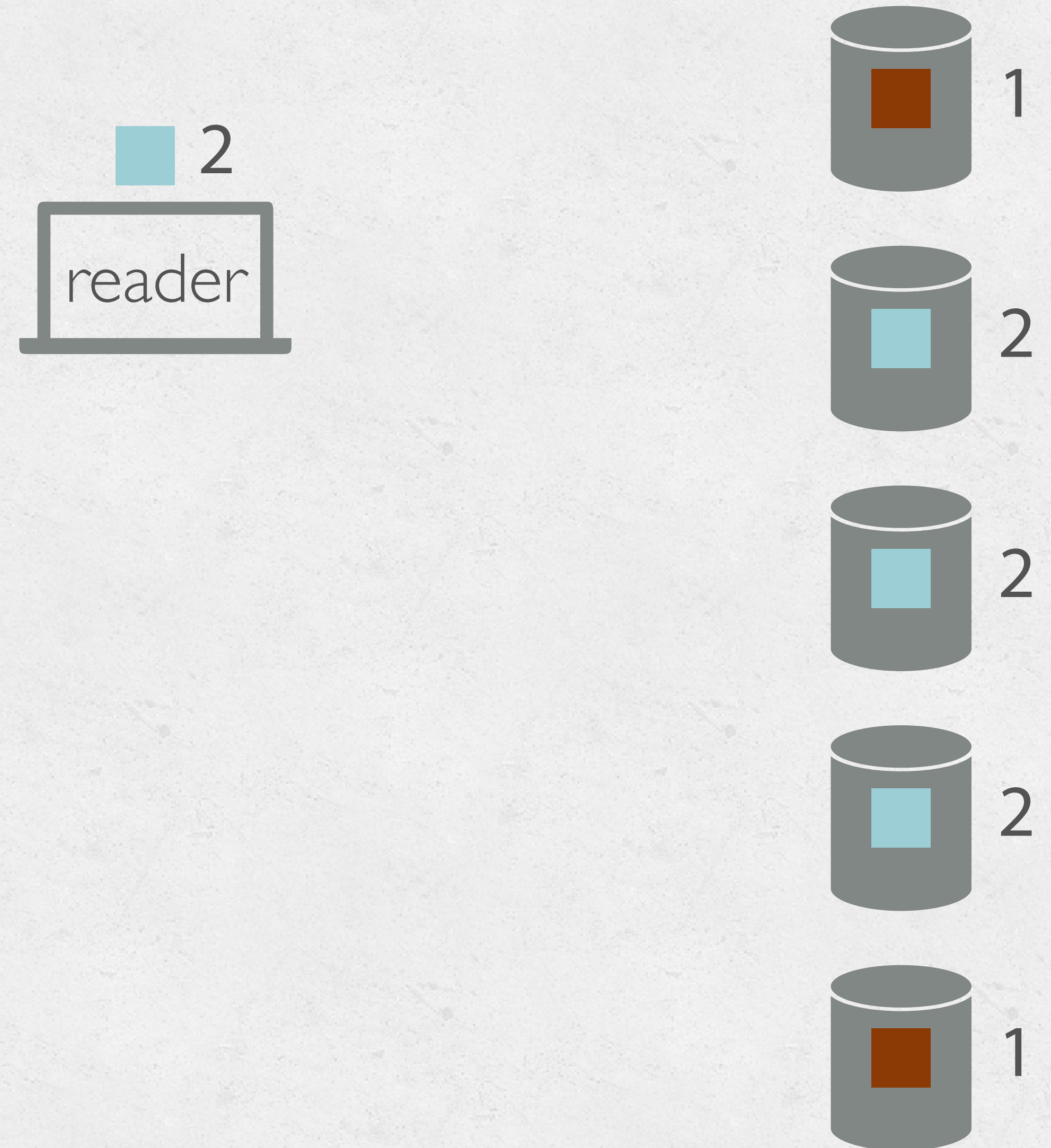
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW II

Suppose a write is ongoing (or the writer died).

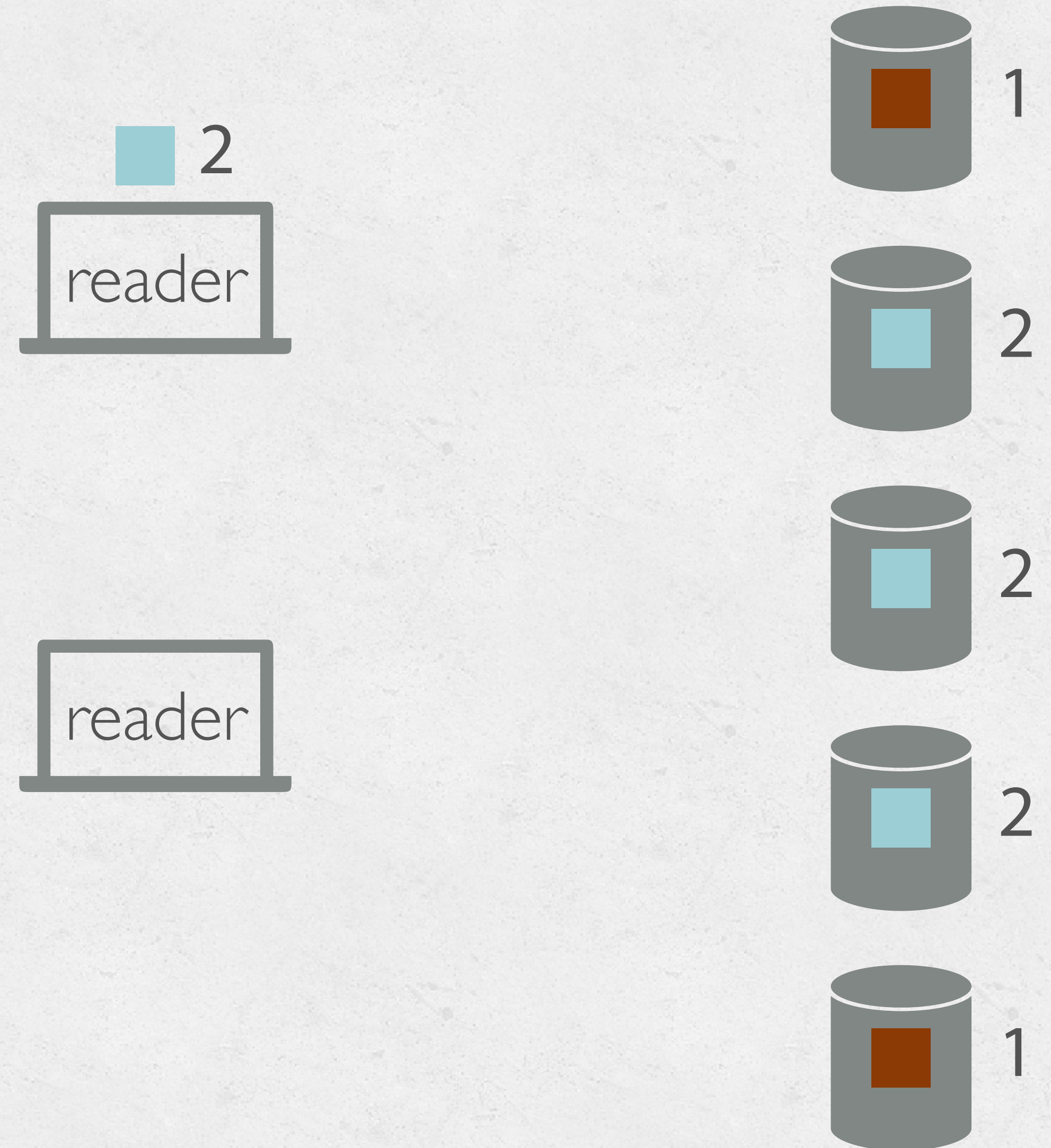
- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



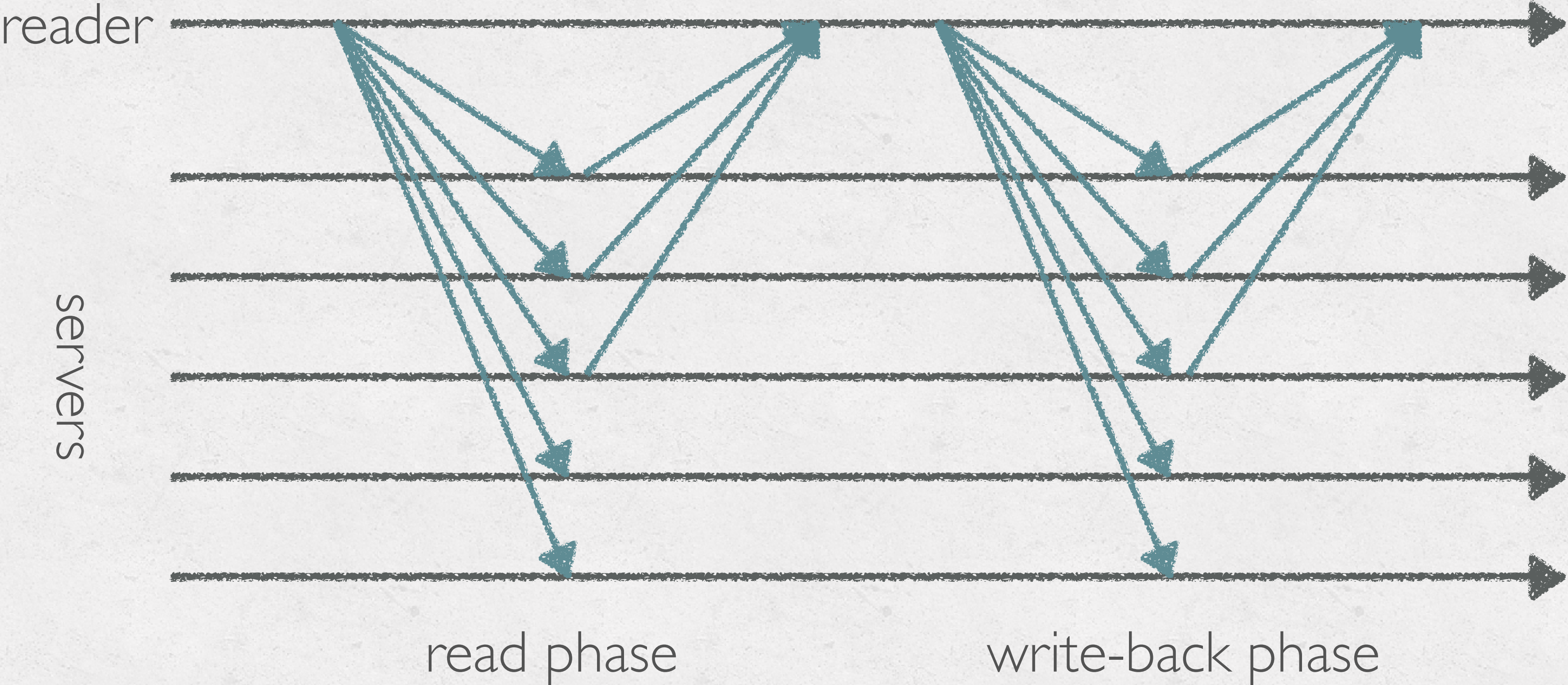
MRSW II

Suppose a write is ongoing (or the writer died).

- Reader reads value from a majority, takes the one with the highest timestamp.
- Reader then performs a **write-back**, writing the value to a majority (not necessarily the same one). Only returns from read after write-back is complete.
- Later readers are guaranteed to read a value **at least as new** as the previously returned one.



MRSW III



Do we always need to execute the write-back phase?

What if we only care about sequential consistency?

Do we care about the write-back phase at all?

PUTTING IT ALL TOGETHER: MRMW

Does the previous solution just work?

PUTTING IT ALL TOGETHER: MRMW

Does the previous solution just work?

What if writers use the **same timestamp**?

PUTTING IT ALL TOGETHER: MRMW

Does the previous solution just work?

What if writers use the **same timestamp**?

Prevented by
breaking ties using
writers ID, same as
PMMC.

PUTTING IT ALL TOGETHER: MRMW

Does the previous solution just work?

What if writers use the **same timestamp**?

What if a write that starts after a previous write ended uses a **smaller timestamp**?

Prevented by breaking ties using writers ID, same as PMMC.

MRMW: UNTIMELY TIMESTAMPS



MRMW: UNTIMELY TIMESTAMPS

Reads from a majority, sees blue value has the highest timestamp.



MRMW: UNTIMELY TIMESTAMPS

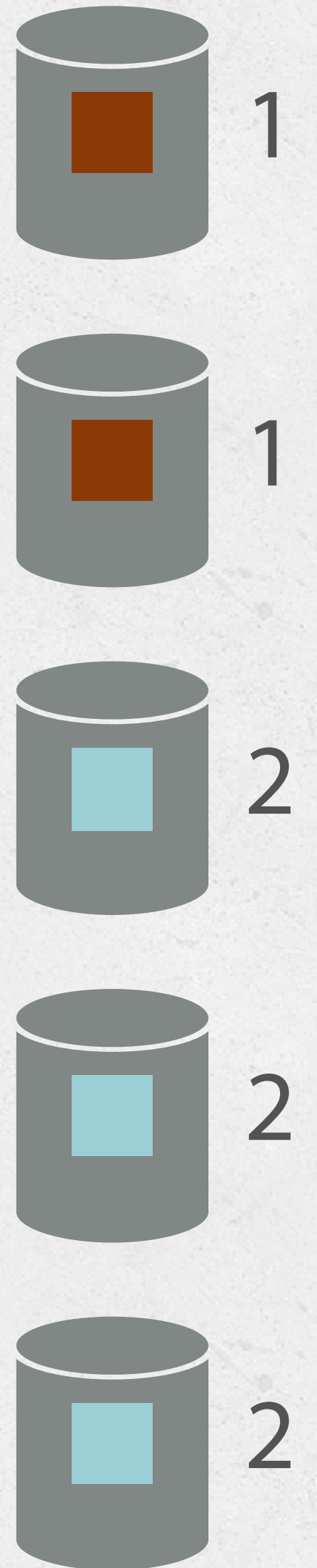
Reads from a majority, sees blue value has the highest timestamp.



Not linearizable!

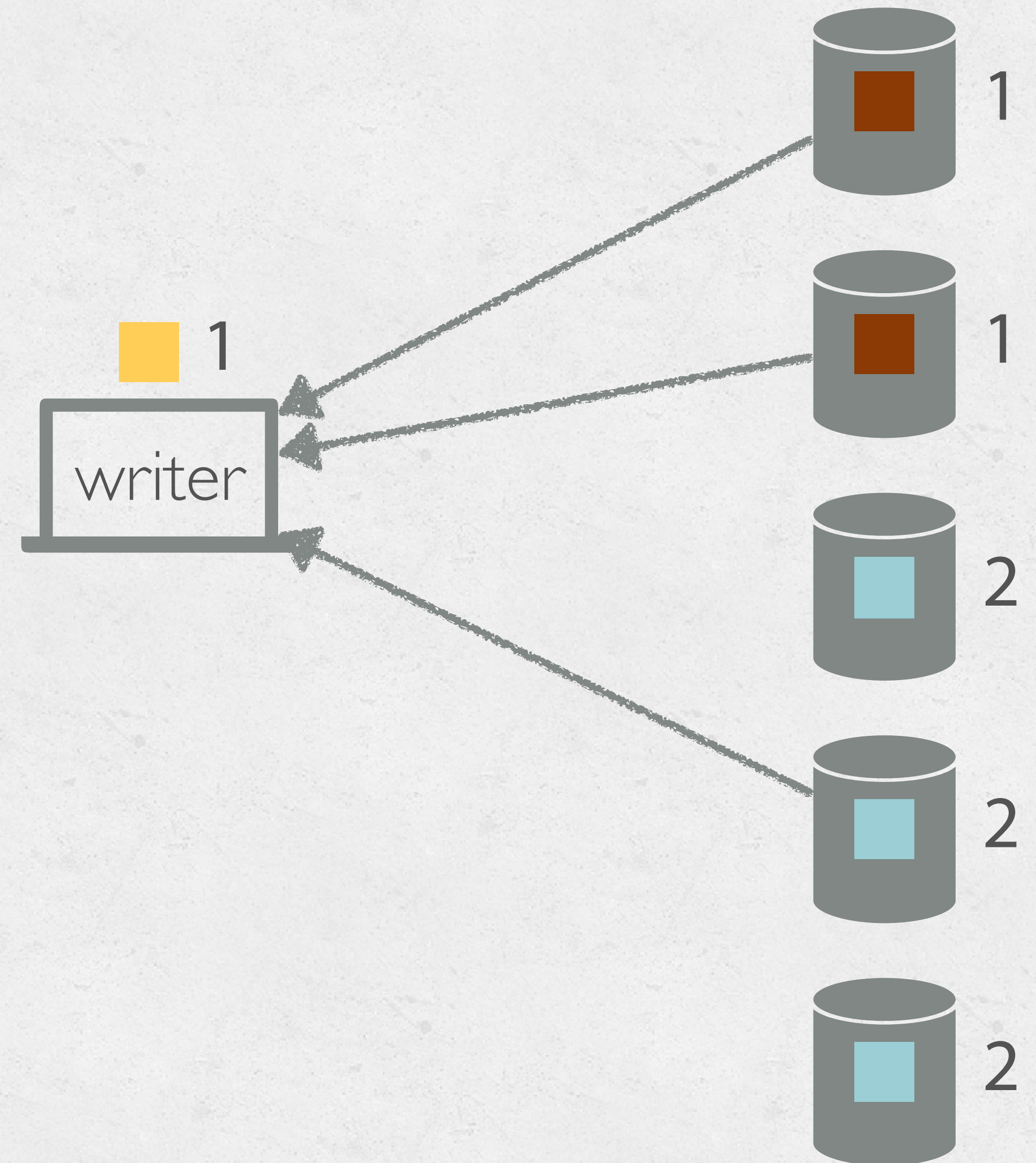
MRMW II: ENSURING TIMESTAMP ORDERING

- Writer first **queries** a majority, updates its timestamp to be larger than largest timestamp found.
- Writer then writes value to majority as usual.
- Written value guaranteed to have a timestamp larger than previously written values, readers will read latest value (again, writer IDs break timestamp ties).



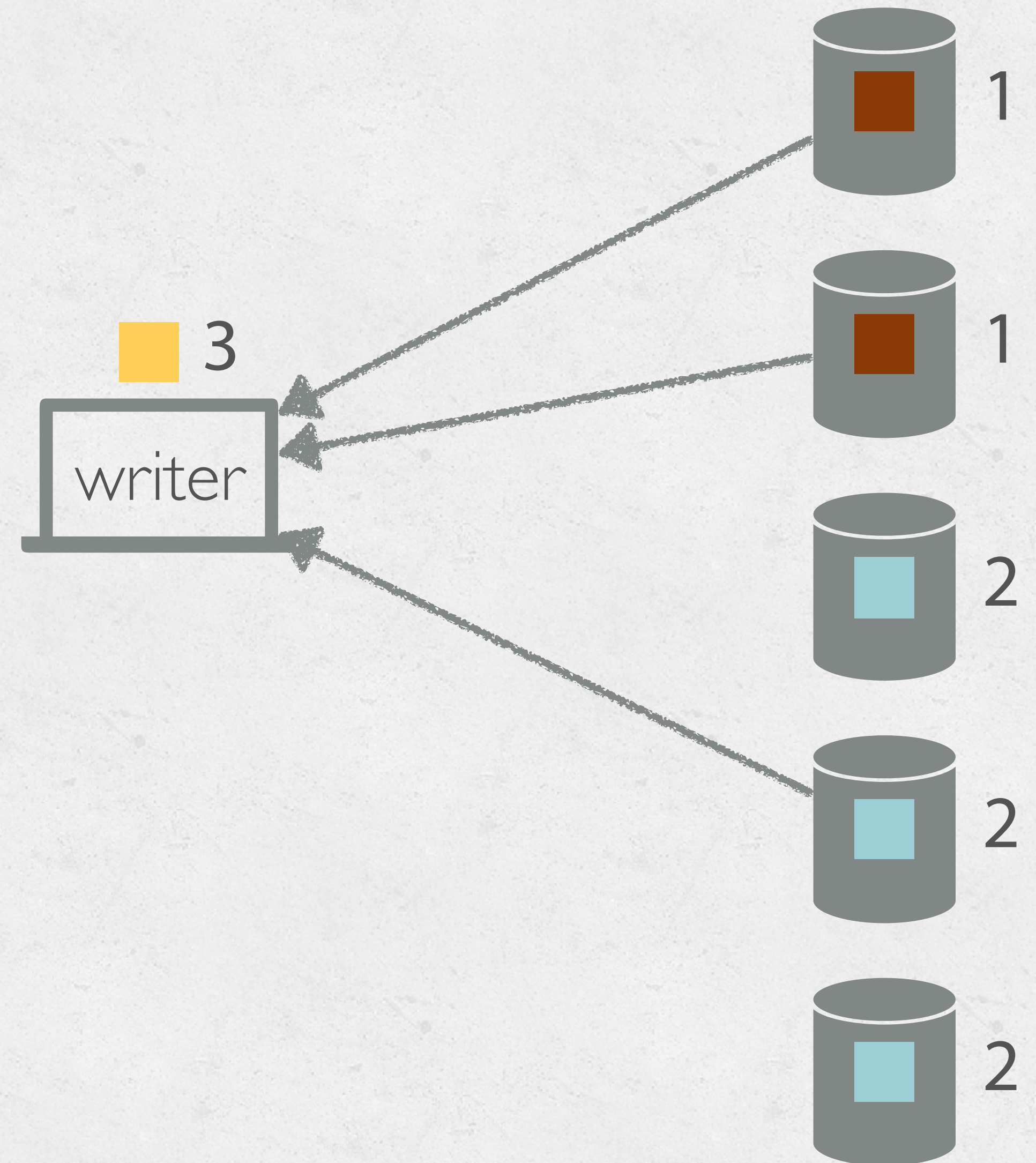
MRMW II: ENSURING TIMESTAMP ORDERING

- Writer first **queries** a majority, updates its timestamp to be larger than largest timestamp found.
- Writer then writes value to majority as usual.
- Written value guaranteed to have a timestamp larger than previously written values, readers will read latest value (again, writer IDs break timestamp ties).



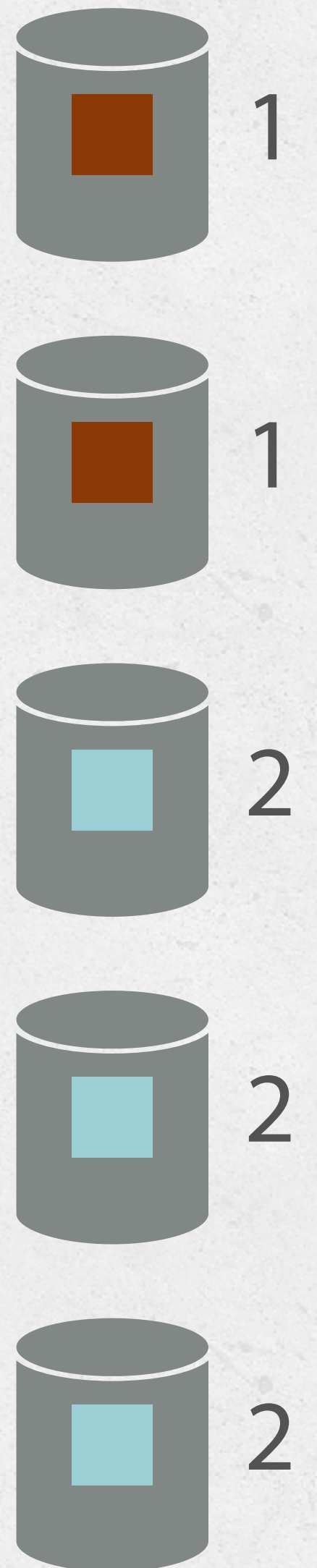
MRMW II: ENSURING TIMESTAMP ORDERING

- Writer first **queries** a majority, updates its timestamp to be larger than largest timestamp found.
- Writer then writes value to majority as usual.
- Written value guaranteed to have a timestamp larger than previously written values, readers will read latest value (again, writer IDs break timestamp ties).



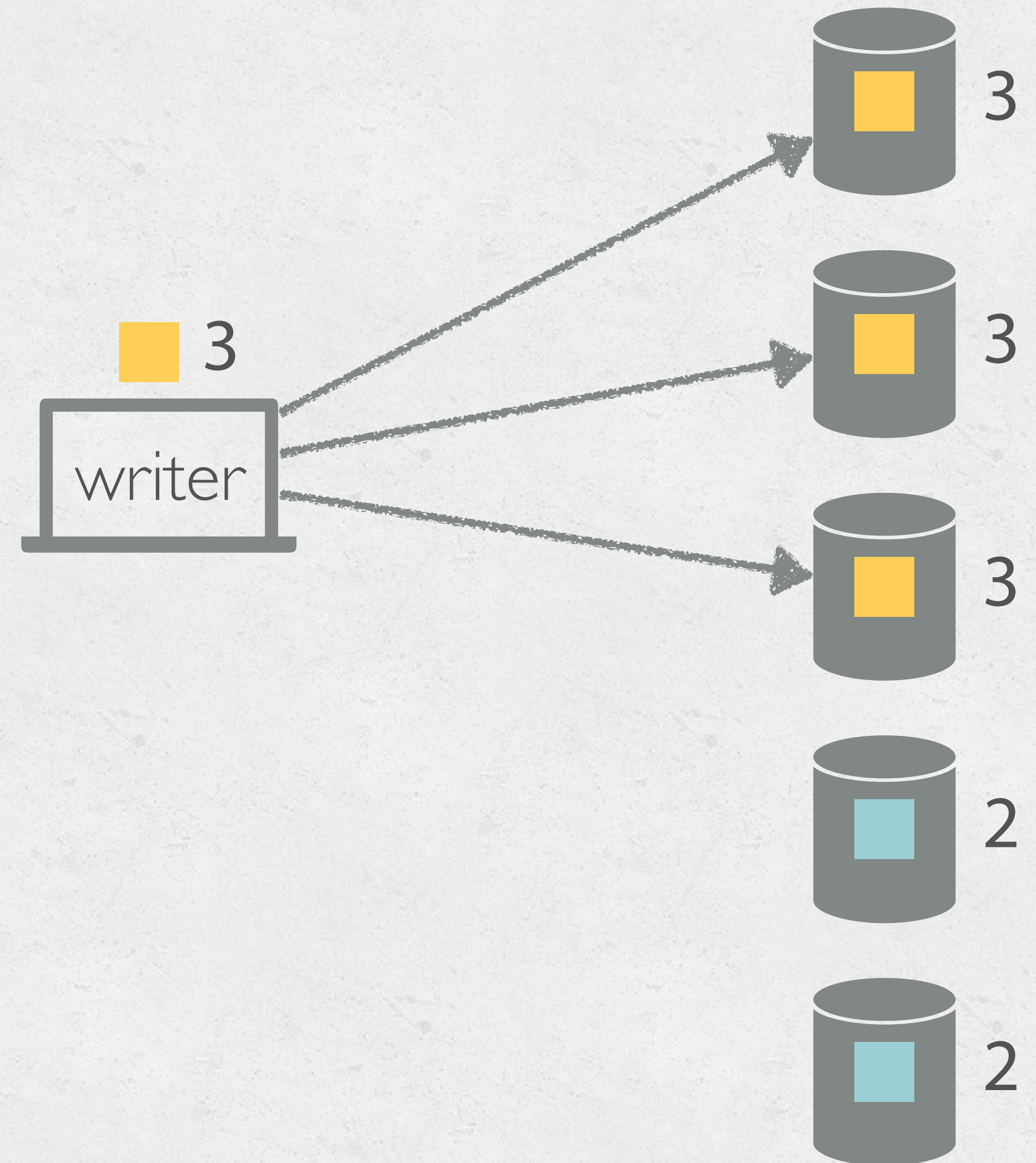
MRMW II: ENSURING TIMESTAMP ORDERING

- Writer first **queries** a majority, updates its timestamp to be larger than largest timestamp found.
- Writer then writes value to majority as usual.
- Written value guaranteed to have a timestamp larger than previously written values, readers will read latest value (again, writer IDs break timestamp ties).

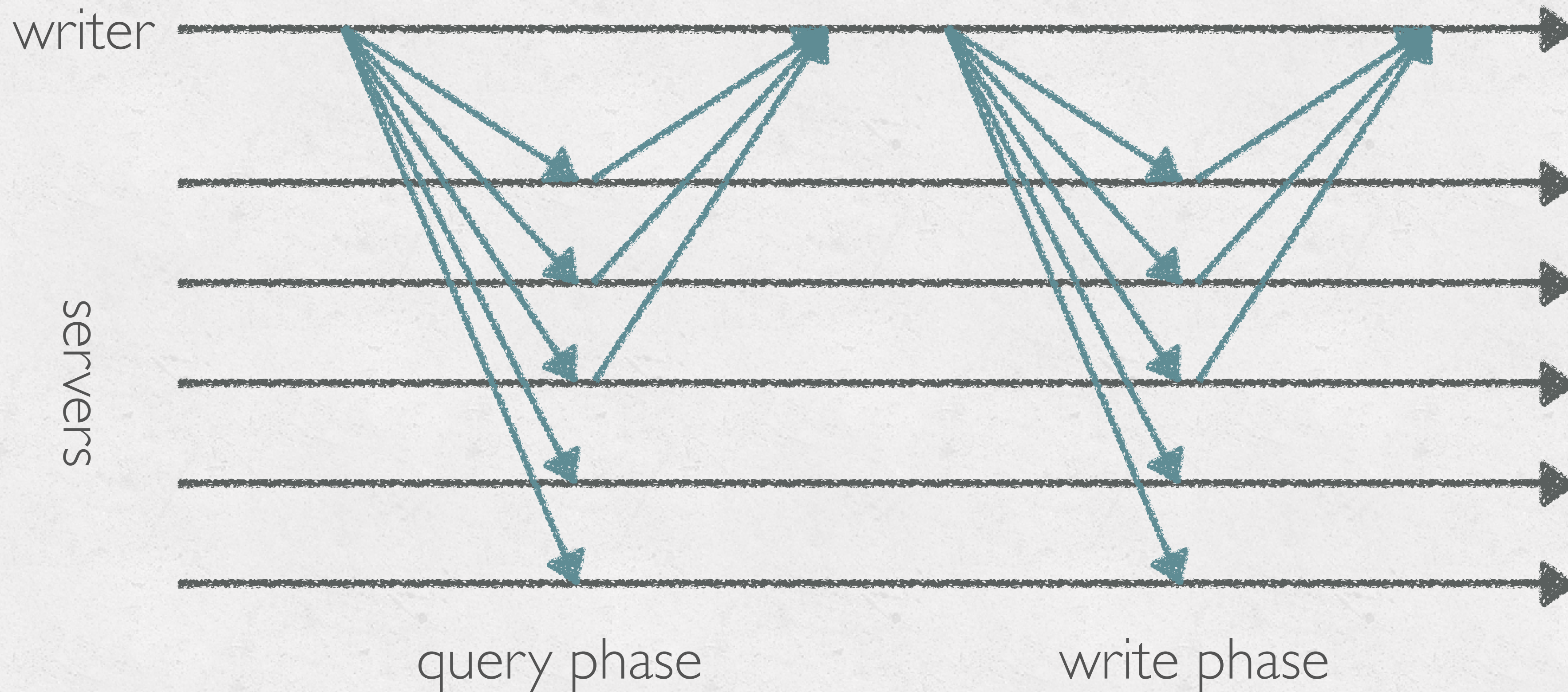


MRMW II: ENSURING TIMESTAMP ORDERING

- Writer first **queries** a majority, updates its timestamp to be larger than largest timestamp found.
- Writer then writes value to majority as usual.
- Written value guaranteed to have a timestamp larger than previously written values, readers will read latest value (again, writer IDs break timestamp ties).



MRMW III



WAIT A SECOND!

- The methods for reading and writing are now the exact same.
- The only difference is that a read writes and returns the value that was read, but a write writes the value to be written.
- Also, for the record, there's no reason that processes can't be both readers and writers.

Attiya, Bar-Noy, Dolev 1995 ABD Algorithm

ABD vs. PAXOS

- Paxos doesn't guarantee liveness when the network is asynchrony. ABD guarantees wait-freedom, even when there are multiple writers.
- Paxos-based state-machine replication (SMR) can support arbitrary state machines. The ABD algorithm only allows a read/write interface.
- ABD removes the leader bottleneck, has the same latency cost as leader-based Paxos.

WHAT CAN WE DO WITH REGISTERS?

- Implement a read/write key-value store.
- Emulate shared memory.

Consensus isn't always the right problem! Don't solve it if you don't have to!