

GroveKV is:

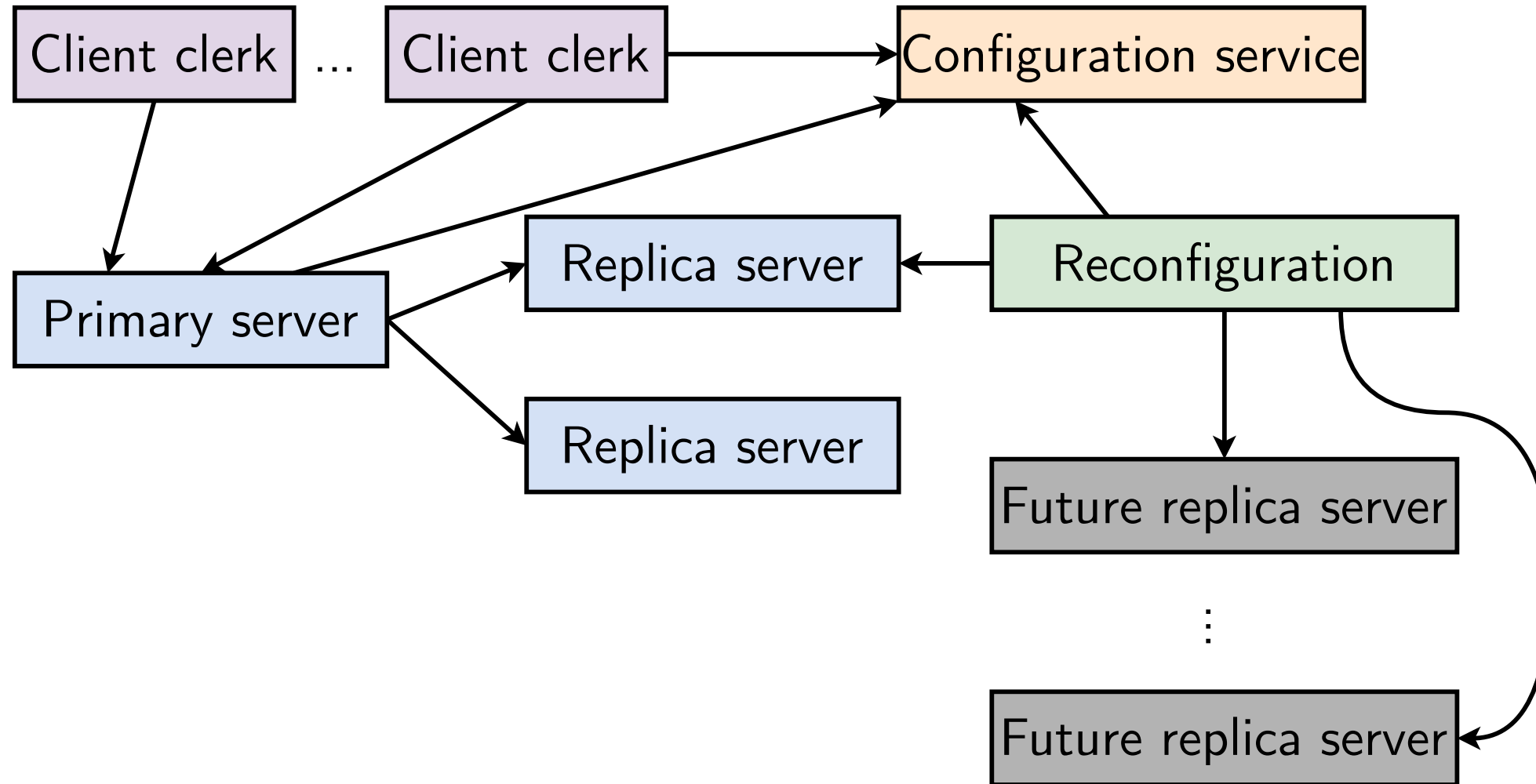
a **fault-tolerant**, **linearizable key-value service**

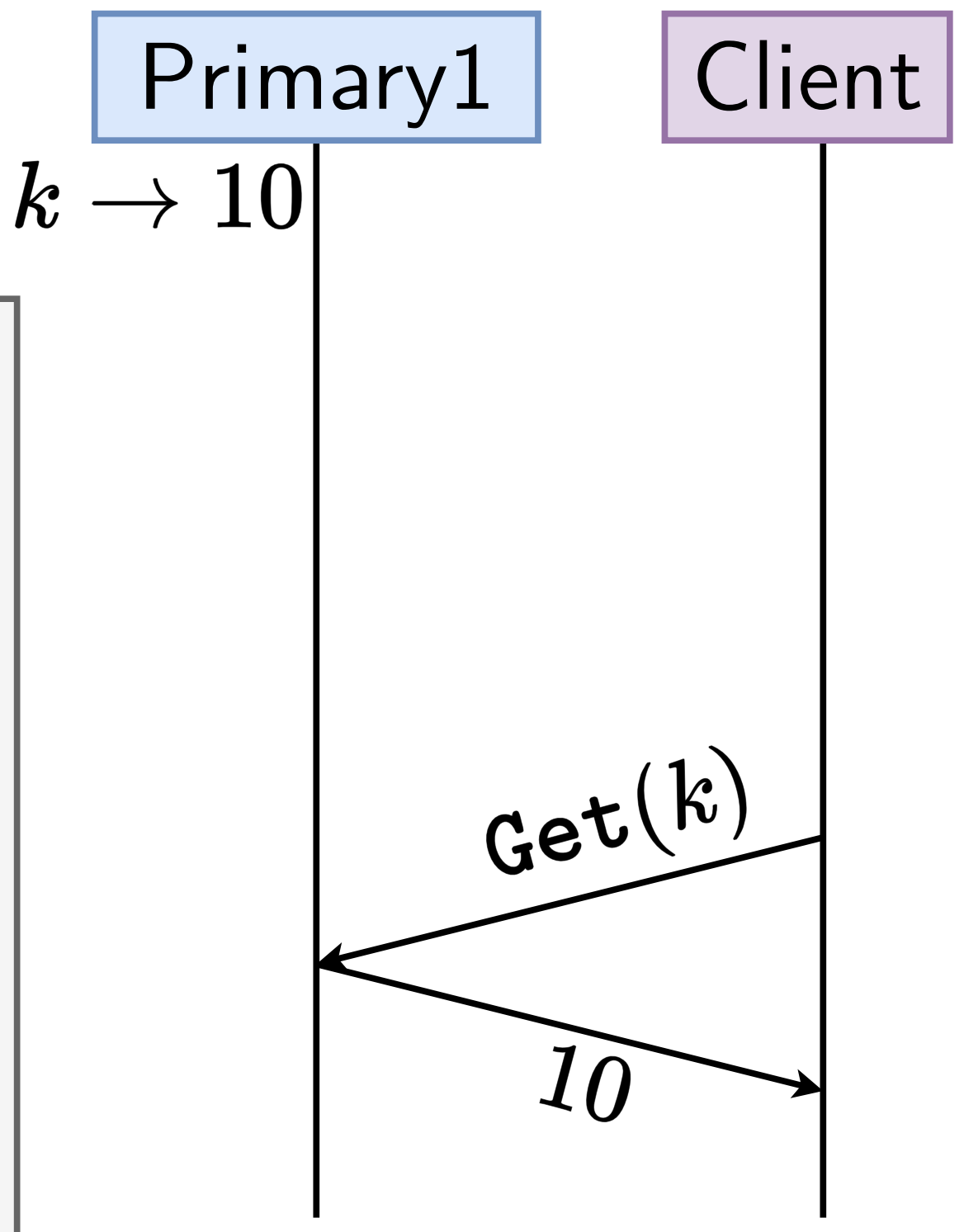
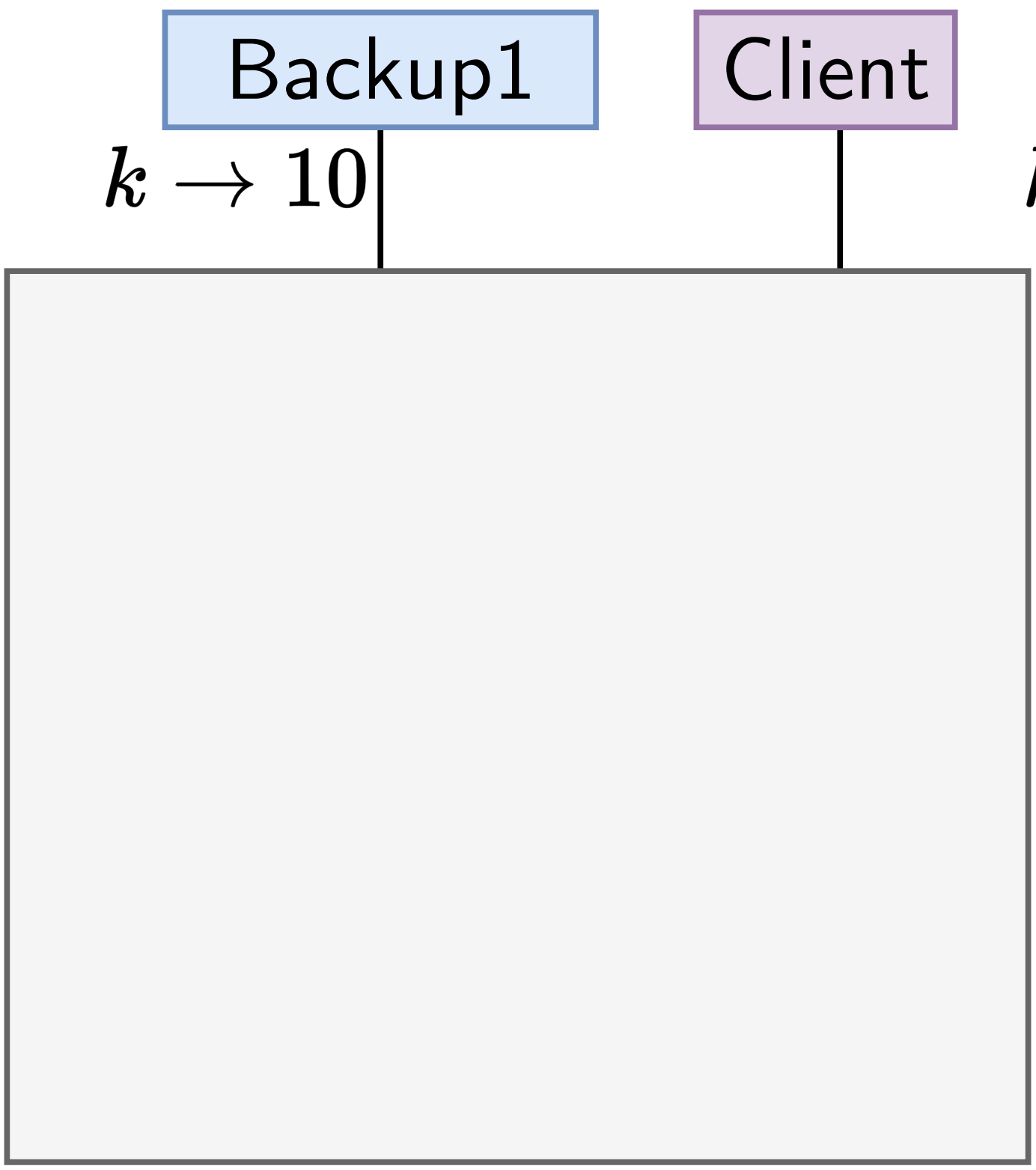
Put (k, v) and Get (k)

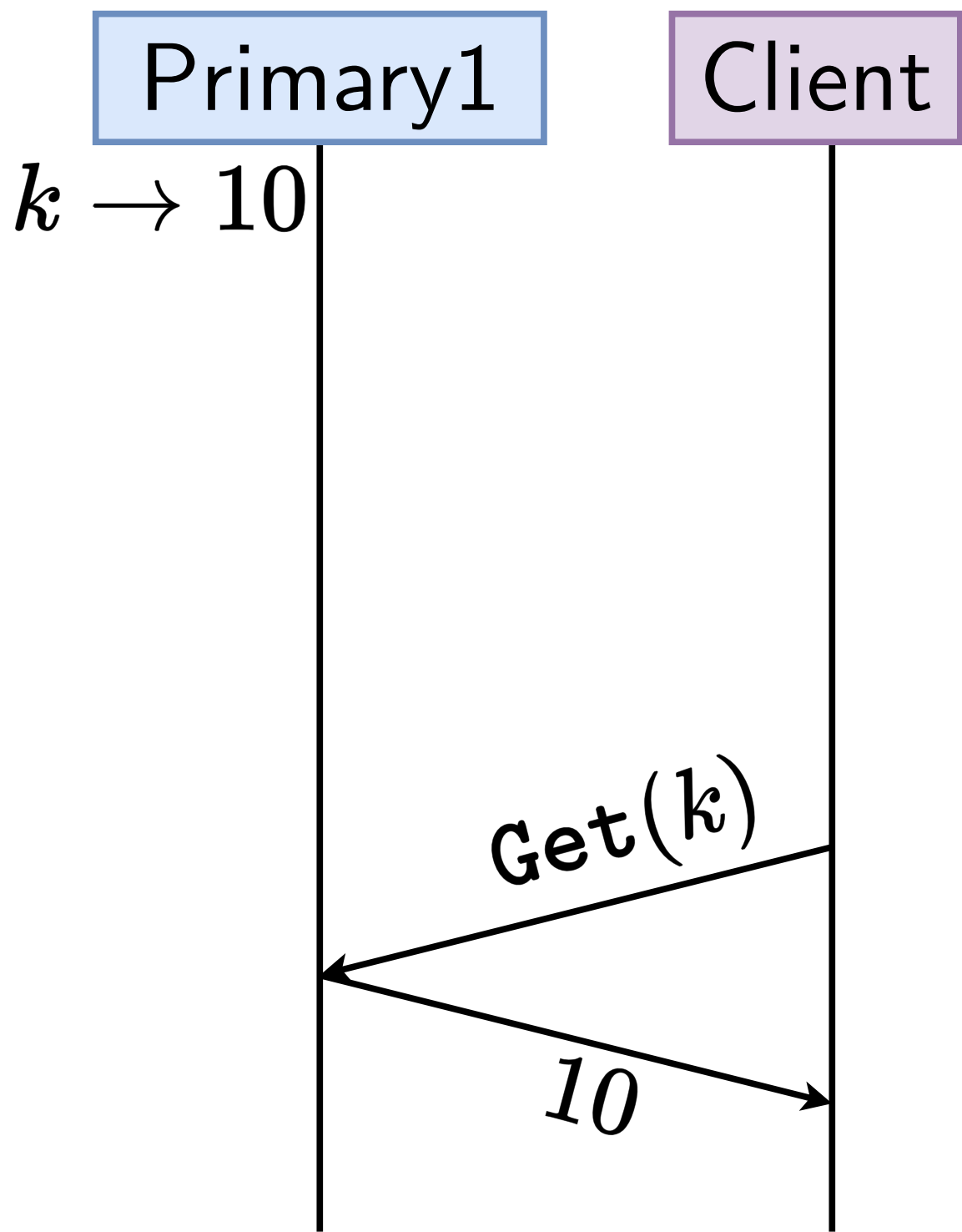
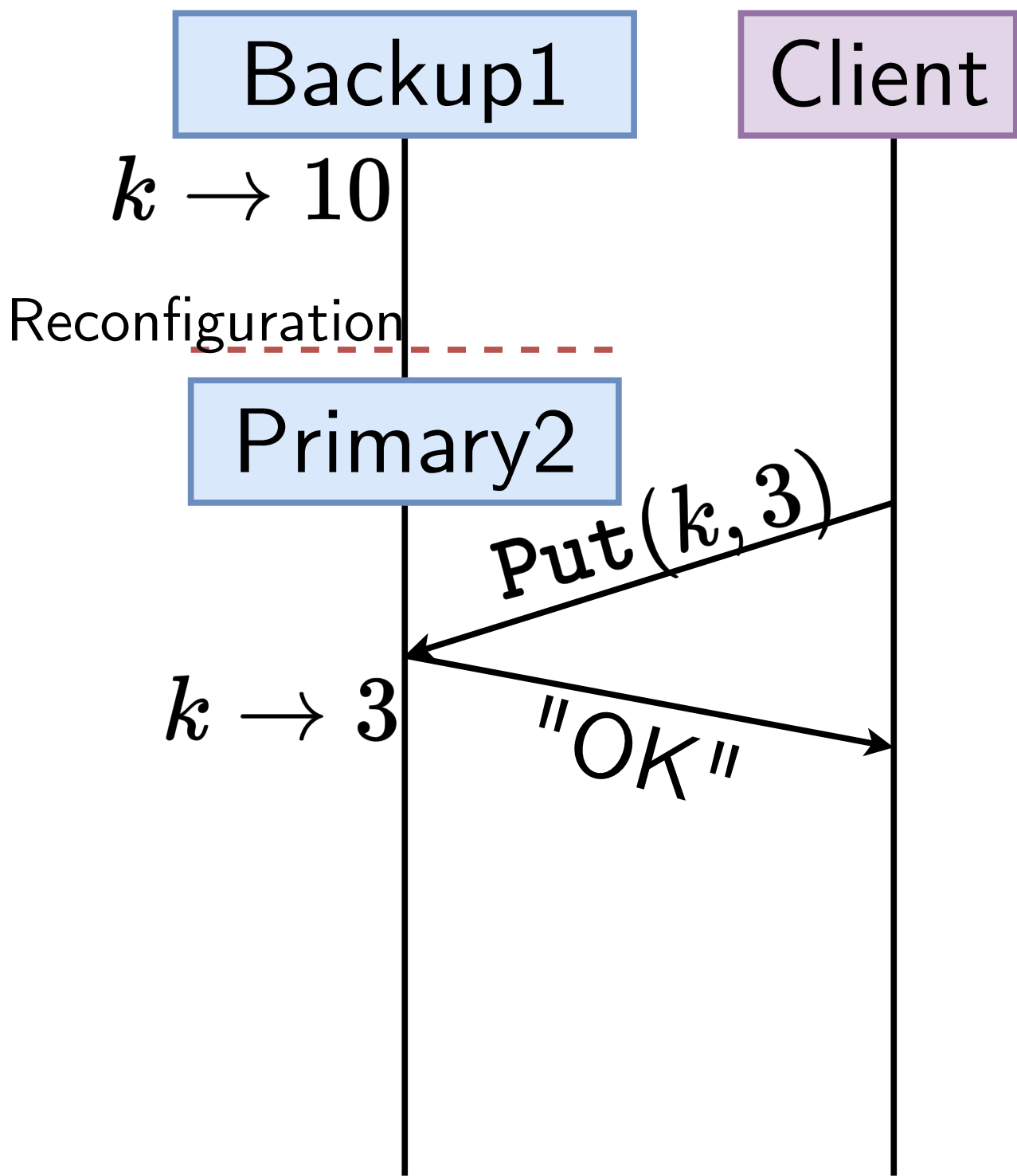
Crash-safe and reconfigurable

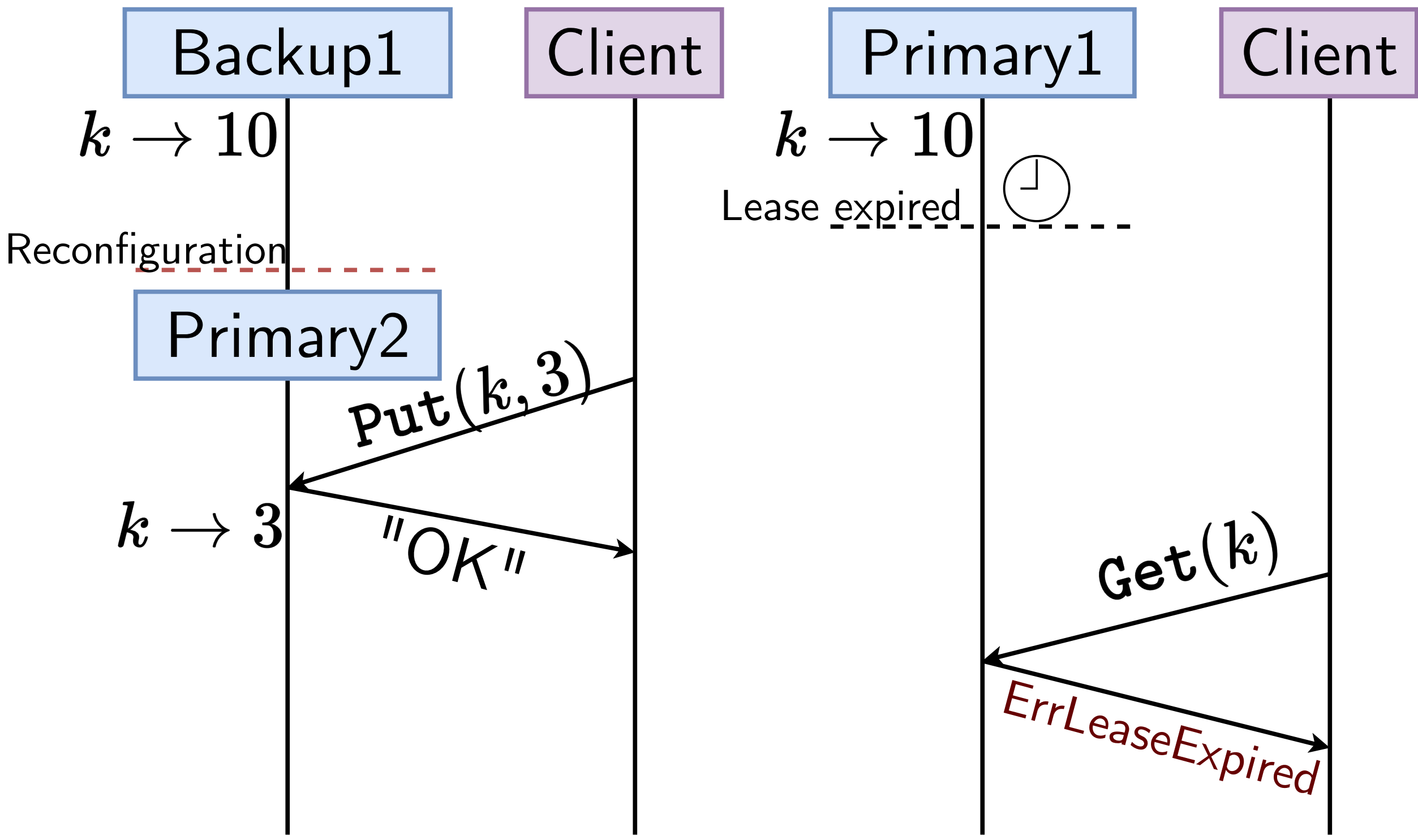
+ Lease-based reads

How does GroveKV work?









Want to do all this correctly

Common approach: tests

E.g.

```
...  
clerk.Put("x", "10")  
clerk.Put("x", "3")  
assert(clerk.Get("x") == "3")
```

Want to do all this correctly

Challenge: lots of "what if" questions

- What if there is a `Get` on a server concurrently with `Put`?
- What if backup crashes and loses operations?
- What if lease expires in the middle of a `Get`?
- What if servers try operations (`Get`? `Put`?) during reconfiguration?
- What if ...?



Want to know:

For any thread interleaving and crashes etc., the code behaves correctly.



\forall thread interleaving and crashes etc., the code behaves correctly.

Formal verification

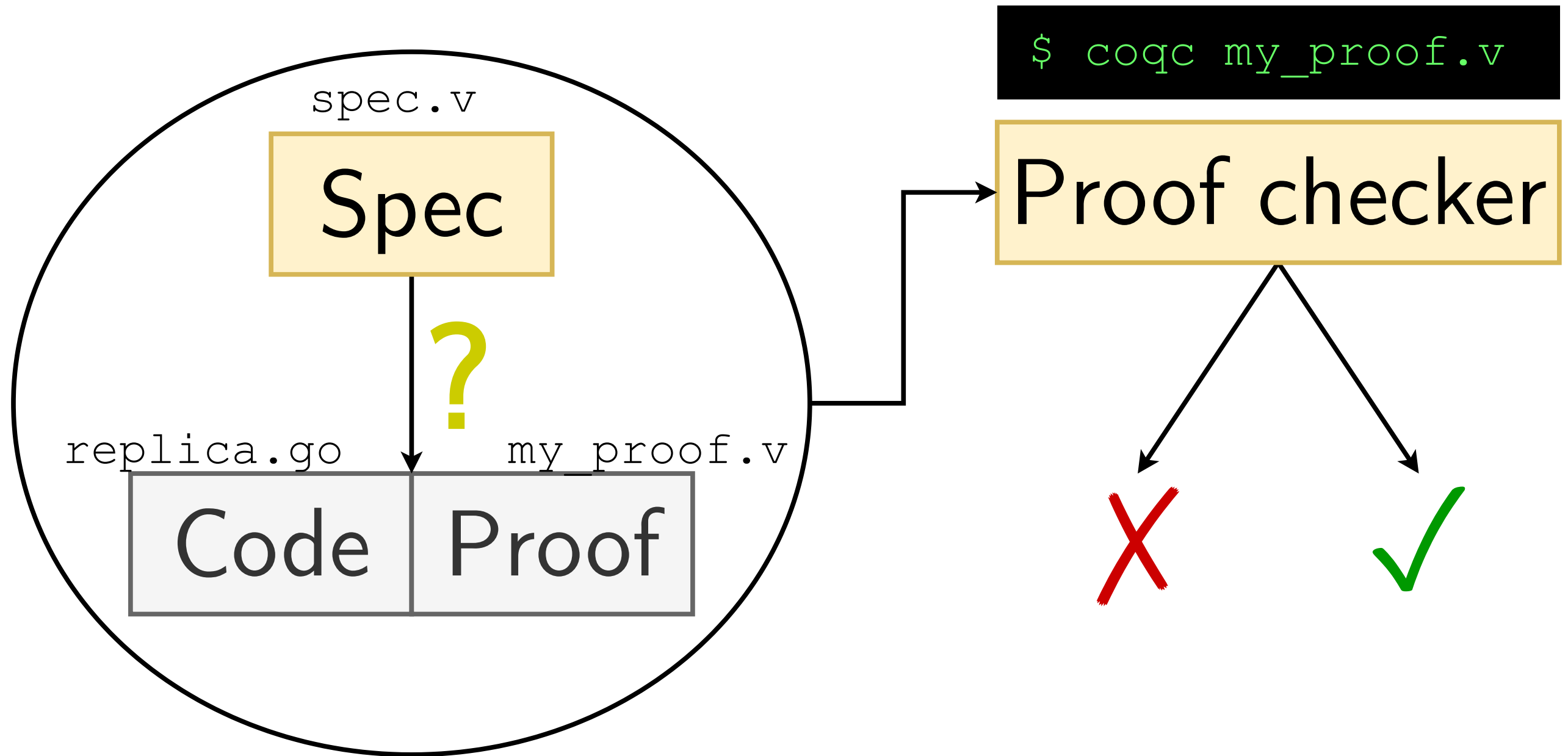
Theorem. \forall thread interleaving and crashes etc., the code behaves correctly.

Need a "mathematical model" for how code executes

Need a definition for "correctly"

Use mechanized proof checker

Formal verification



Grove

A framework for proving " \forall theorems",
without explicitly writing proofs about
interleavings, etc.

Verify system by going line-by-line 🧐

Grove

Specs are pre+postconditions (Hoare logic)

E.g. $\{\top\}$ `sort(a)` $\{\text{RET } b, \text{is_sorted}(b)\}$

Precondition

Code

Postcondition

$\{\text{is_sorted } a\}$ `bsearch(a, x)` $\{\text{RET } i, a[i] = x\}$

How to specify GroveKV?

Maybe:

$\{\top\}$ `clerk.Get(k)` $\{\text{RET } v, \text{kvmap}[k] = v\}$?

$\{\top\}$ `clerk.Put(k, v)` $\{\text{kvmap}[k] = v\}$?

What exactly is `kvmap`? Does this deal with concurrent operations?

Concurrent separation logic: *ownership of resources*

Concurrent Separation Logic

Pre/post condition describe what is logically *owned* by the running code

"Heap points-to"

$x \mapsto v$ denotes ownership of address x

$x \mapsto v * y \mapsto w$ denotes separate ownership

Similar to ownership in Rust

GroveKV top-level spec

$\{k \mapsto_{\text{kv}} w\}$ `clerk.Put(k, v)` $\{k \mapsto_{\text{kv}} v\}$

$\{k \mapsto_{\text{kv}} v\}$ `clerk.Get(k)` $\{\text{ret } v, k \mapsto_{\text{kv}} v\}$

$x \mapsto_{\text{kv}} \text{"0"} * y \mapsto_{\text{kv}} \text{"0"}$

`ck.Put(x, "1")`

`ck.Put(y, "2")`

`go f(...)`

`ck.Put(x, "5")`

`assert(ck.Get(x) == "5")`

`ck.Put(y, "3")`

`assert(ck.Get(y) == "3")`



$x \mapsto_{kv} "0"$ * $y \mapsto_{kv} "0"$

$ck.Put(x, "1") \longrightarrow x \mapsto_{kv} "1" * y \mapsto_{kv} "0"$

$ck.Put(y, "2")$
 $go f(...)$

$ck.Put(x, "5")$

$assert(ck.Get(x) == "5")$

$ck.Put(y, "3")$

$assert(ck.Get(y) == "3")$

$x \mapsto_{kv} "0" * y \mapsto_{kv} "0"$

`ck.Put(x, "1")`

$x \mapsto_{kv} "1" * y \mapsto_{kv} "0"$

`ck.Put(y, "2")`

$x \mapsto_{kv} "1" * y \mapsto_{kv} "2"$

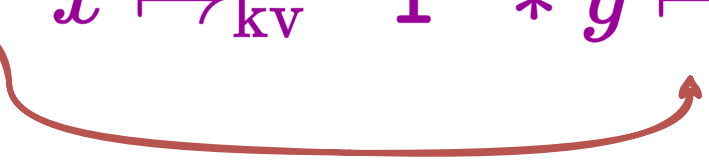
`go f(...)`

`ck.Put(x, "5")`

`ck.Put(y, "3")`

`assert(ck.Get(x) == "5")`

`assert(ck.Get(y) == "3")`



$x \mapsto_{kv} "0" * y \mapsto_{kv} "0"$

`ck.Put(x, "1")` $x \mapsto_{kv} "1" * y \mapsto_{kv} "0"$

`ck.Put(y, "2")` $x \mapsto_{kv} "1" * y \mapsto_{kv} "2"$

`go f(...)`

$x \mapsto_{kv} "1"$

`ck.Put(x, "5")`

$x \mapsto_{kv} "5"$

`assert(ck.Get(x) == "5")`

$x \mapsto_{kv} "5"$

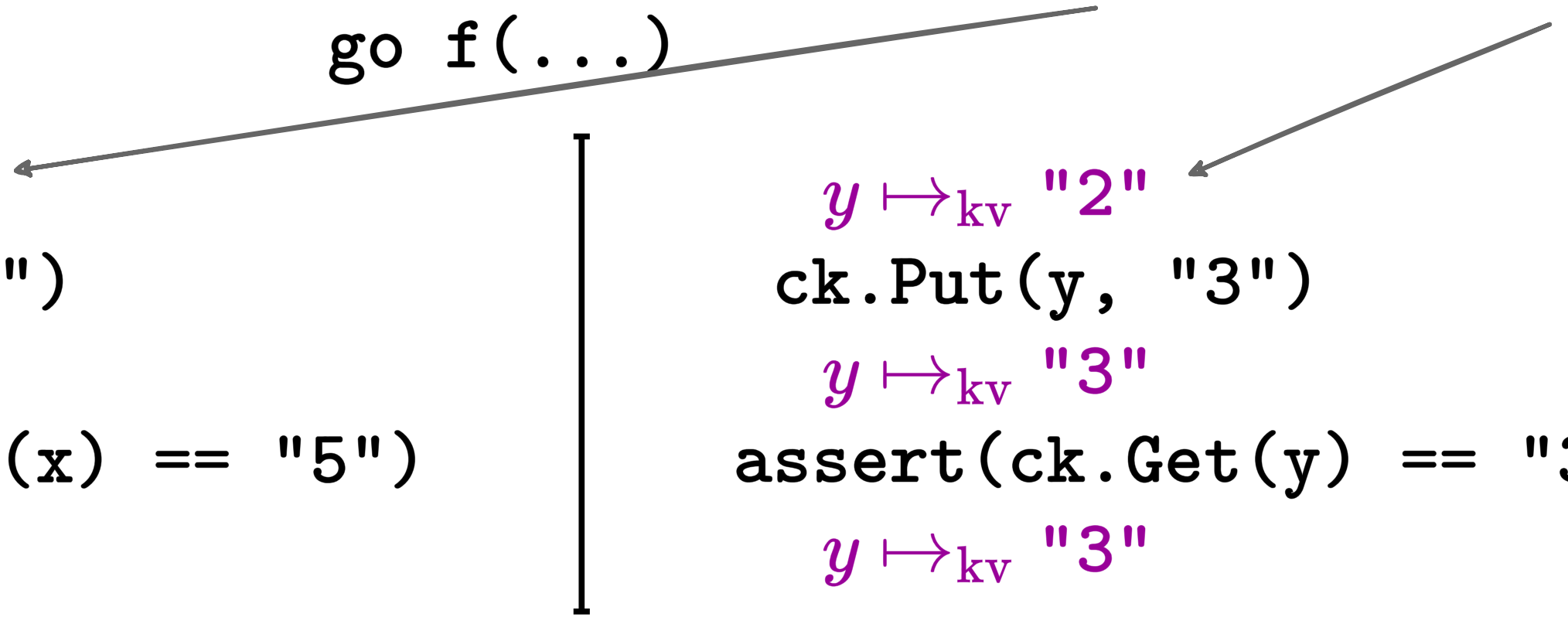
$y \mapsto_{kv} "2"$

`ck.Put(y, "3")`

$y \mapsto_{kv} "3"$

`assert(ck.Get(y) == "3")`

$y \mapsto_{kv} "3"$



What about specs for RPCs?

`{...} ApplyAsBackupRPC(e, index, op) {???`

Idea from CSL: *ghost resources*

Ghost state useful for specification & proof

Append-only list ghost resource

$a \stackrel{\text{list}}{\vdash} \ell$: ownership of append-only list

$a \stackrel{\text{list}}{\sqsubseteq} \ell$: knowledge of list lower-bound

$a \stackrel{\text{list}}{\vdash} \square \ell$: knowledge of frozen list

Primary/backup resources

Server j owns $\text{accepted}_j[e] \stackrel{\text{list}}{\mapsto} ops$

Replication invariant:

$$\exists l \ e, \text{committed} \stackrel{\text{list}}{\mapsto} l^*$$
$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} l^* \cdots * \text{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} l$$

ApplyAsBackupRPC spec

{...}

`server[j].ApplyAsBackupRPC(e, index, op)`

$\left\{ \text{accepted}_j[e] \stackrel{\text{list}}{=} ops + [op] \right\}$

Captures the "promise" made by RPC

Proof code.

```
mutex_lock.  
....  
update (a:=ops+[op]).  
get_lower_bound a.  
mutex_unlock with a.  
applyBackup_spec.  
applyBackup_spec.  
...  
applyBackup_spec.  
open  $I_{rep}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{rep}$ .  
done.
```

Qed.

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 0; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```


Proof code.
mutex lock.

```
....  
update (a:=ops+[op]).  
get_lower_bound a.  
mutex_unlock with a.  
applyBackup_spec.  
applyBackup_spec.  
...  
applyBackup_spec.  
open  $I_{rep}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{rep}$ .  
done.
```

Qed.

$\text{accepted}_0[e] \stackrel{\text{list}}{\mapsto} ops$

...

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 0; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```

```
Proof code.  
mutex_lock.  
....
```

```
update (a:=ops+[op]).  
get_lower_bound a.  
mutex_unlock with a.  
applyBackup_spec.  
applyBackup_spec.  
...  
applyBackup_spec.  
open  $I_{\text{rep}}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{\text{rep}}$ .  
done.
```

Qed.

$\text{accepted}_0[e] \stackrel{\text{list}}{\mapsto} ops$

...

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 0; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```

```
Proof code.  
mutex_lock.
```

```
....  
update (a:=ops+[op]).
```

```
get_lower_bound a.  
mutex_unlock with a.
```

```
applyBackup_spec.  
applyBackup_spec.
```

```
...
```

```
applyBackup_spec.
```

```
open  $I_{rep}$ .
```

```
update (c:=ops+[op]).
```

```
get_lower_bound c.
```

```
close  $I_{rep}$ .
```

```
done.
```

```
Qed.
```

```
accepted0[e]  $\stackrel{\text{list}}{\mapsto}$  ops + [op]
```

```
...
```

```
s.mutex.Lock()
```

```
nextIndex := s.nextIndex
```

```
e := s.epoch
```

```
s.nextIndex += 1
```

```
res := s.stateLogger.LocalApply(op)
```

```
s.mutex.Unlock()
```

```
for j := 0; j < len(s.backupClerks); j++ {
```

```
    s.backupClerks[j].ApplyAsBackupRPC(e,
```

```
nextIndex, op)
```

```
}
```

```
return res
```

```
Proof code.  
mutex_lock.  
....  
update (a:=ops+[op]).  
get_lower_bound a.
```

```
mutex_unlock with a.  
applyBackup_spec.  
applyBackup_spec.  
...  
applyBackup_spec.  
open  $I_{rep}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{rep}$ .  
done.  
Qed.
```

$\text{accepted}_0[e] \stackrel{\text{list}}{\mapsto} ops + [op]$

$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$

...

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 0; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```

```
Proof code.  
mutex_lock.  
....  
update (a:=ops+[op]).  
get_lower_bound a.  
mutex_unlock with a.
```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\equiv} ops + [op]$$

```
applyBackup_spec.  
applyBackup_spec.  
...  
applyBackup_spec.  
open  $I_{\text{rep}}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{\text{rep}}$ .  
done.  
Qed.
```

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 0; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```

```
Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
```

```
applyBackup_spec.
...
applyBackup_spec.
open  $I_{\text{rep}}$ .
update (c:=ops+[op]).
get_lower_bound c.
close  $I_{\text{rep}}$ .
done.
Qed.
```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} \text{ops} + [op]$$
$$\text{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} \text{ops} + [op]$$

...

```
s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 1; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res
```

```
Proof code.  
mutex_lock.  
....  
update (a:=ops+[op]).  
get_lower_bound a.  
mutex_unlock with a.  
applyBackup_spec.  
applyBackup_spec.
```

```
...  
applyBackup_spec.  
open  $I_{rep}$ .  
update (c:=ops+[op]).  
get_lower_bound c.  
close  $I_{rep}$ .  
done.  
Qed.
```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$
$$\text{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$
$$\text{accepted}_2[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```
s.mutex.Lock()  
nextIndex := s.nextIndex  
e := s.epoch  
s.nextIndex += 1  
res := s.stateLogger.LocalApply(op)  
s.mutex.Unlock()  
  
for j := 2; j < len(s.backupClerks); j++ {  
    s.backupClerks[j].ApplyAsBackupRPC(e,  
nextIndex, op)  
}  
  
return res
```

```

Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.

```

```

open  $I_{\text{rep}}$ .
update (c:=ops+[op]).
get_lower_bound c.
close  $I_{\text{rep}}$ .
done.
Qed.

```

$$\mathbf{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\mathbf{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\dots$$

$$\mathbf{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res

```



```

Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.
open  $I_{\text{rep}}$ .

```

```

update (c:=ops+[op]).
get_lower_bound c.
close  $I_{\text{rep}}$ .
done.
Qed.

```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \text{committed} \stackrel{\text{list}}{\mapsto} ops$$

$$\text{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\dots$$

$$\text{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

```

```

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

```

```

return res

```

```

Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.
open  $I_{\text{rep}}$ .
update (c:=ops+[op]).
get_lower_bound c.
close  $I_{\text{rep}}$ .
done.
Qed.

```

$$\mathbf{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \mathbf{committed} \stackrel{\text{list}}{\mapsto} ops + [op]$$

$$\mathbf{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\dots$$

$$\mathbf{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res

```

```

Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.
open  $I_{\text{rep}}$ .
update (c:=ops+[op]).
get_lower_bound c.

```

```

close  $I_{\text{rep}}$ .
done.
Qed.

```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \text{committed} \stackrel{\text{list}}{\mapsto} ops + [op]$$

$$\text{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \text{committed} \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\dots$$

$$\text{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res

```

```

Proof code.
mutex_lock.
...
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.
open  $I_{\text{rep}}$ .
update (c:=ops+[op]).
get_lower_bound c.

```

```

close  $I_{\text{rep}}$ .
done.
Qed.

```

$$\text{accepted}_0[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \text{committed} \stackrel{\text{list}}{\mapsto} ops + [op]$$

$$\text{accepted}_1[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op] \quad \text{committed} \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

$$\dots$$

$$\text{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res

```

Proof code.

mutex_lock.

....

update (a:=ops+[op]).

get_lower_bound a.

mutex_unlock with a.

applyBackup_spec.

applyBackup_spec.

...

applyBackup_spec.

open I_{rep} .

update (c:=ops+[op]).

get_lower_bound c.

close I_{rep} .

done.

Qed.

committed $\stackrel{\text{list}}{\sqsubseteq} ops + [op]$

```
s.mutex.Lock()
```

```
nextIndex := s.nextIndex
```

```
e := s.epoch
```

```
s.nextIndex += 1
```

```
res := s.stateLogger.LocalApply(op)
```

```
s.mutex.Unlock()
```

```
for j := 0; j < len(s.backupClerks); j++ {
```

```
    s.backupClerks[j].ApplyAsBackupRPC(e,
```

```
nextIndex, op)
```

```
}
```

```
return res
```

```

Proof code.
mutex_lock.
....
update (a:=ops+[op]).
get_lower_bound a.
mutex_unlock with a.
applyBackup_spec.
applyBackup_spec.
...
applyBackup_spec.
open  $I_{rep}$ .
update (c:=ops+[op]).
get_lower_bound c.
close  $I_{rep}$ .
done.
Qed.

```

committed $\stackrel{\text{list}}{\sqsubseteq} ops + [op]$

```

s.mutex.Lock()
nextIndex := s.nextIndex
e := s.epoch
s.nextIndex += 1
res := s.stateLogger.LocalApply(op)
s.mutex.Unlock()

for j := 0; j < len(s.backupClerks); j++ {
    s.backupClerks[j].ApplyAsBackupRPC(e,
nextIndex, op)
}

return res

```

Time-bounded invariants

CurrentEpoch $\mapsto e$ L

L ^{expires}
 $\geq t$

\implies **CurrentEpoch** $\mapsto e$

GetTimeRange().latest $< t$

Time-bounded invariants

CurrentEpoch $\mapsto e$

accepted₀[e] $\stackrel{\text{list}}{\mapsto} ops$

\implies **committed** $\stackrel{\text{list}}{\mapsto} ops'$
 $ops' \preceq ops$

I_{rep}

$\exists l, e, \text{committed} \stackrel{\text{list}}{\mapsto} l*$

accepted₀[e] $\stackrel{\text{list}}{\sqsupseteq} l * \dots * \text{accepted}_n[e] \stackrel{\text{list}}{\sqsupseteq} l$

What about "what if" questions?

Let's look at a few "what if" scenarios, and see how the proof handles them

What if backup loses operations?

Would be buggy:

`ApplyAsBackupRPC` promises $\text{accepted}_j[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$

Backup can't roll back $\text{accepted}_j[e] \stackrel{\text{list}}{\mapsto} ops + [op]$

What if old primary does a Put?

Would be buggy:

At least one replica promises $\text{accepted}_j[e] \stackrel{\text{list}}{\mapsto} \square ops$

New op requires $\text{accepted}_j[e] \stackrel{\text{list}}{\sqsupseteq} ops + [op]$

Replica can't modify $\text{accepted}_j[e] \stackrel{\text{list}}{\mapsto} \square ops$

What if lease expires during Get?

Not buggy:

Get access to `CurrentEpoch` $\mapsto e$ and `committed` $\overset{\text{list}}{\mapsto} ops'$
at the moment of `GetTimeRange()`

Benefits of verifying GroveKV

Eliminate the need for tests?

Not quite... GroveKV's spec rules out (some) safety bugs

Liveness + performance bugs still possible

(When) are proofs worth it?

Component	Lines of Code	Lines of Proof
Network library	120	Trusted
Filesystem library	50	Trusted
RPC library	161	1,311
Replica server	574	7,986
Reconfiguration	65	803
Configuration server	200	2,048
Clerk	156	878
Append-only file	90	1,345
State logger	134	1,732
Exactly-once operations	128	2,186
Key-value mappings	157	916
Time bounded invariants	–	142
Total	1,835	–
Total verified	1,665	19,347

Figure 5: Lines of Go code and Coq proof for GroveKV.