



# **Argosy: Verifying layered storage systems with recovery refinement**

**Tej Chajed**, Joseph Tassarotti, Frans Kaashoek, Nickolai Zeldovich

MIT


logical disk



disk<sub>1</sub>

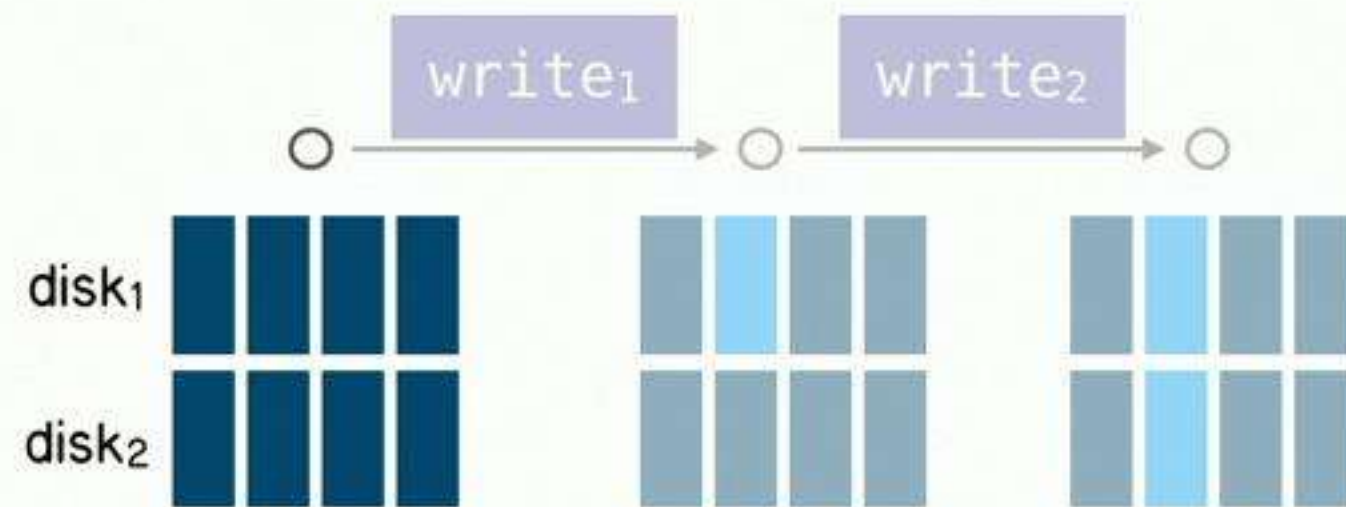


disk<sub>2</sub>



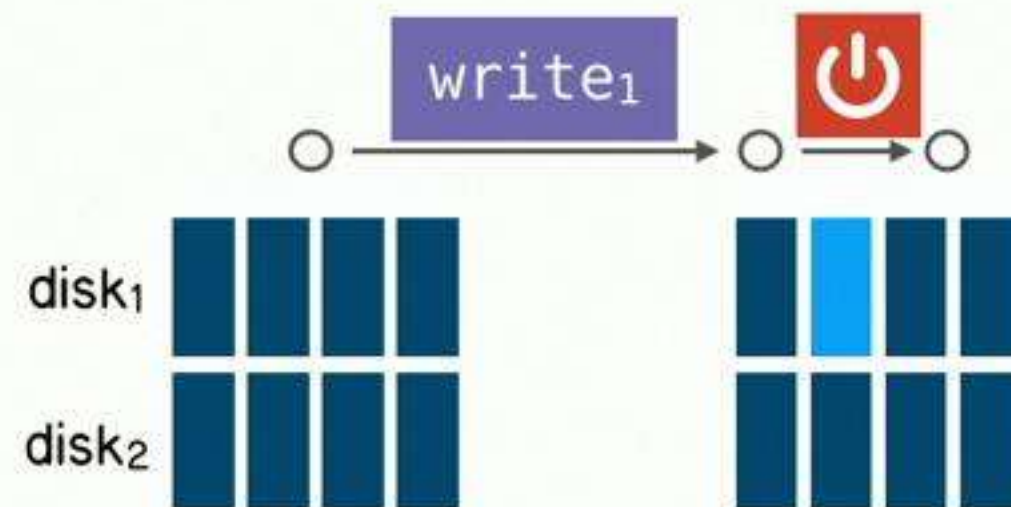
Bob writes a replication system

logical disk

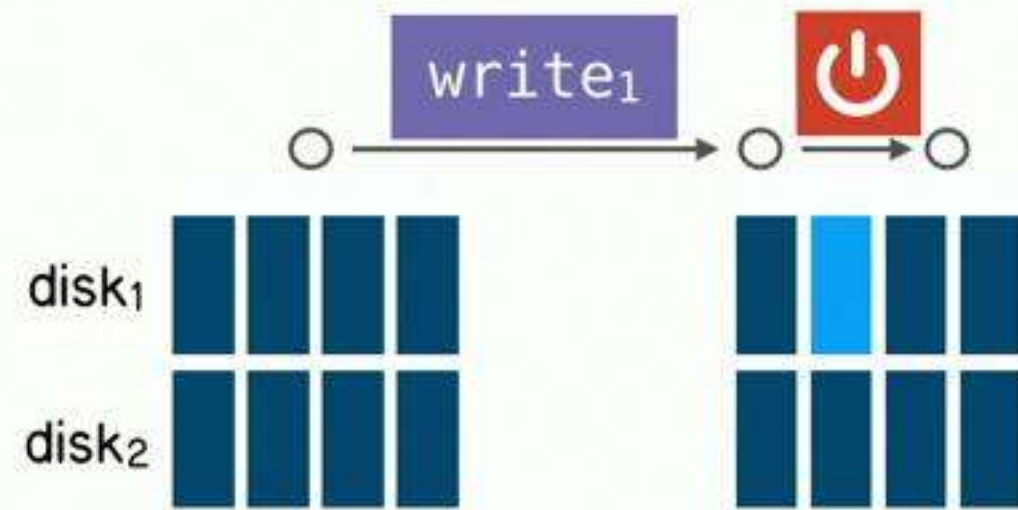


Bob writes a replication system

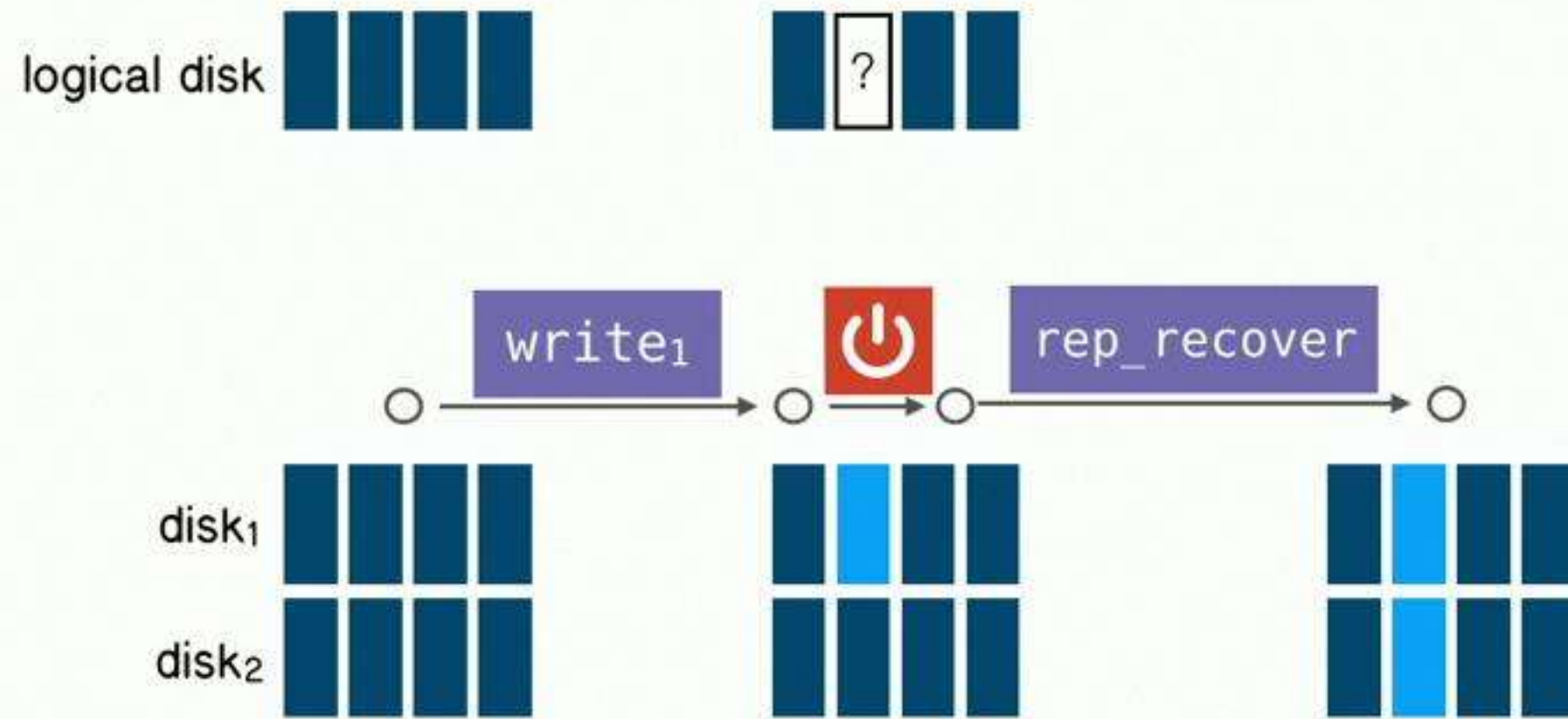
logical disk



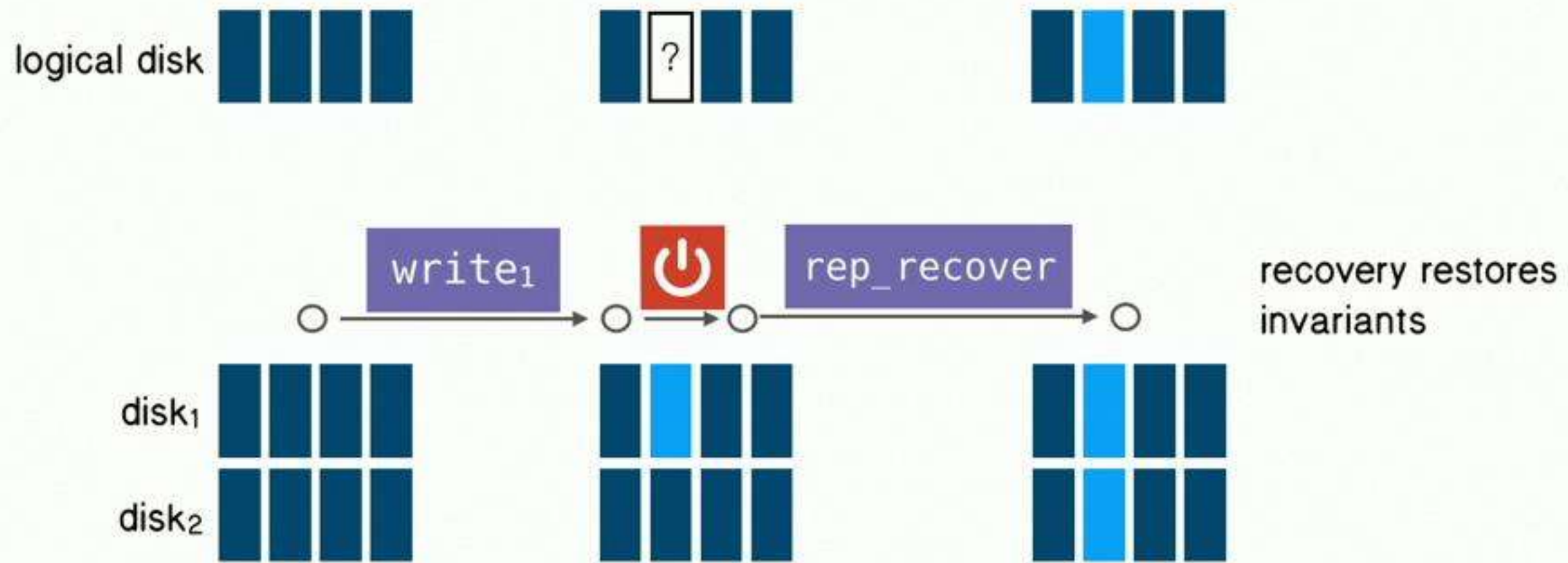
Bob writes a replication system



Bob writes a replication system

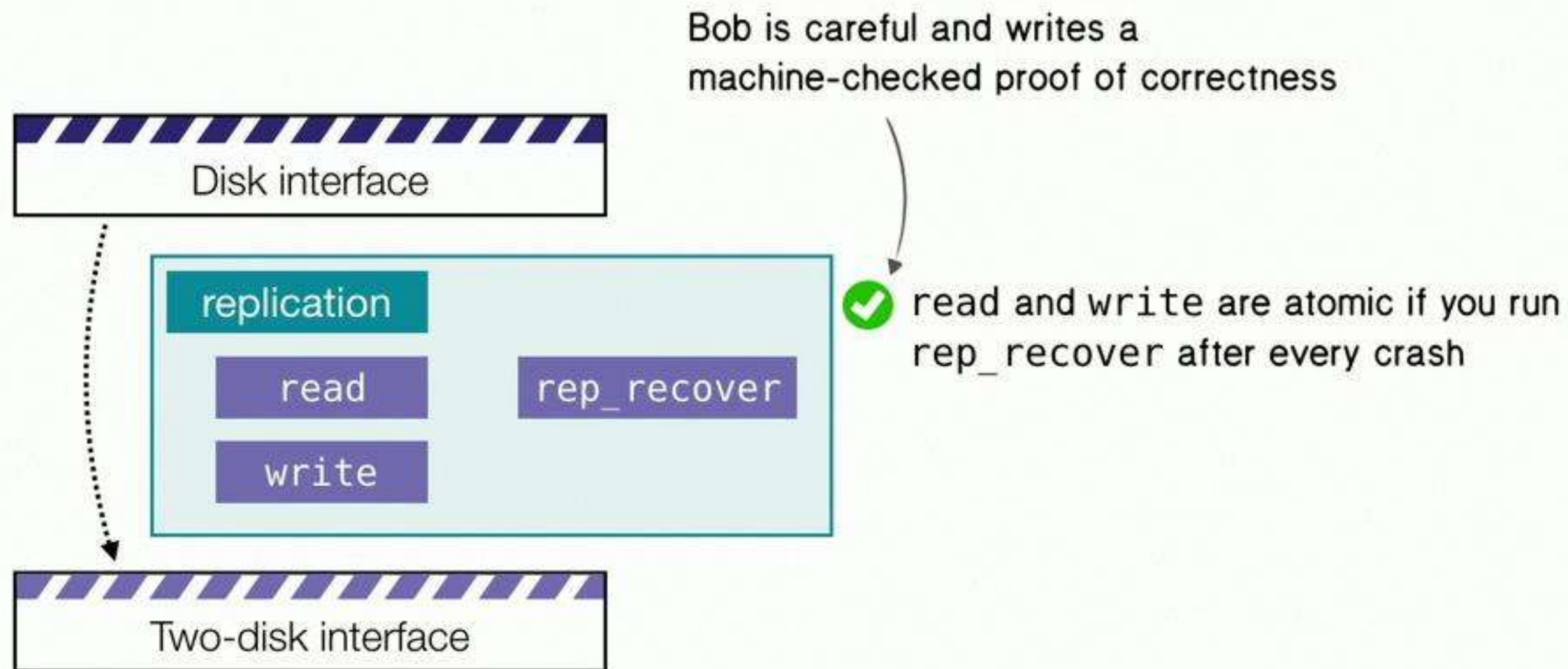


Bob writes a replication system and implements its recovery procedure

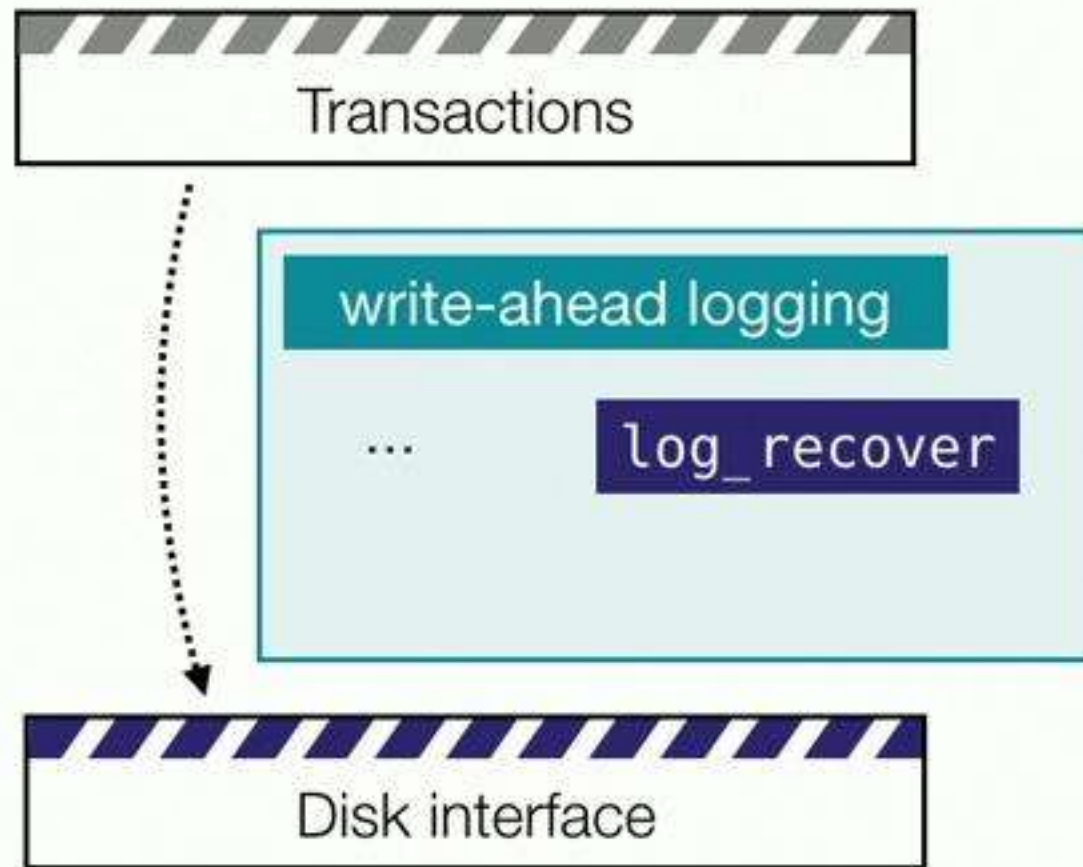


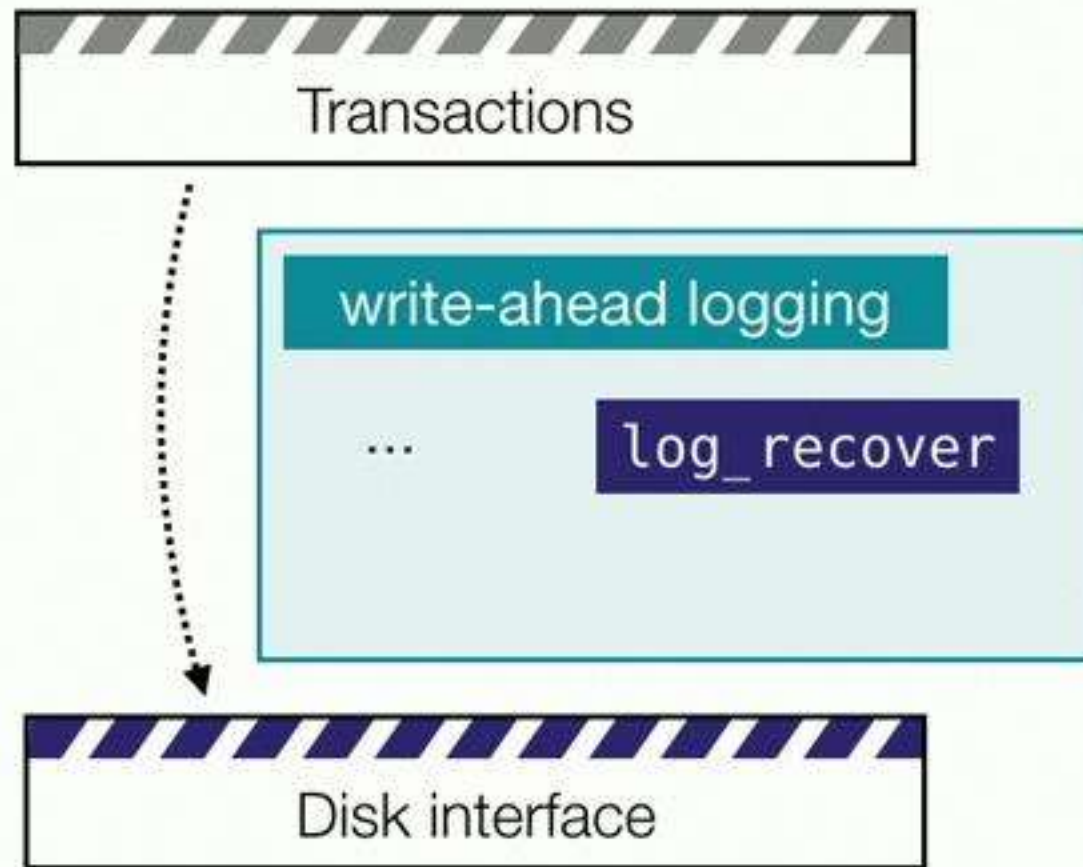
Bob writes a replication system and implements its recovery procedure



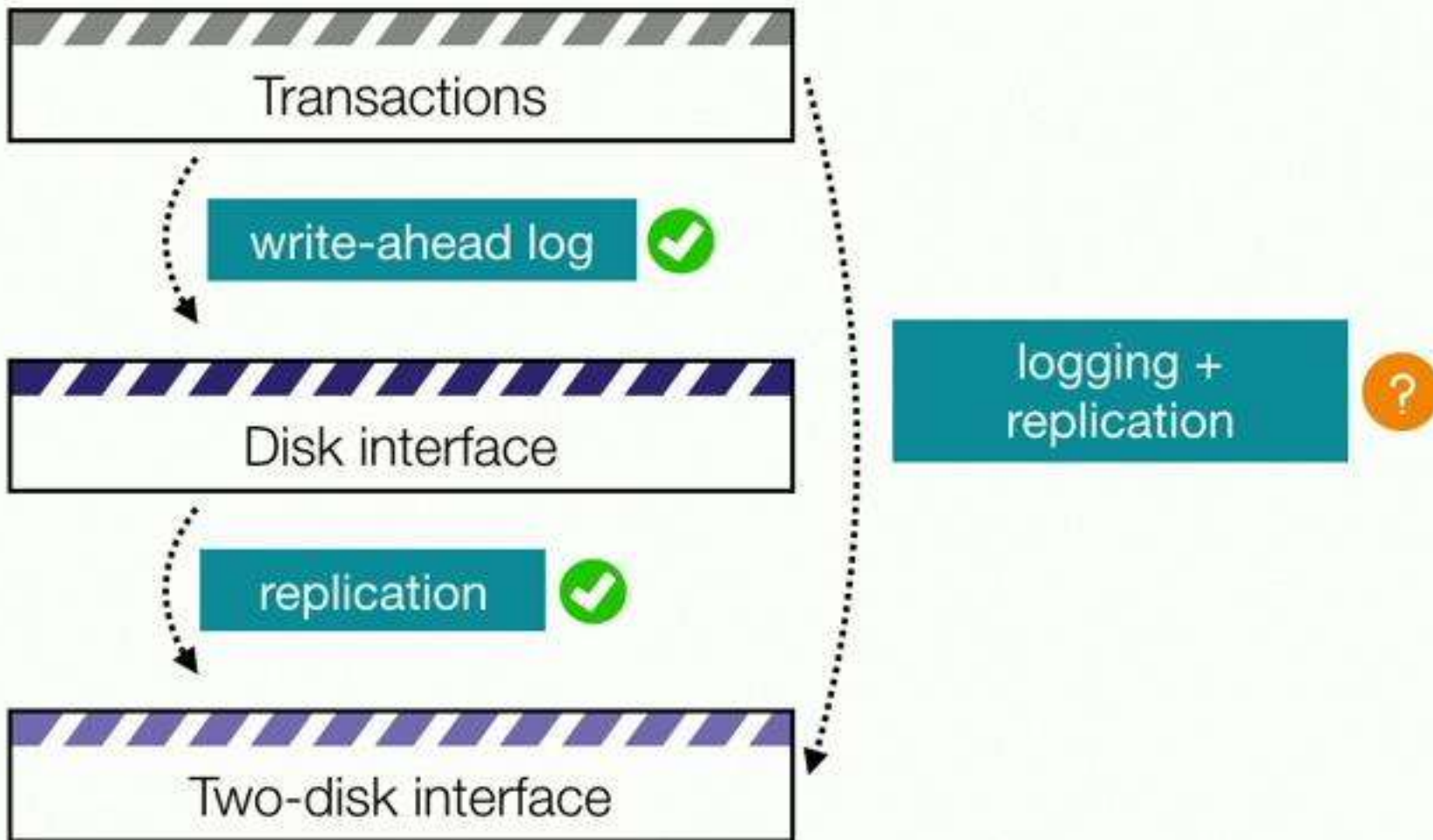


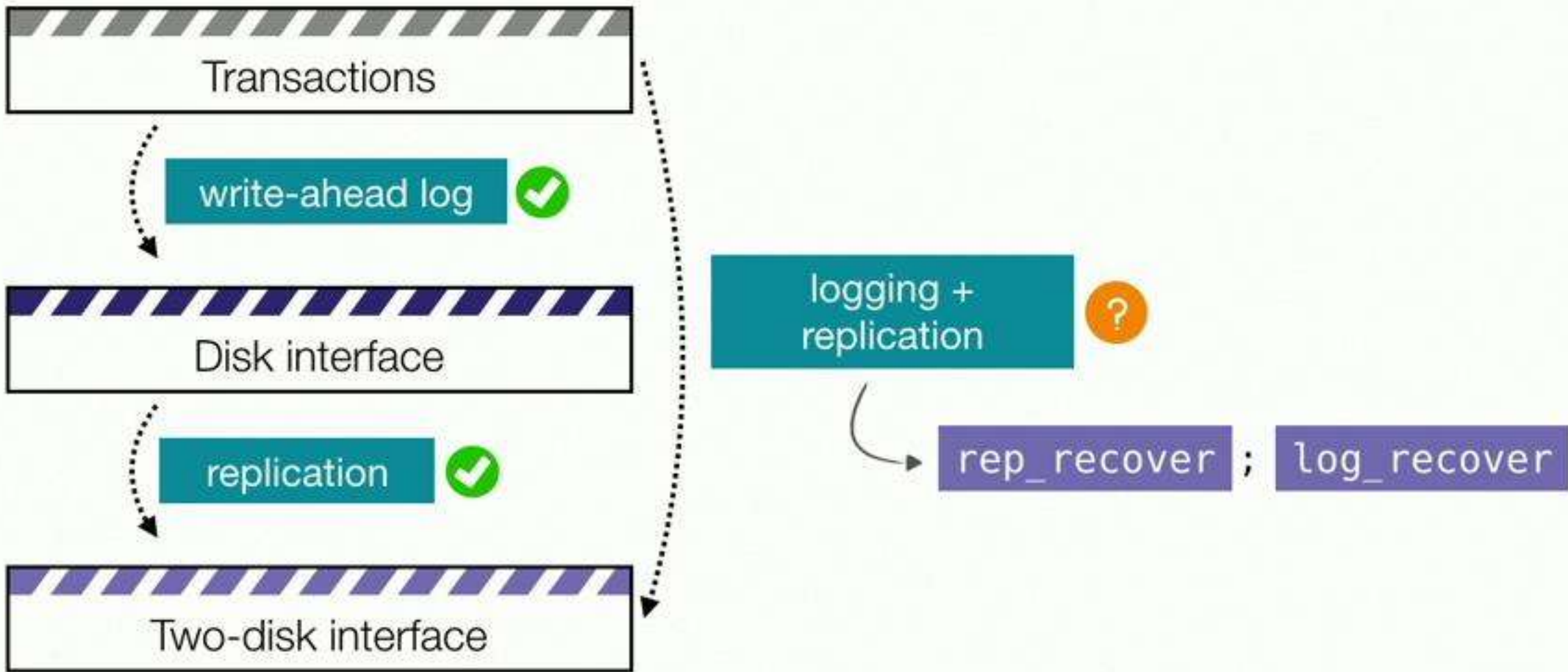






✓ ops are atomic if you run  
log\_recover after every crash





# Challenge: crashes during composed recovery

rep\_recover ✓ under crashes

log\_recover ✓ under crashes

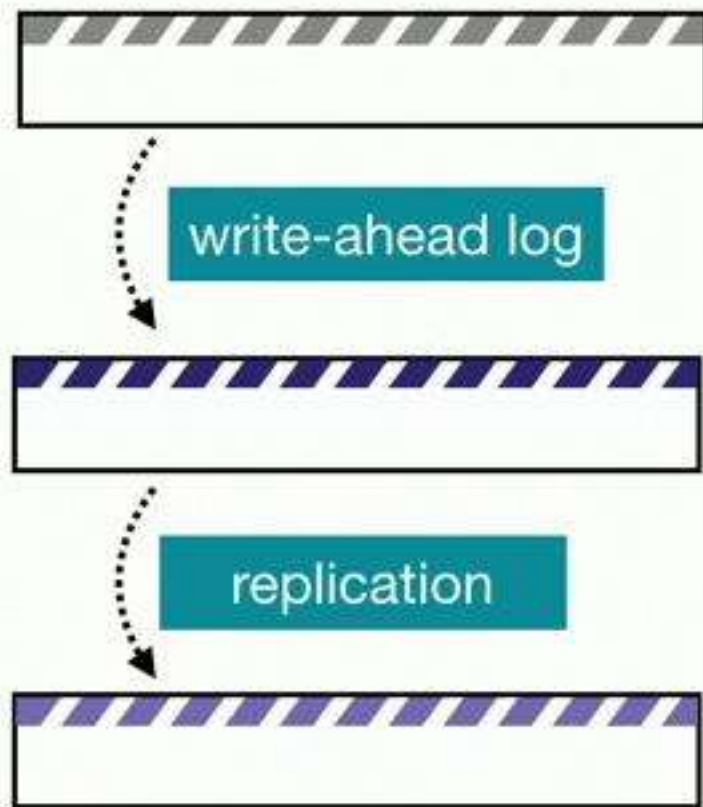
---

rep\_recover ; log\_recover ?

how do we prove correctness  
under crashes using the existing proofs?



# Prior work cannot handle multiple recovery procedures



**CHL** [SOSP '15]

not modular

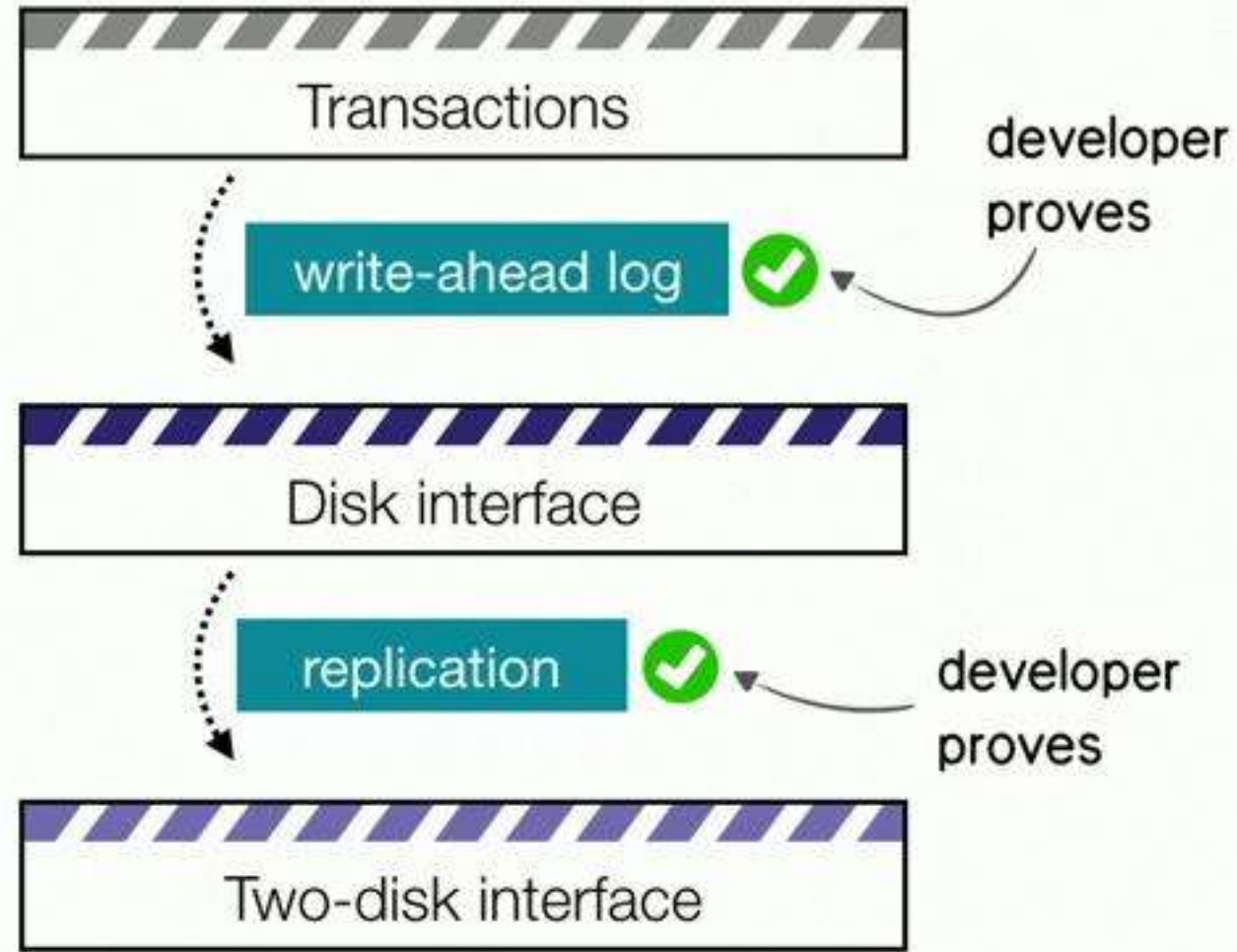
**Yggdrasil** [OSDI '16]

single recovery

**Flashix** [SCP '16]

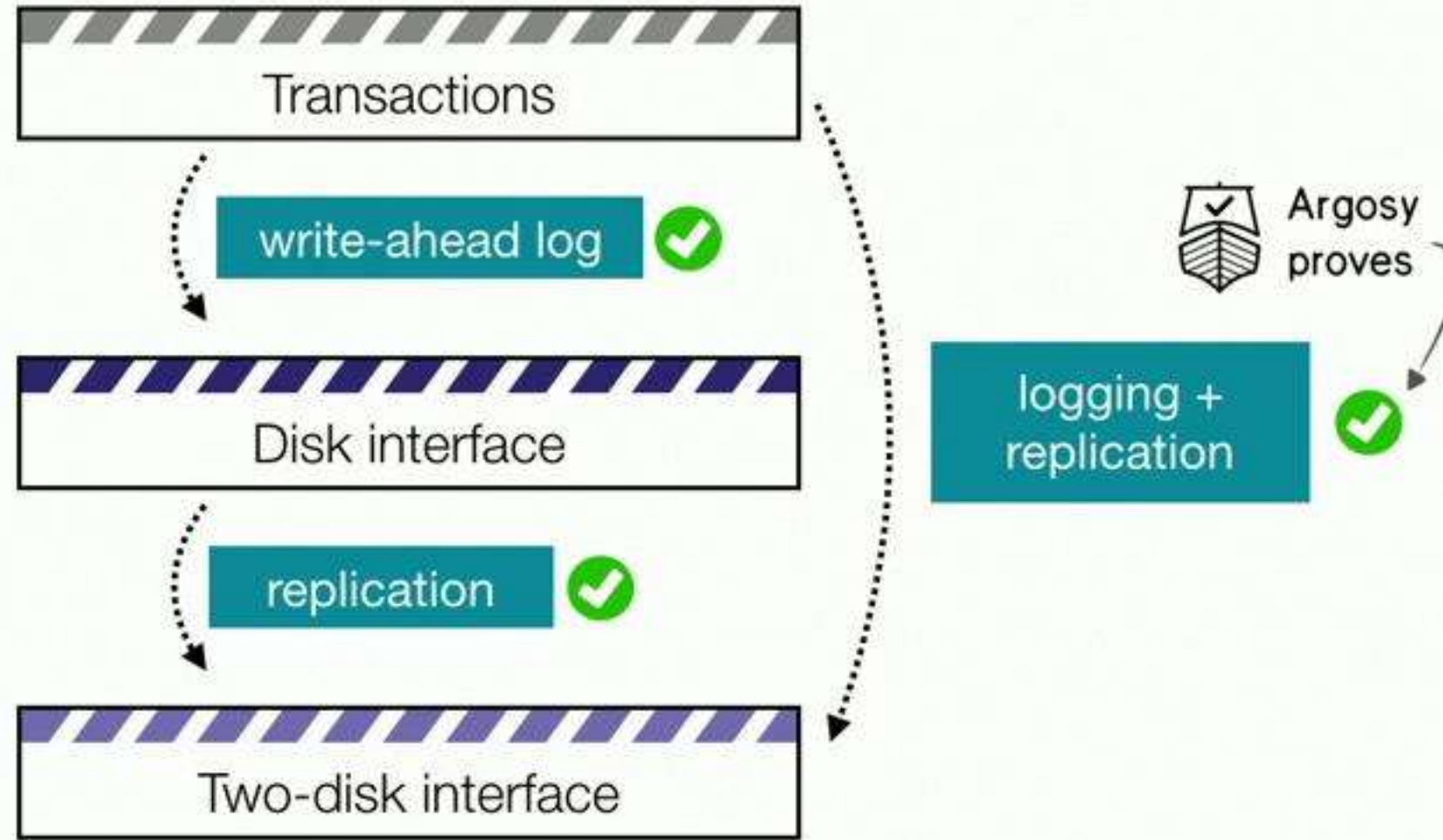
restricted recovery  
procedures

# Argosy supports modular recovery proofs

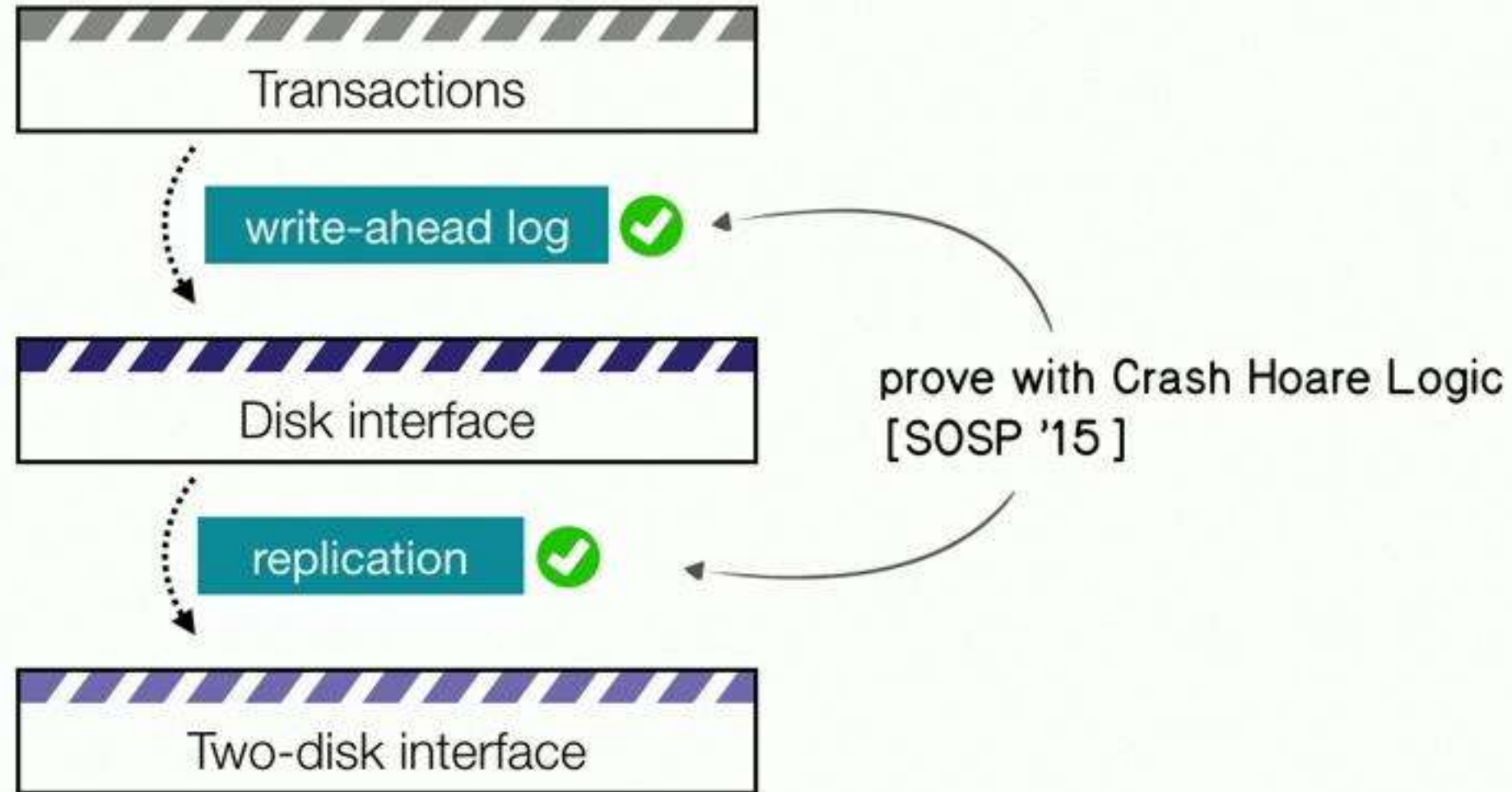




# Argosy supports modular recovery proofs



# Argosy is compatible with existing techniques



# Contributions

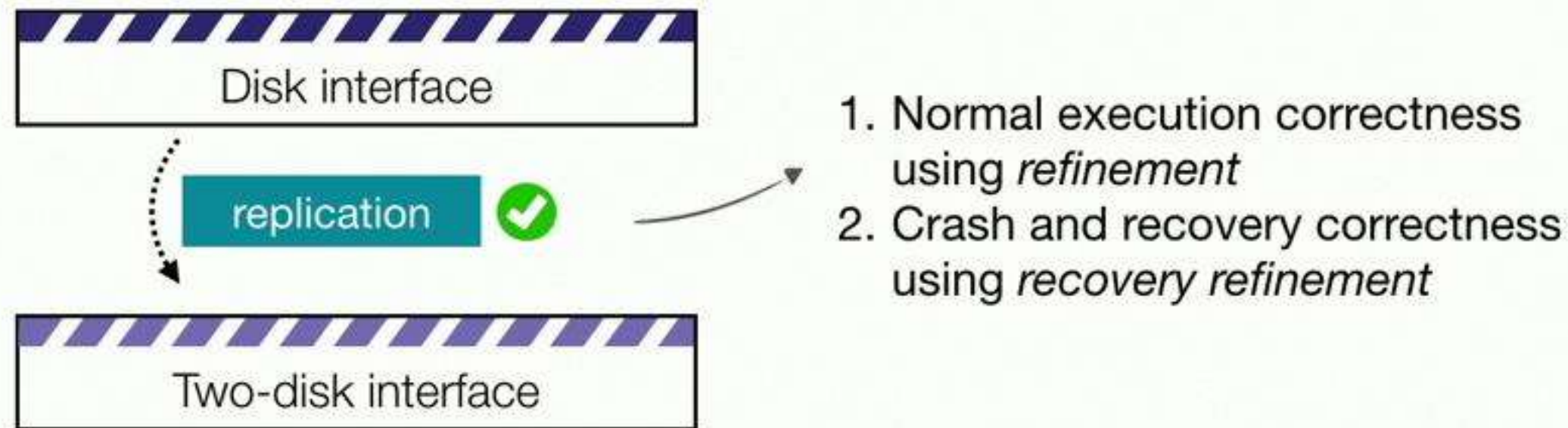
Recovery refinement for modular proofs

CHL for proving recovery refinement

see paper Verified example: logging + replication

see code Machine-checked proofs in Coq 🐔

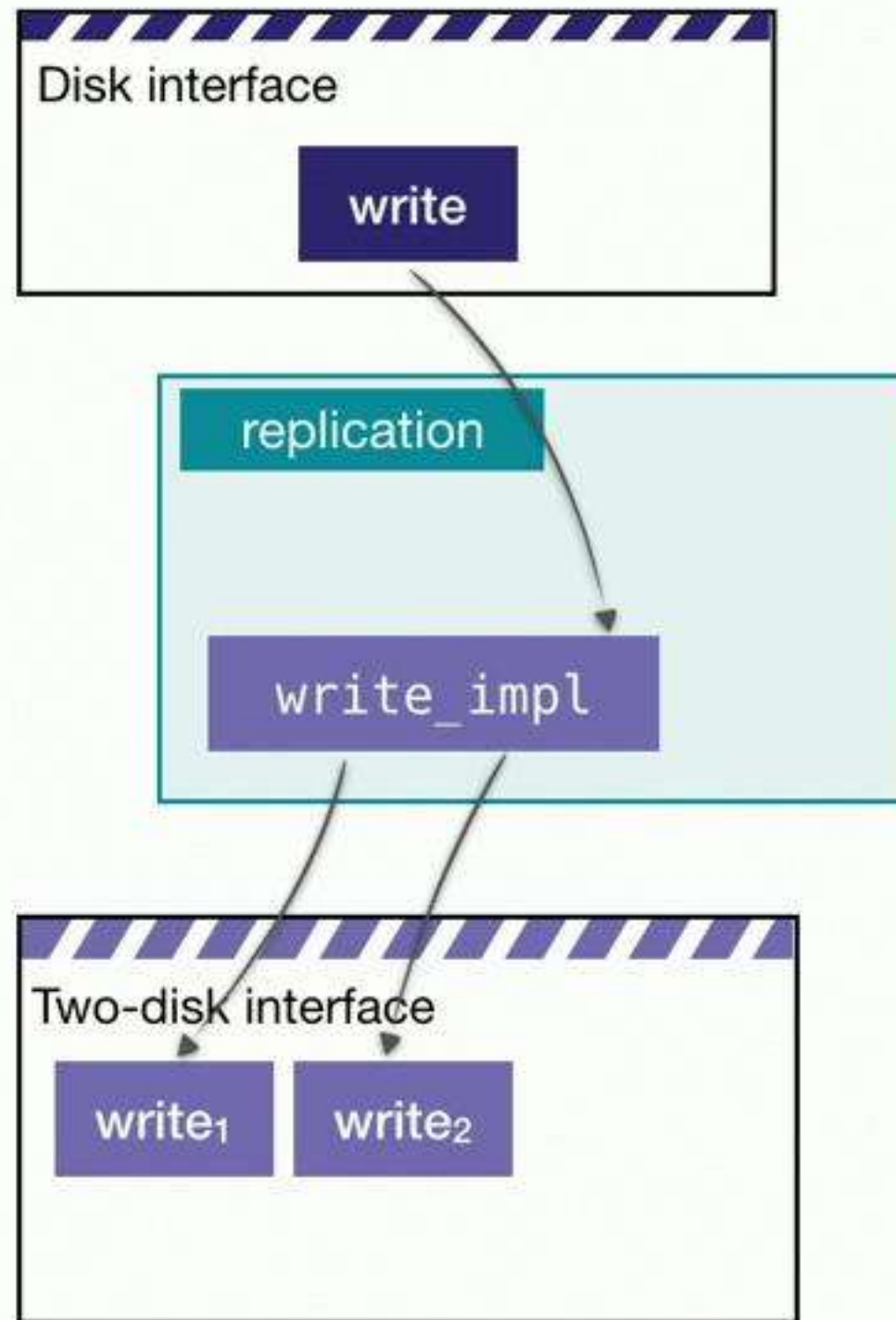
# Preview: recovery refinement



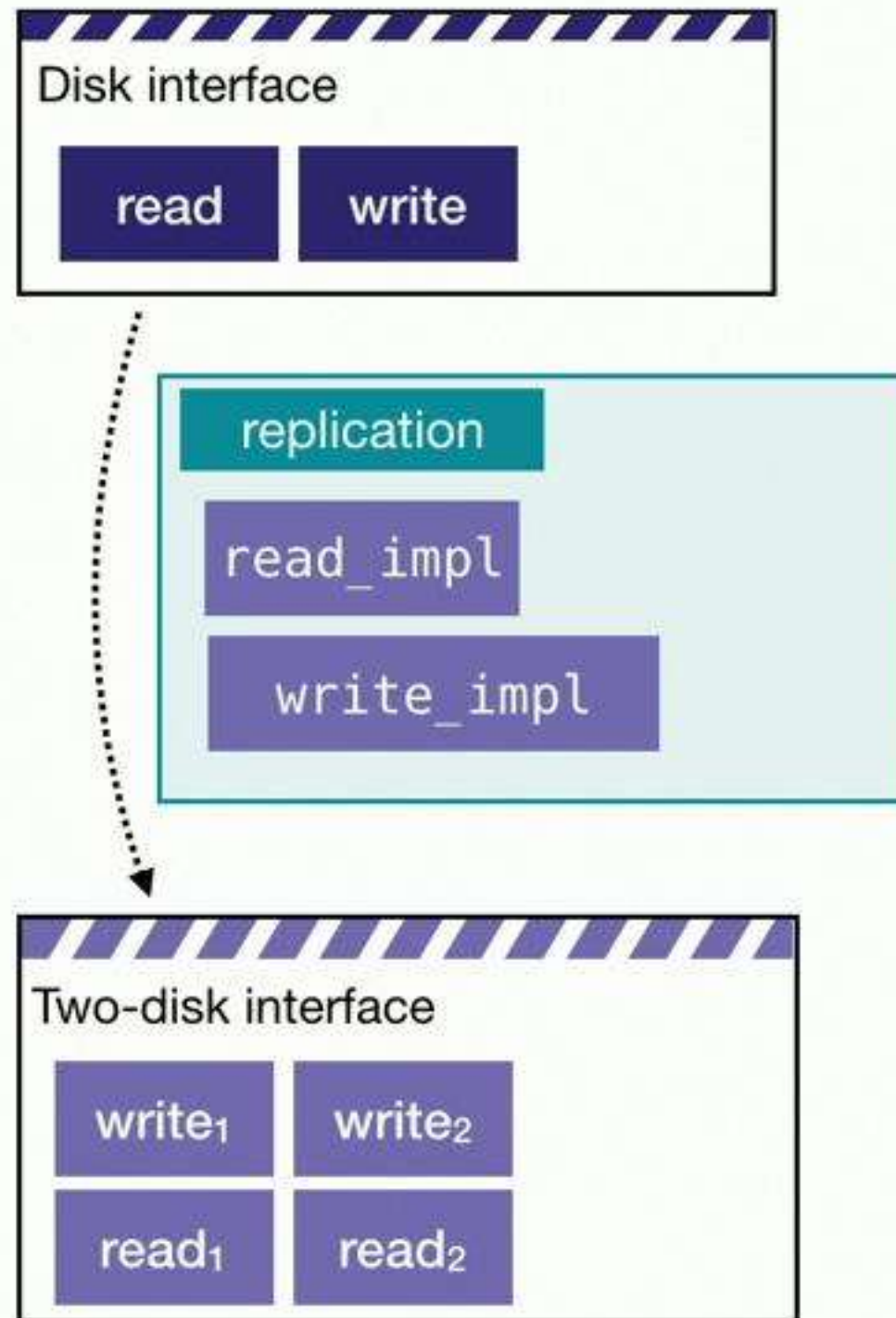


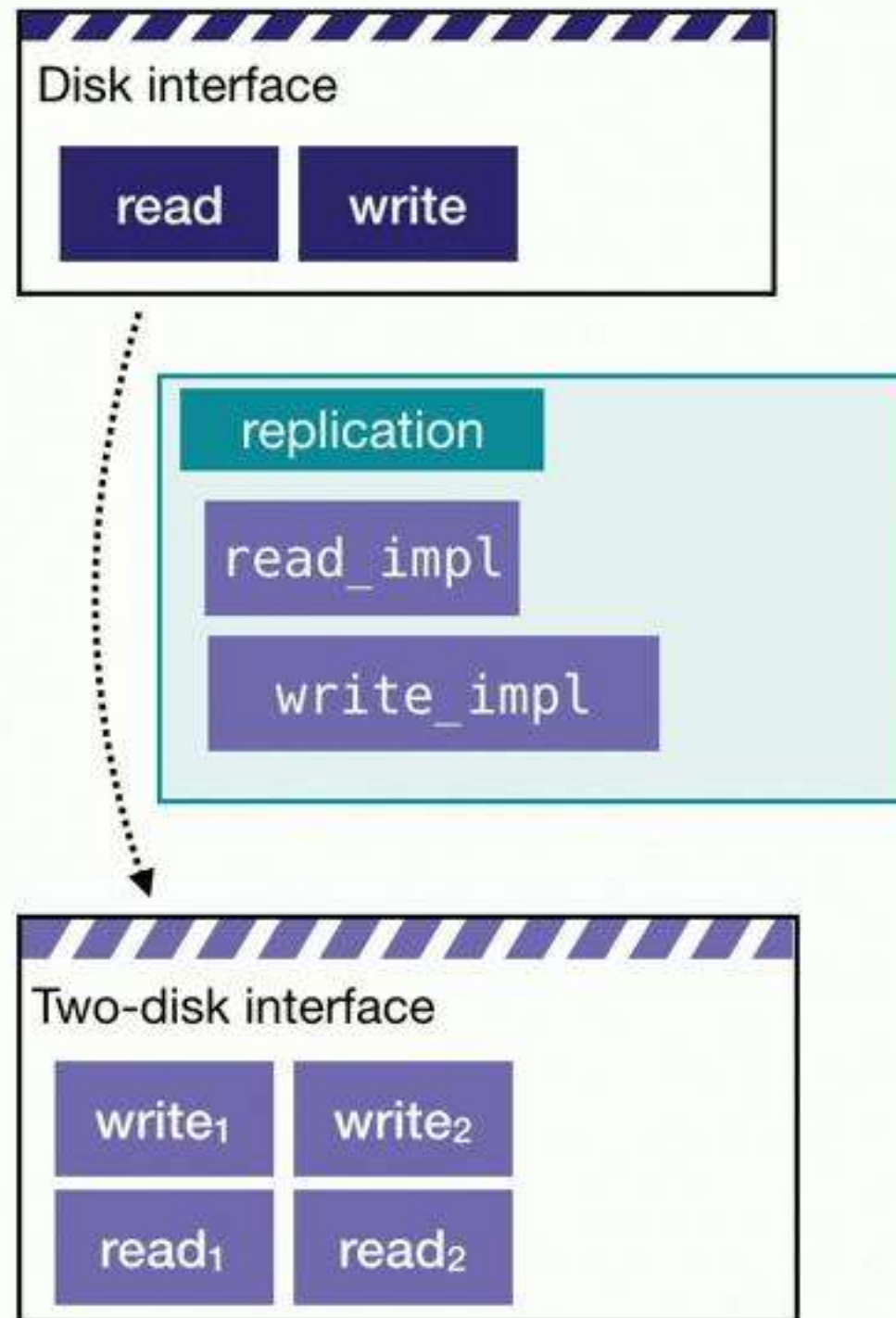
# Refinement



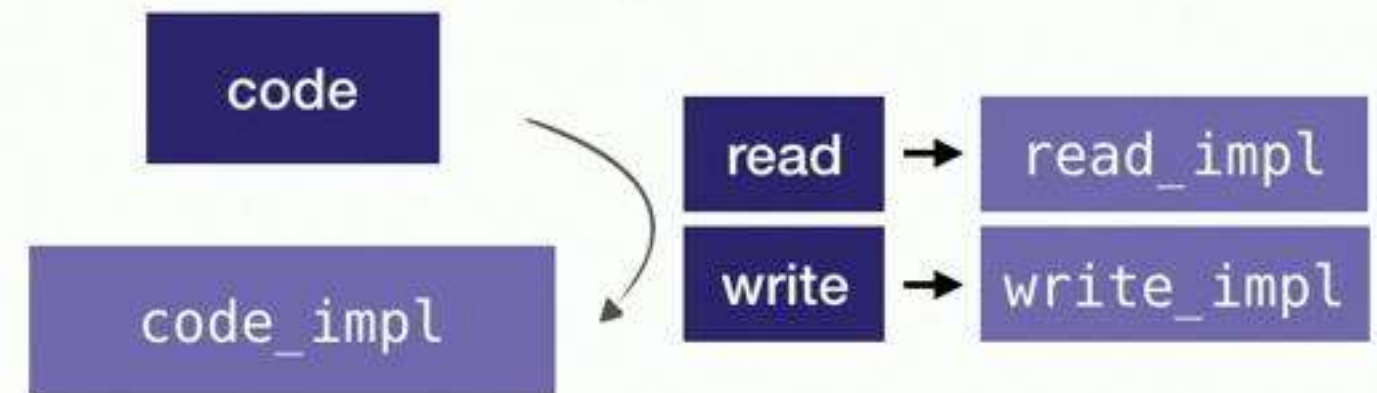




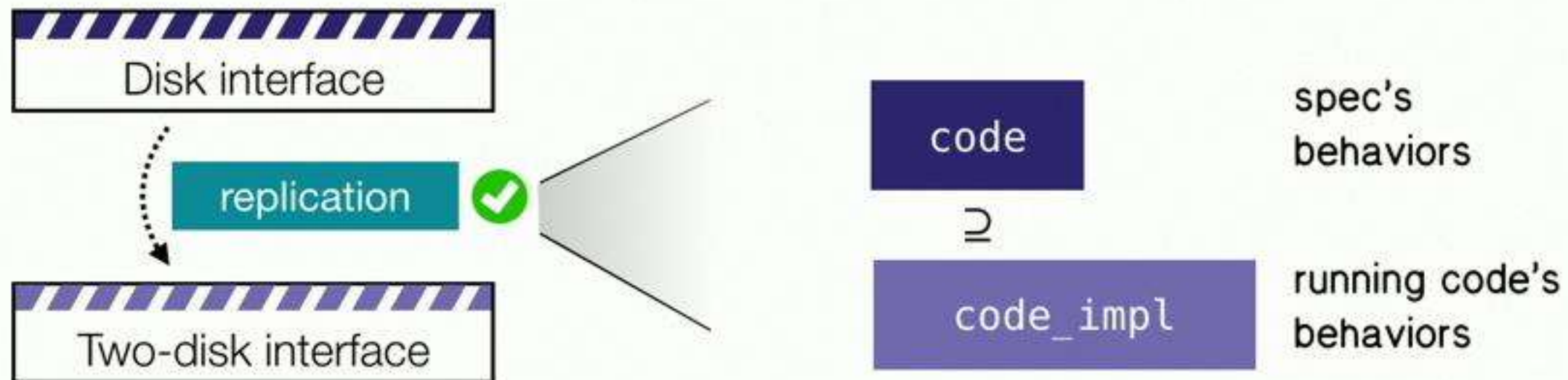




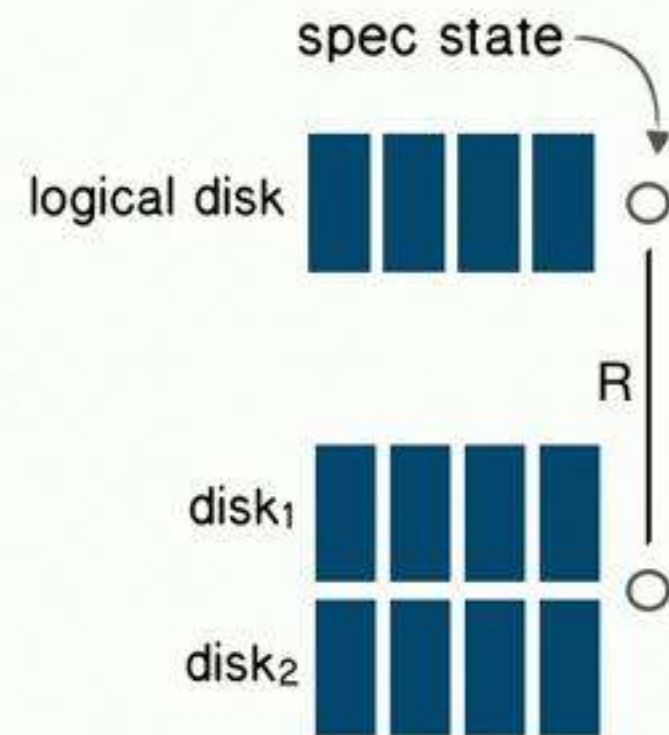
correctness is based on how we use replication:  
run code using Disk interface on top of two disks



# Correctness: trace inclusion



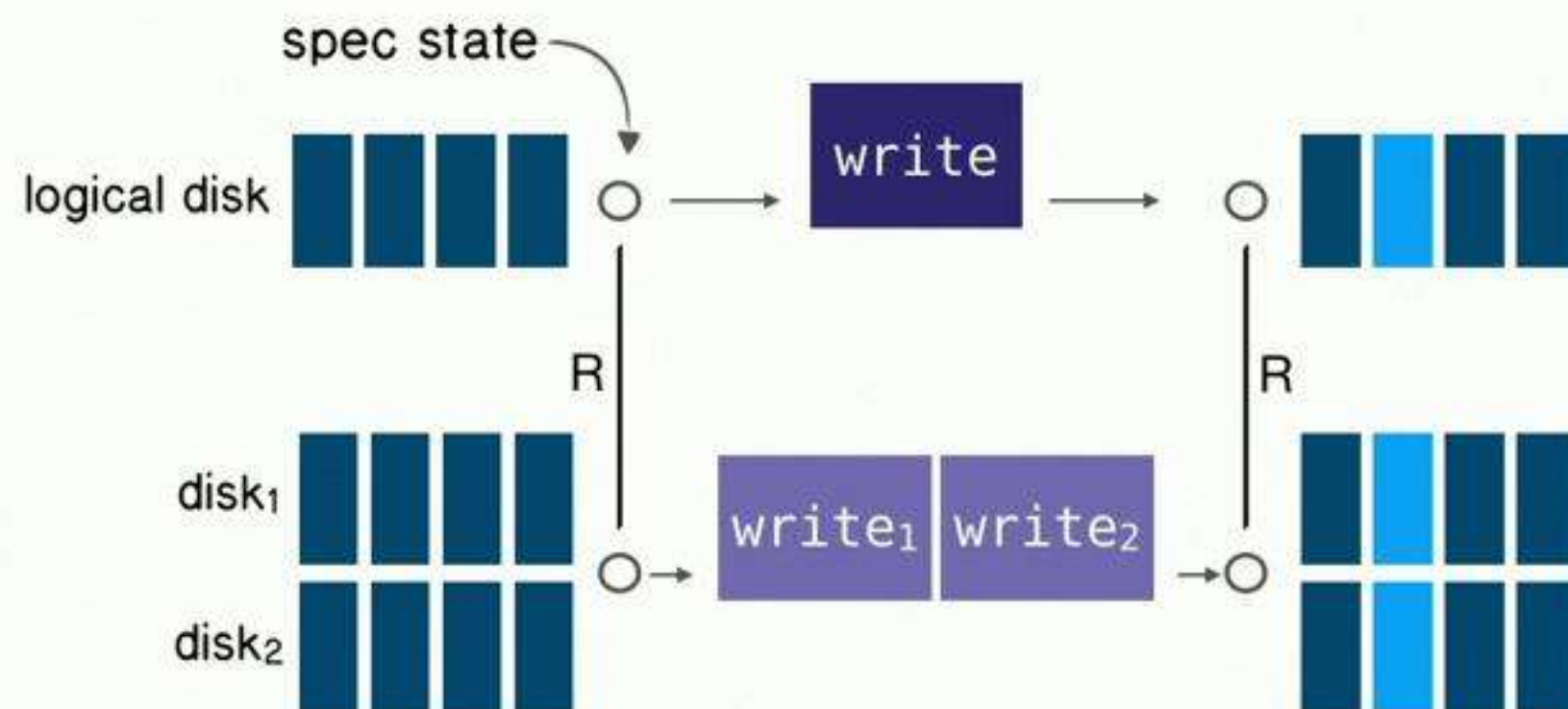
# Proving correctness with an abstraction relation



1. developer provides abstraction relation  $R$

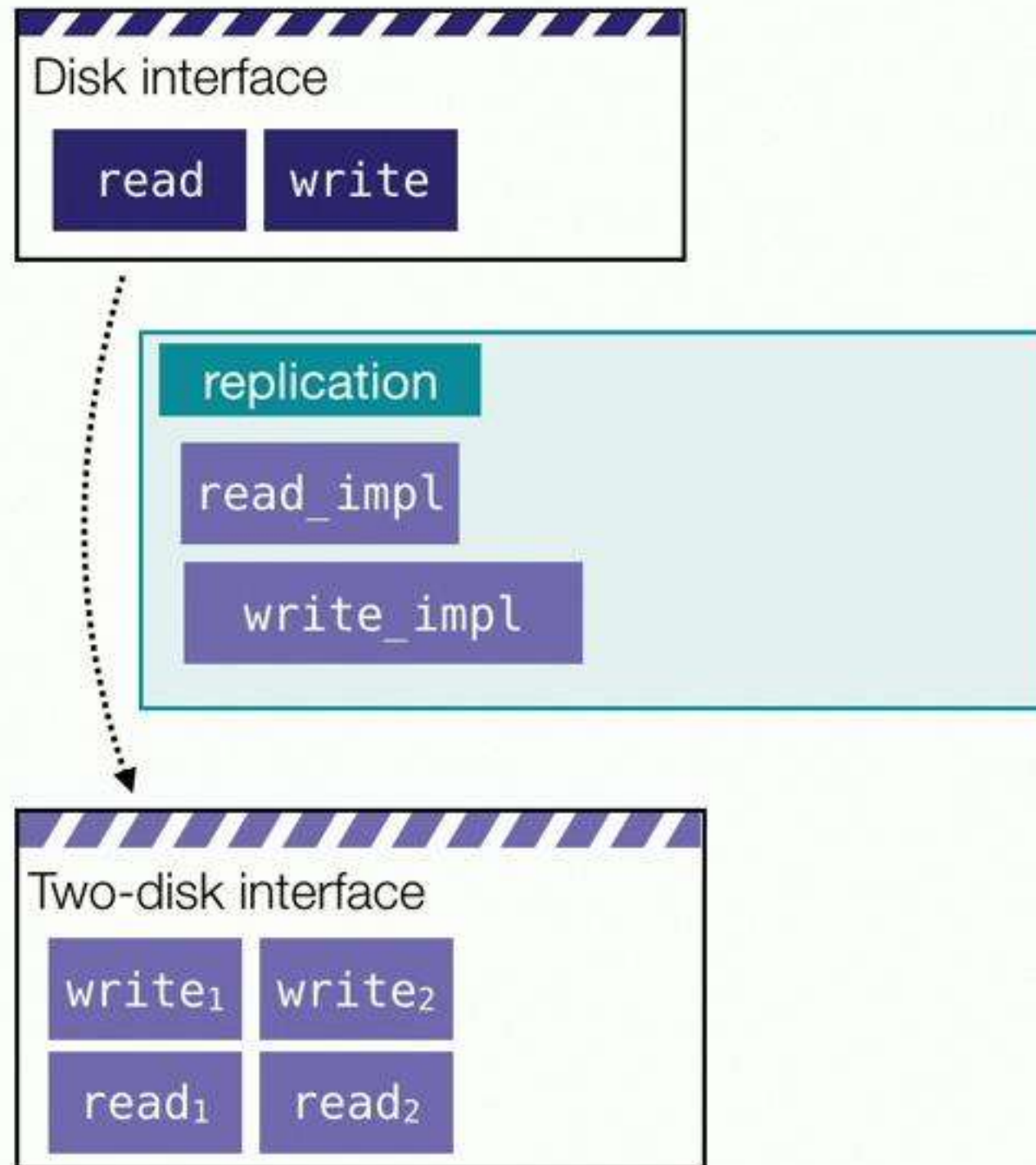


# Proving correctness with an abstraction relation

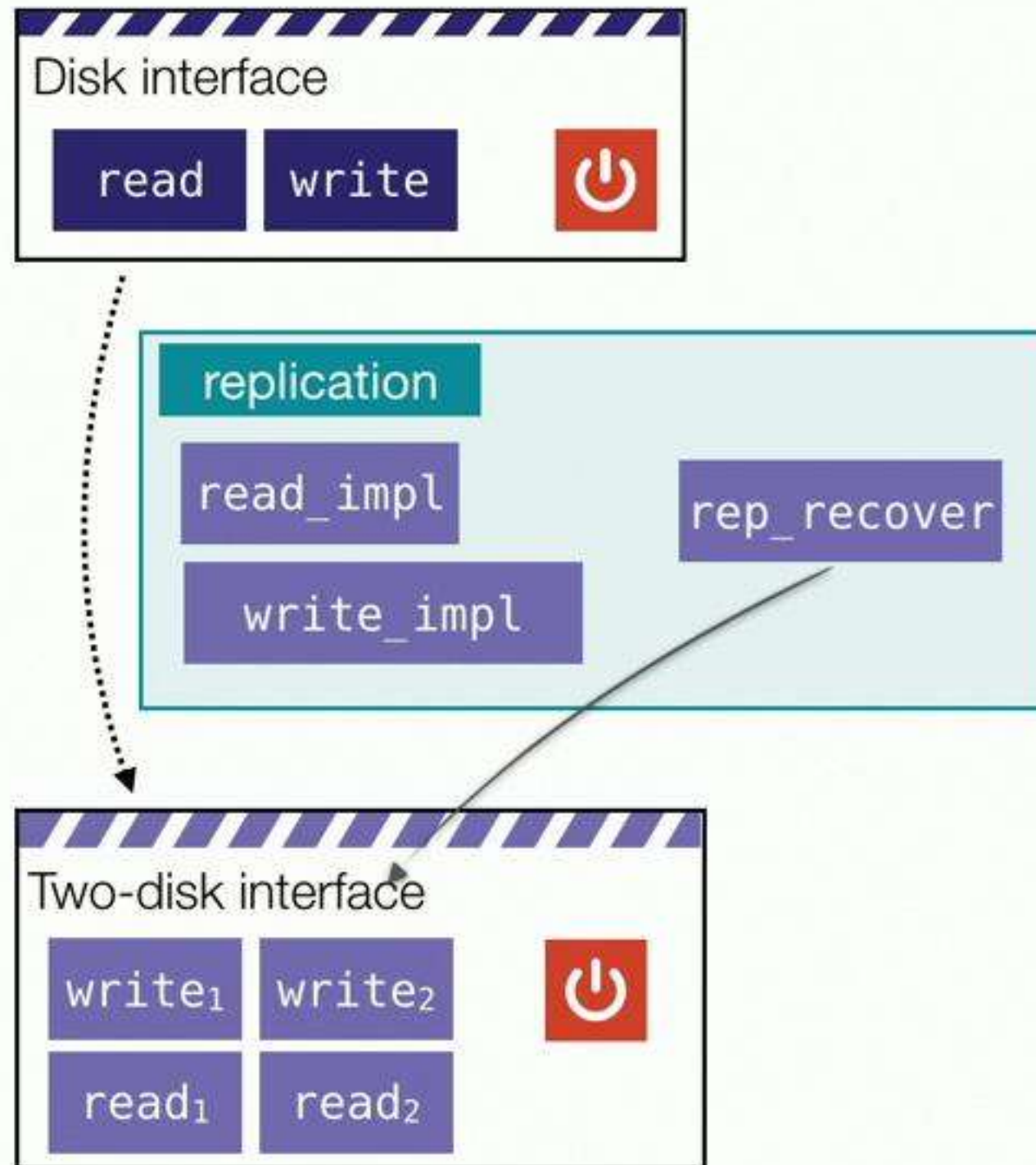


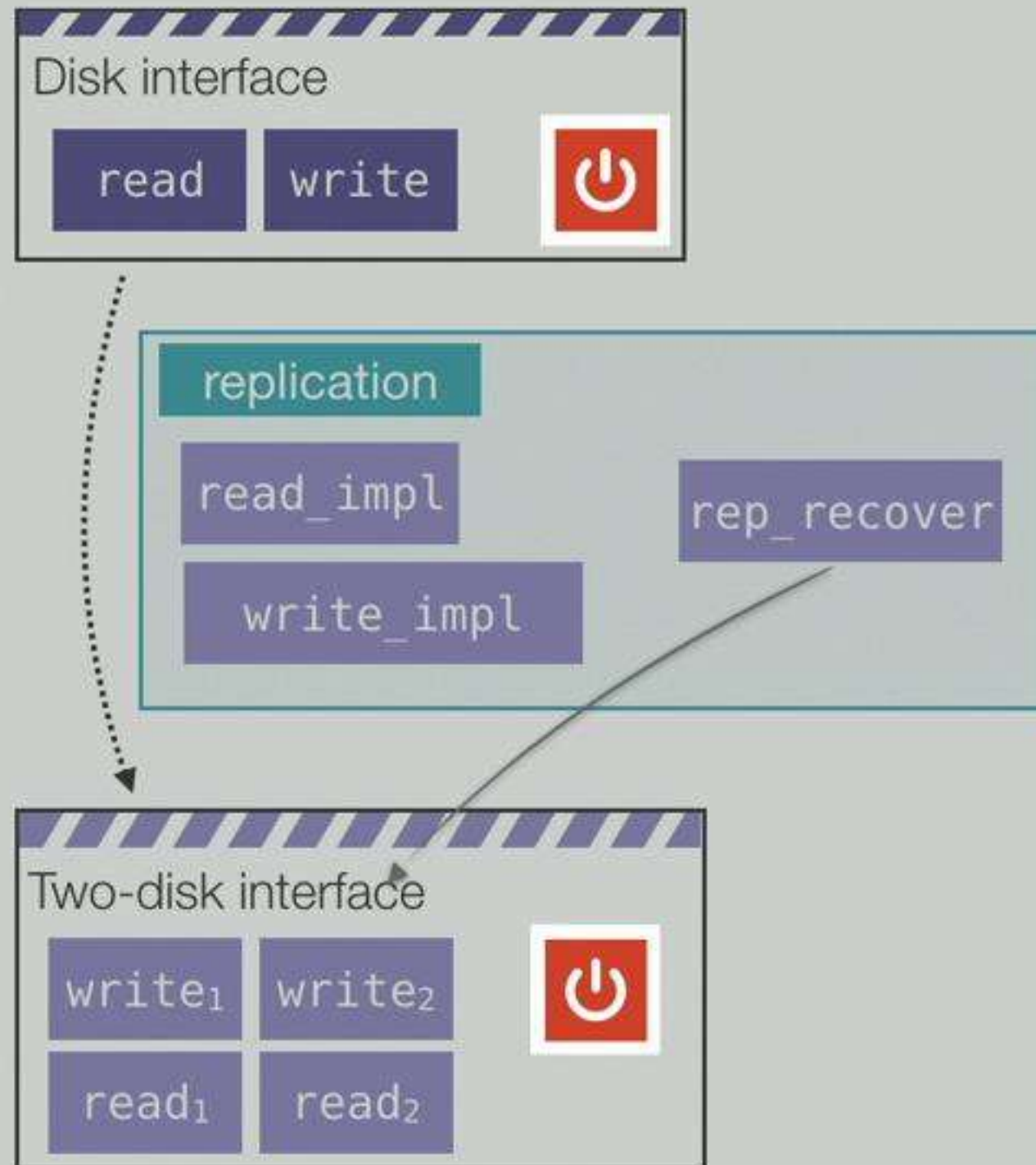
1. developer provides abstraction relation  $R$
2. prove spec execution exists
3. and abstraction relation is preserved

# Recovery refinement

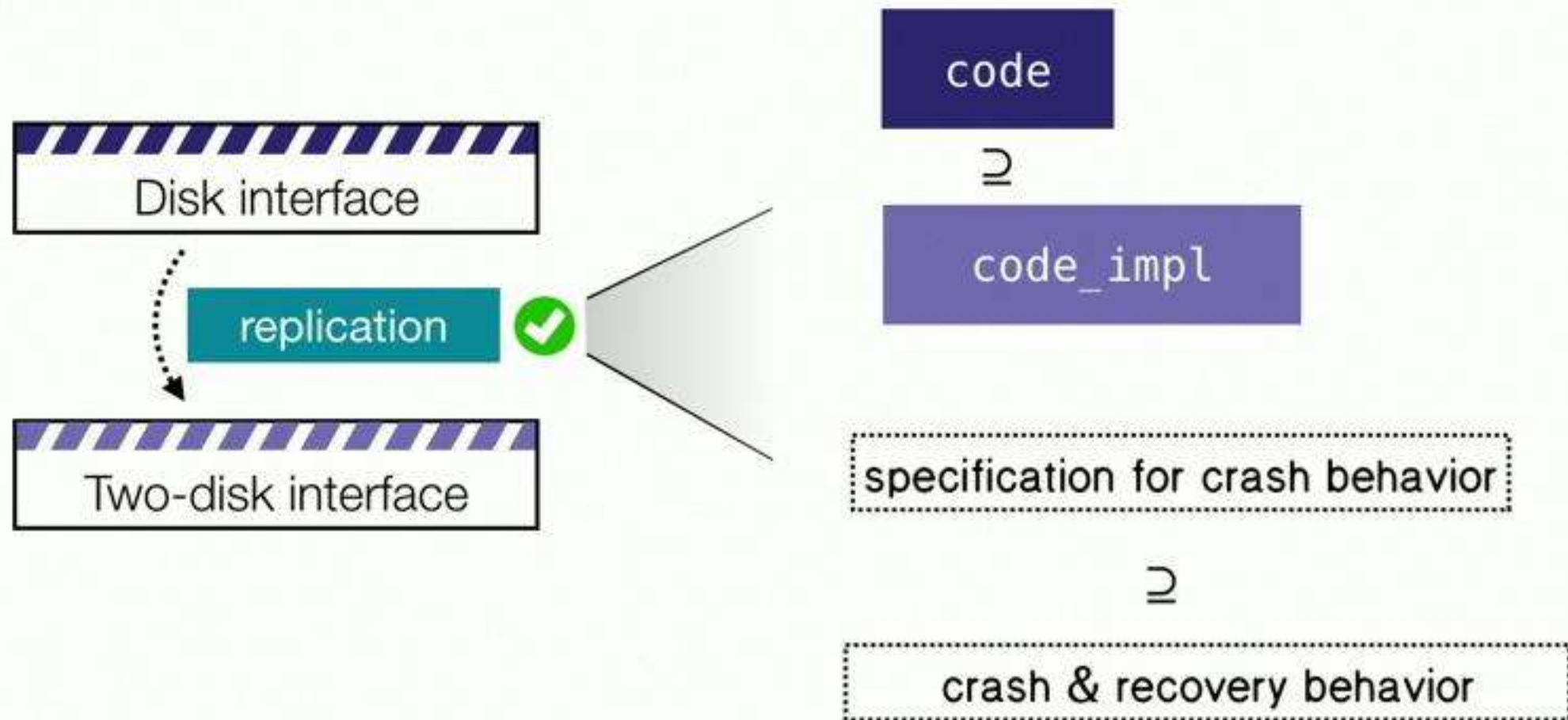




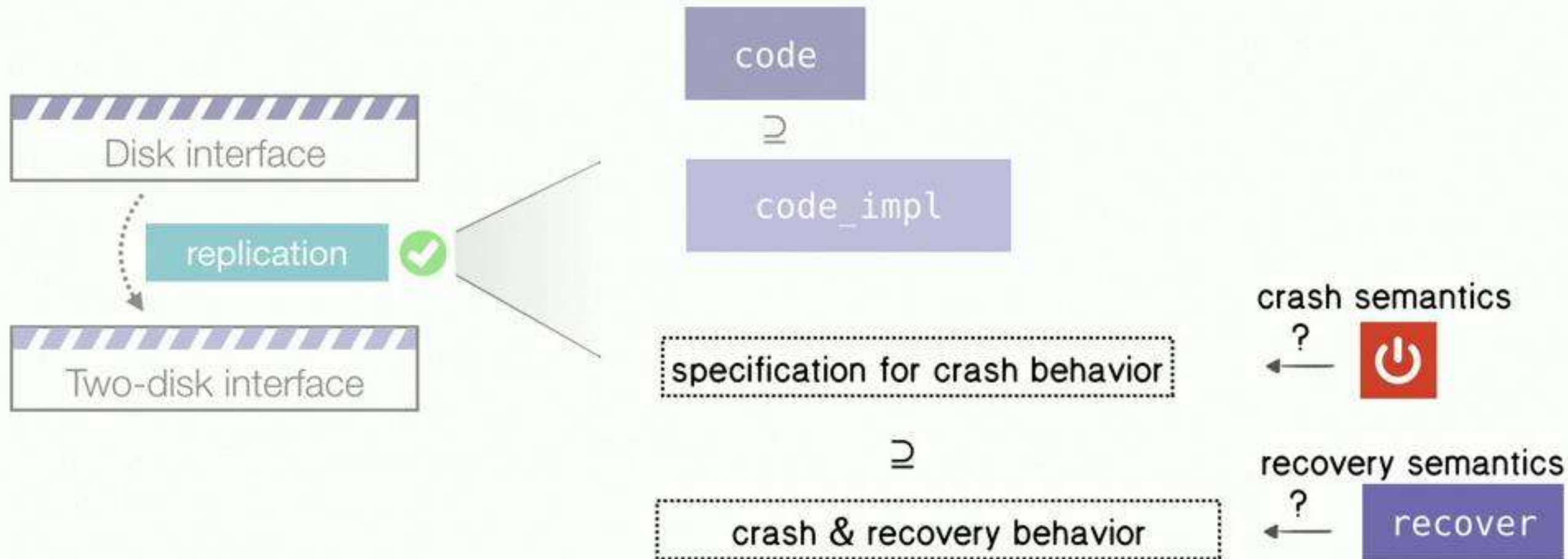




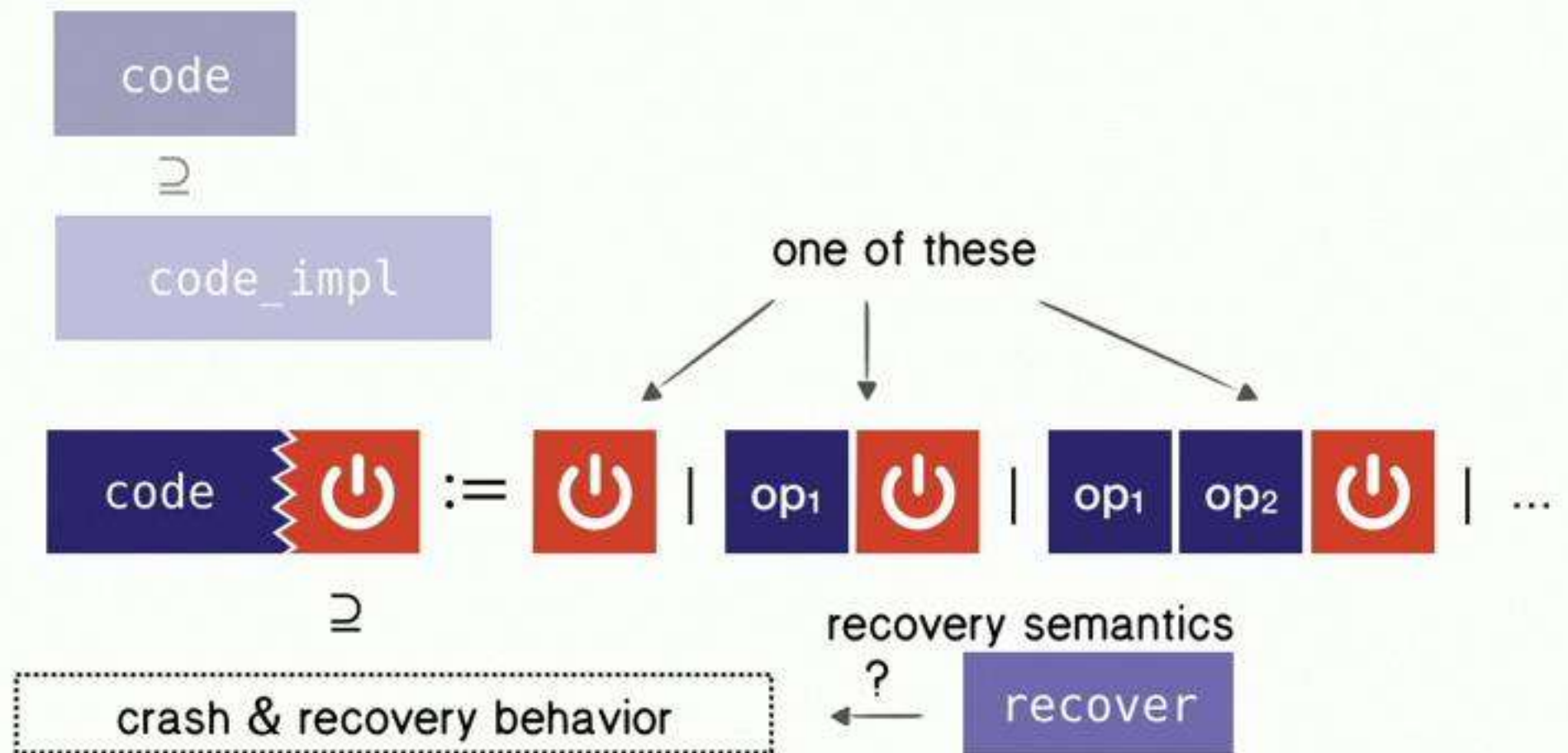
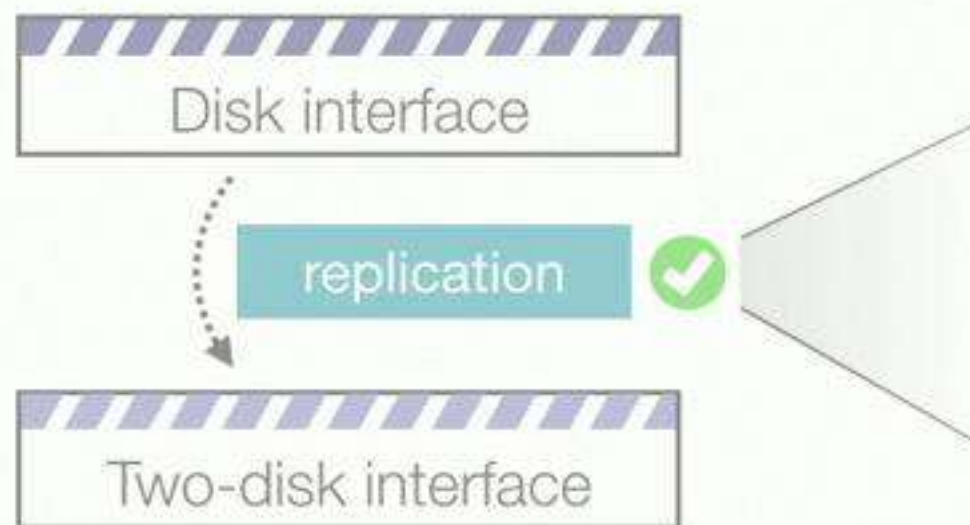
# Extending trace inclusion with recovery

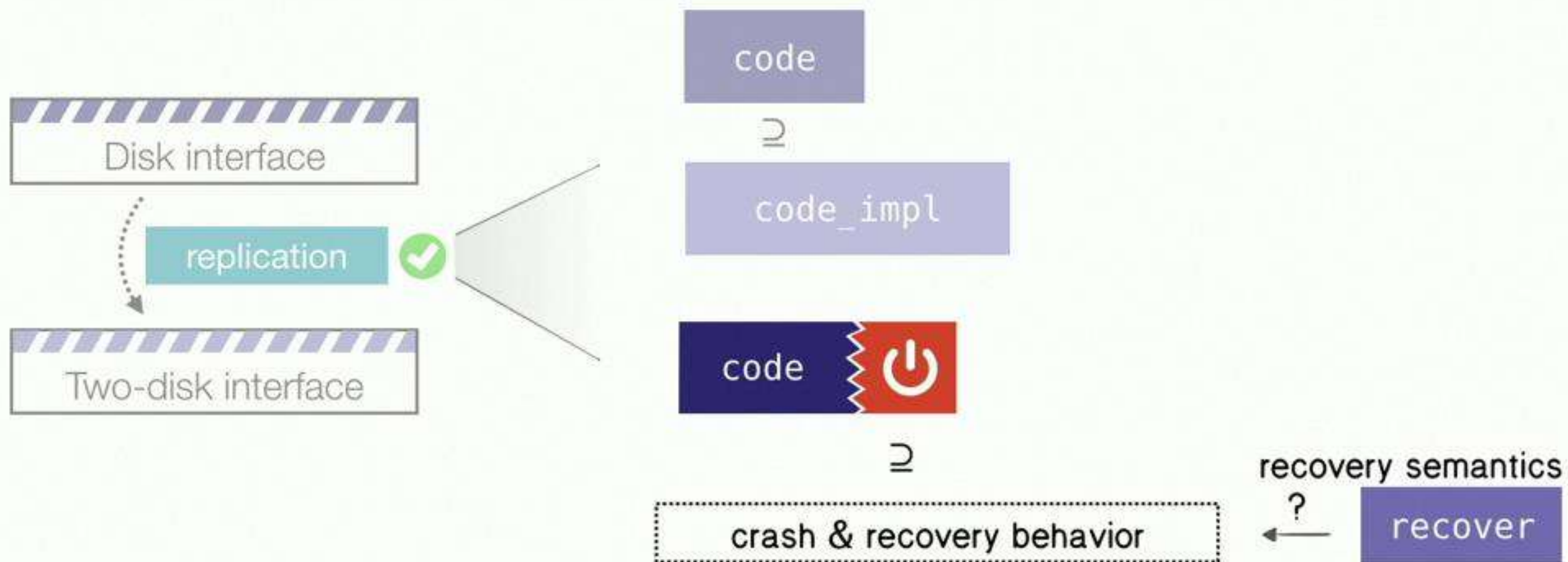


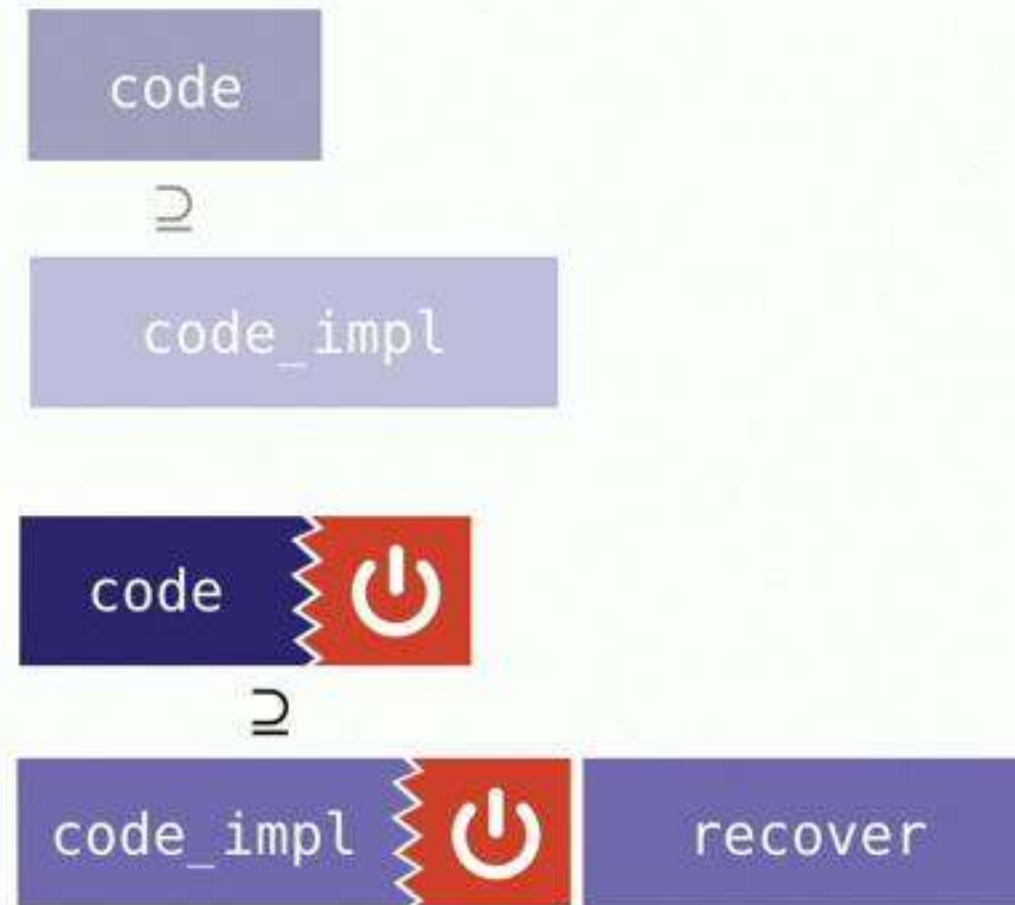
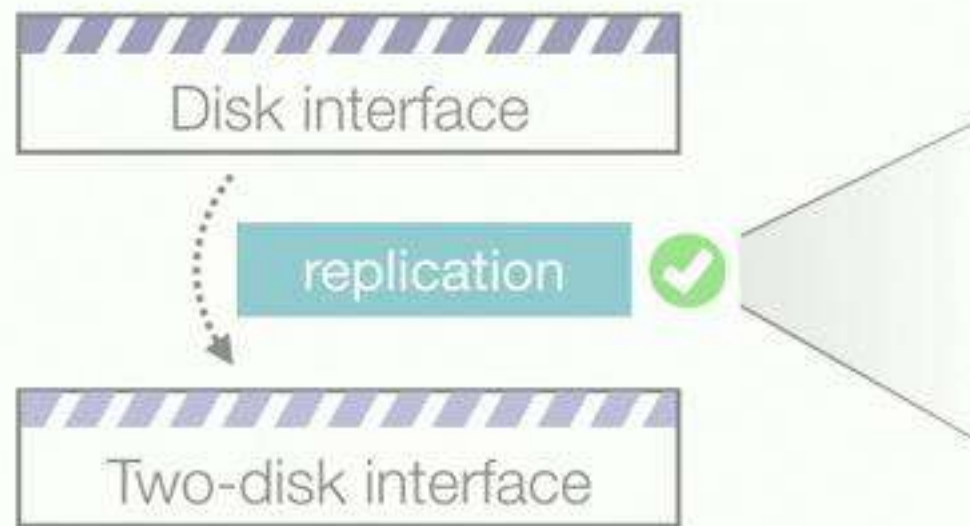
# Extending trace inclusion with recovery



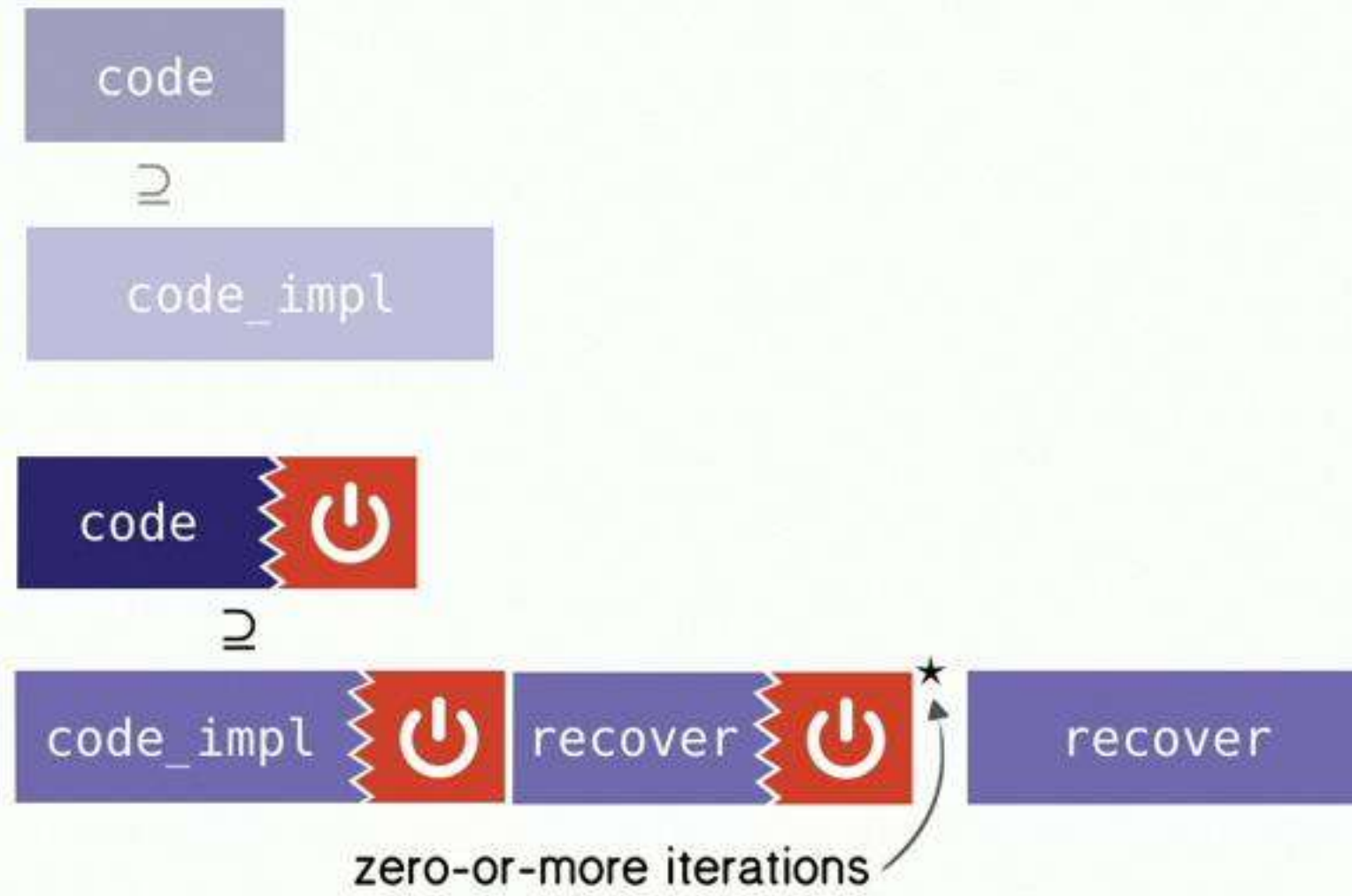
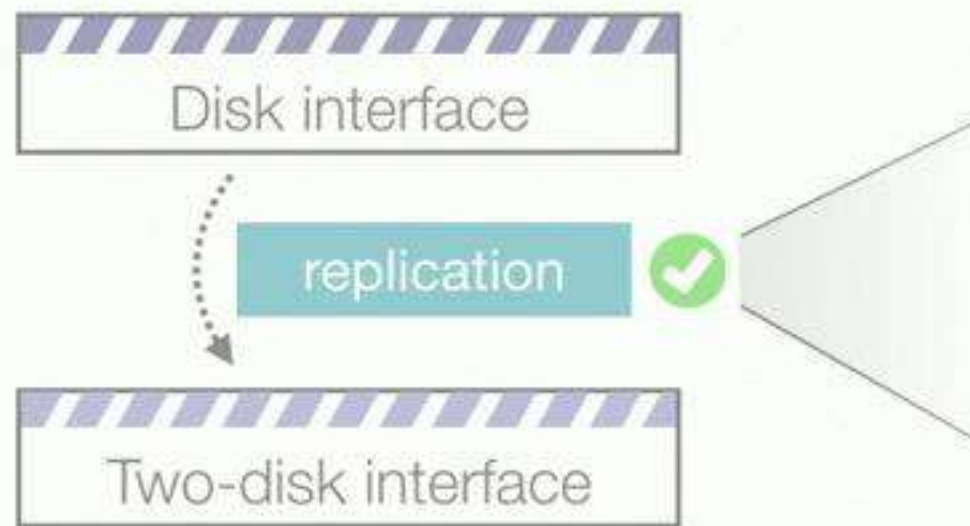




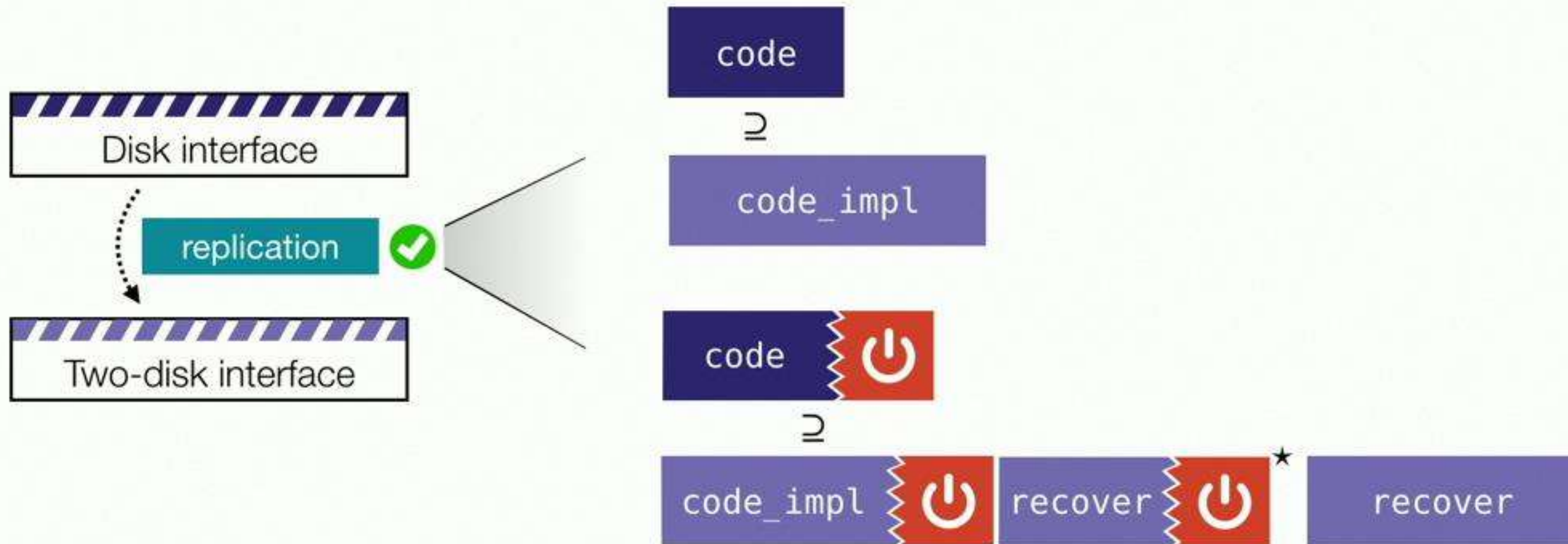








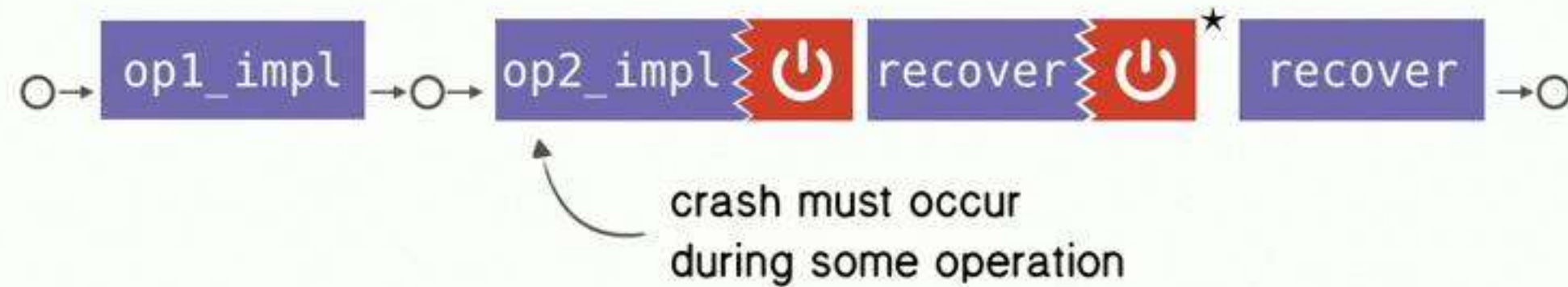
# Trace inclusion, with recovery



# Proving trace inclusion, with recovery



# Proving trace inclusion, with recovery



# Proving trace inclusion, with recovery





# Proving trace inclusion, with recovery

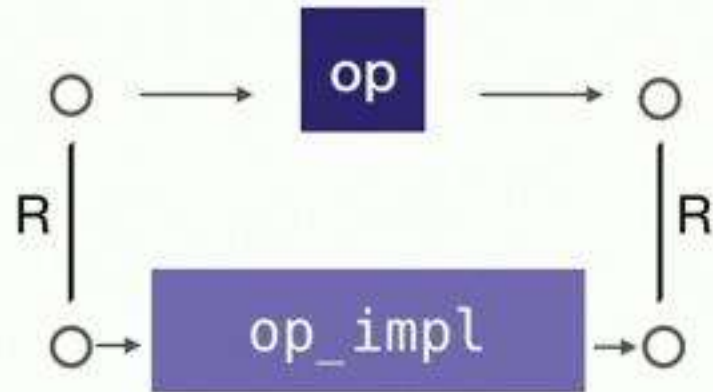


# Proving trace inclusion, with recovery



# Recovery refinement

non-crash execution

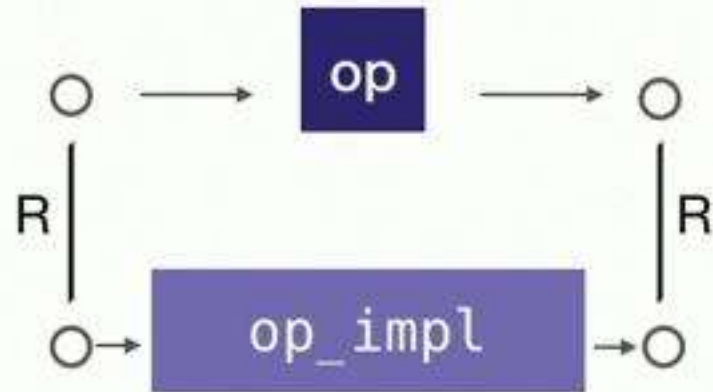


crash and recovery execution



# Recovery refinement

non-crash execution



crash and recovery execution



implies

✓ Trace inclusion

specification behavior  
 $\supseteq$   
running code behavior

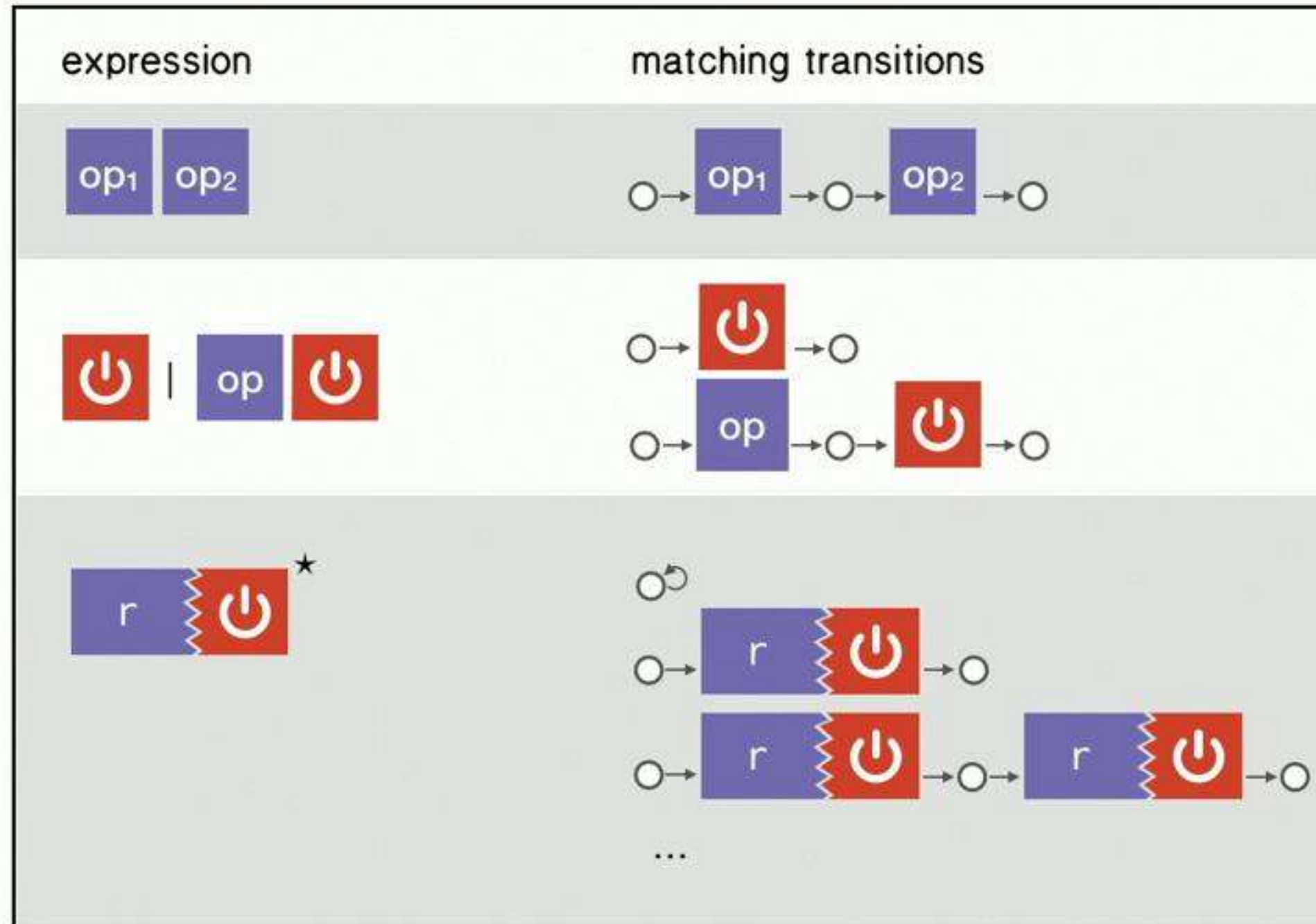
# Composition theorem



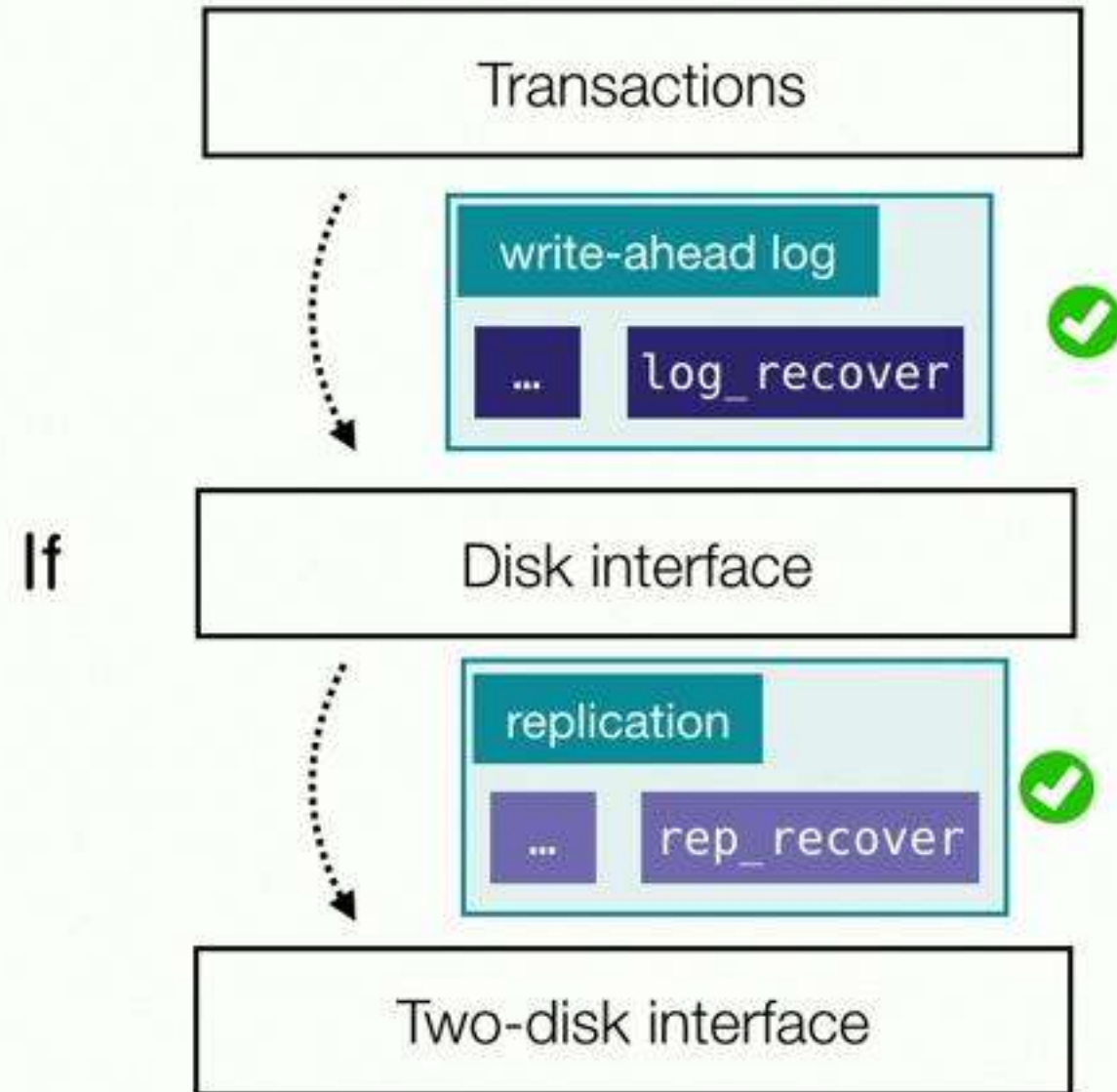
# Kleene algebra for transition relations



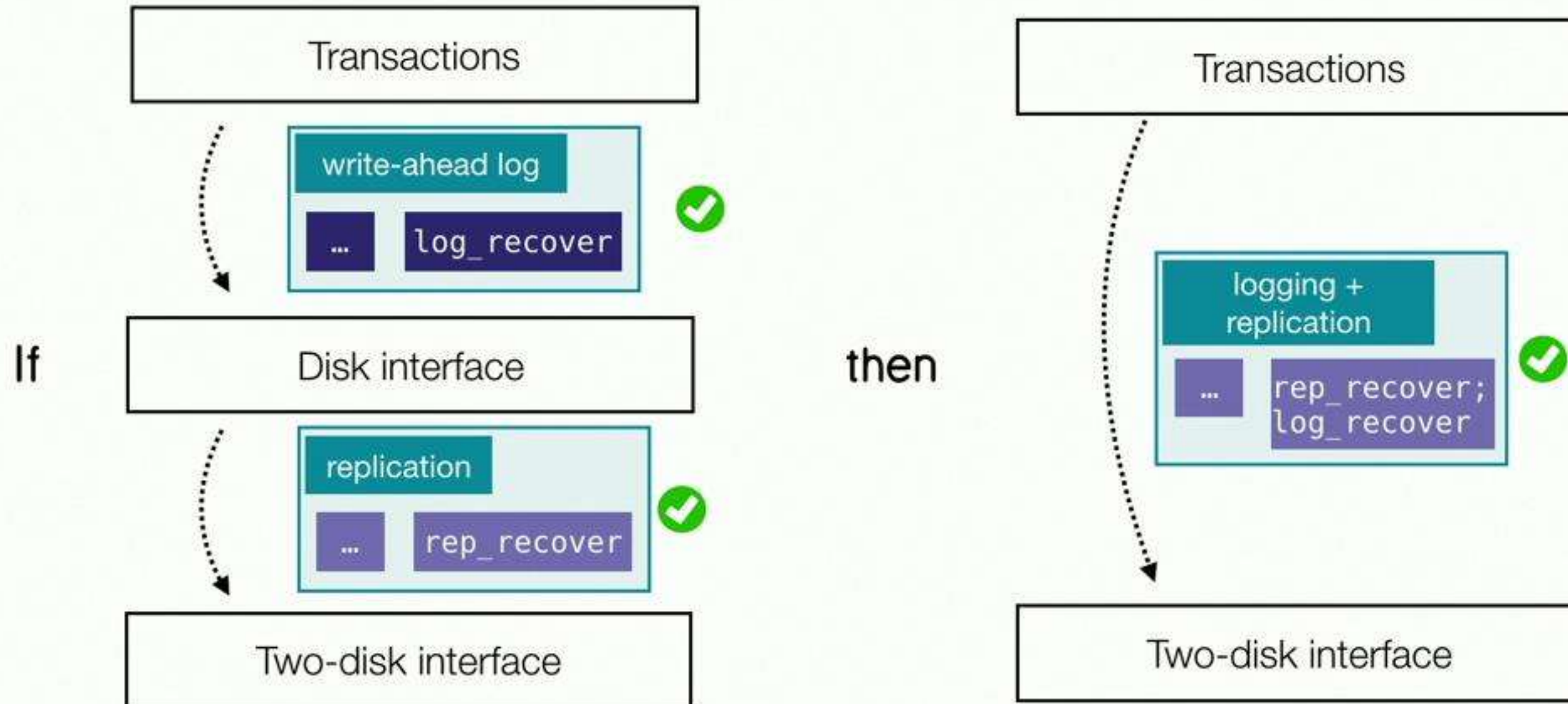
# Kleene algebra for transition relations



# Theorem: recovery refinements compose



# Theorem: recovery refinements compose





# Goal: prove composed recovery correct

rep\_recover



under crashes

log\_recover



under crashes

---

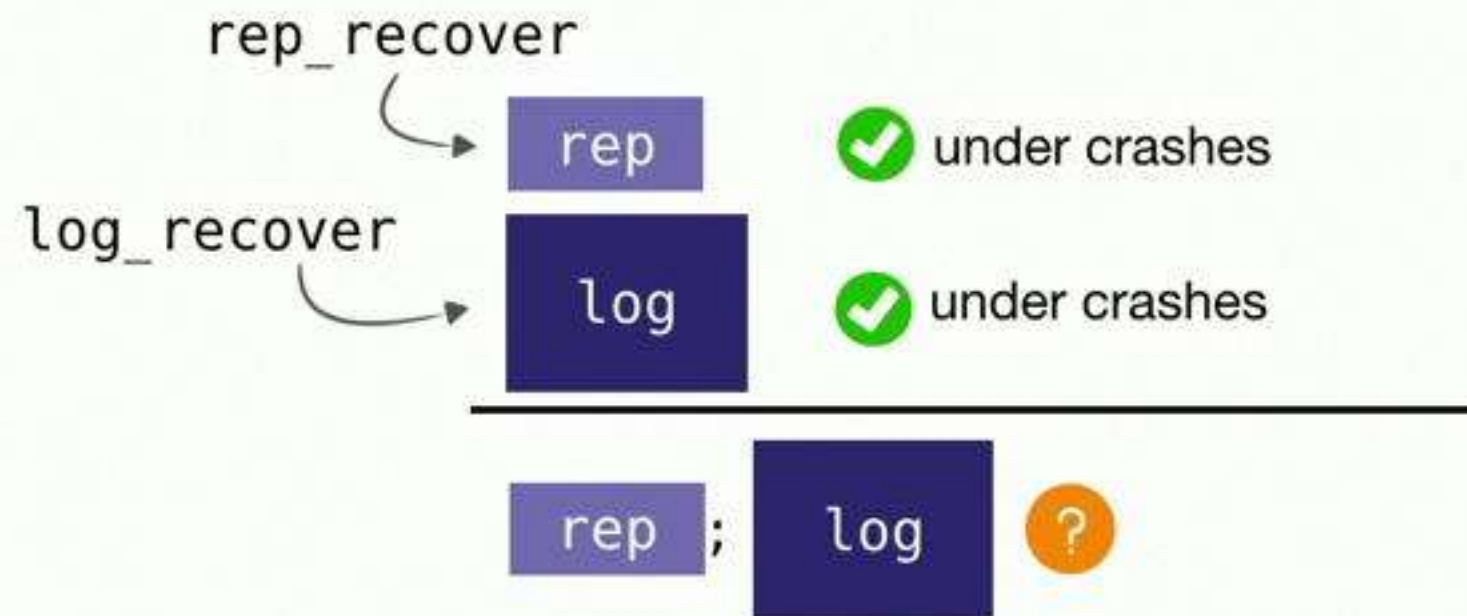
rep\_recover ;

log\_recover

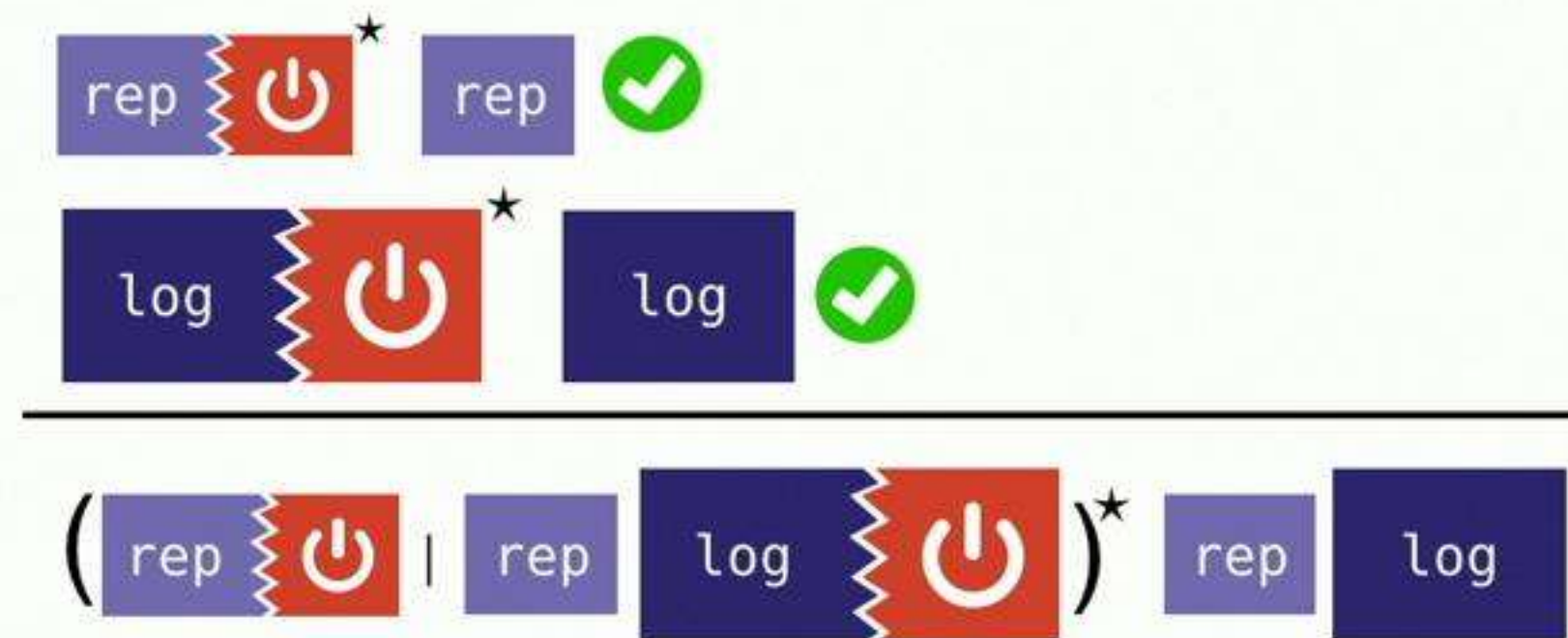




# Goal: prove composed recovery correct









how to re-use recovery proofs here?

# Using Kleene algebra for reasoning



after de-nesting  $(p \mid q)^* = p^*(qp^*)^*$



# Using Kleene algebra for reasoning

$$\left( \text{rep} \text{ } \text{⚡} \mid \text{rep} \text{ } \text{log} \text{ } \text{⚡} \right)^* \text{rep} \text{ } \text{log}$$

after de-nesting  $(p \mid q)^* = p^*(qp^*)^*$

$$= \text{rep} \text{ } \text{⚡}^* \left( \text{rep} \text{ } \text{log} \text{ } \text{⚡} \text{rep} \text{ } \text{⚡}^* \right)^* \text{rep} \text{ } \text{log}$$

after sliding  $(pq)^*p = p(qp)^*$

$$= \text{rep} \text{ } \text{⚡}^* \text{rep} \left( \text{log} \text{ } \text{⚡} \text{rep} \text{ } \text{⚡}^* \text{rep} \right)^* \text{log}$$

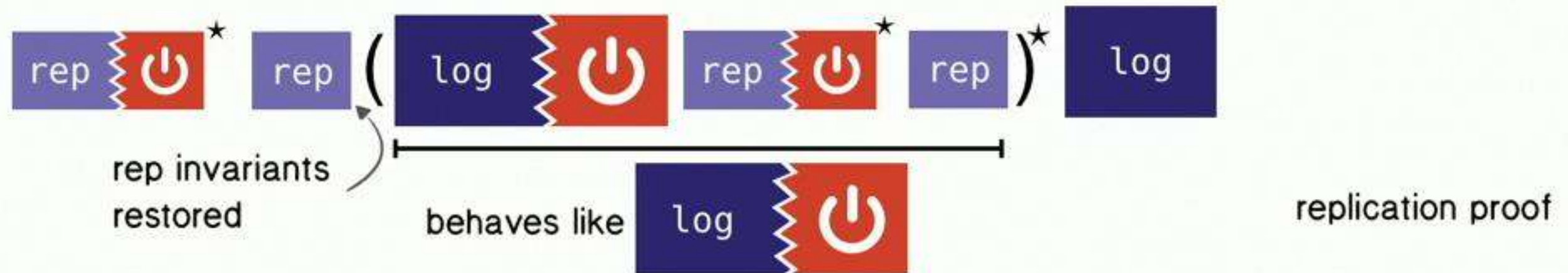
# After rewrite both proofs apply



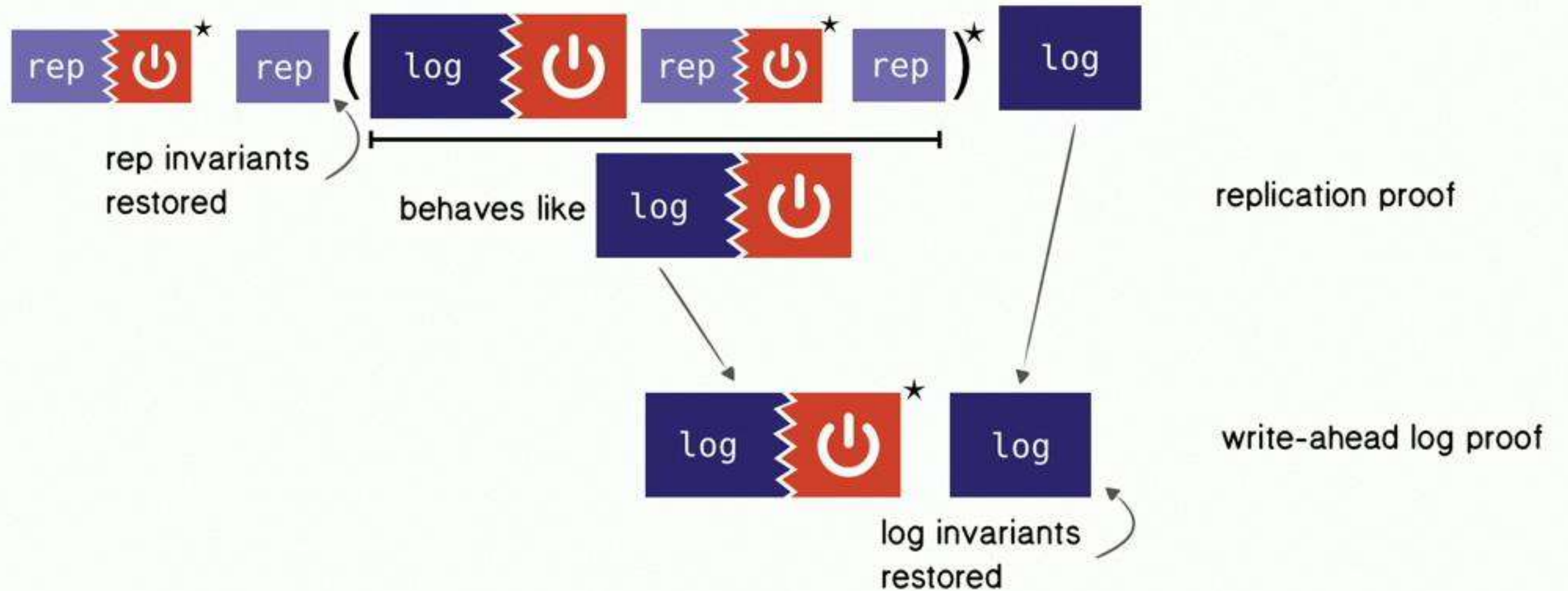
rep invariants  
restored

replication proof

# After rewrite both proofs apply



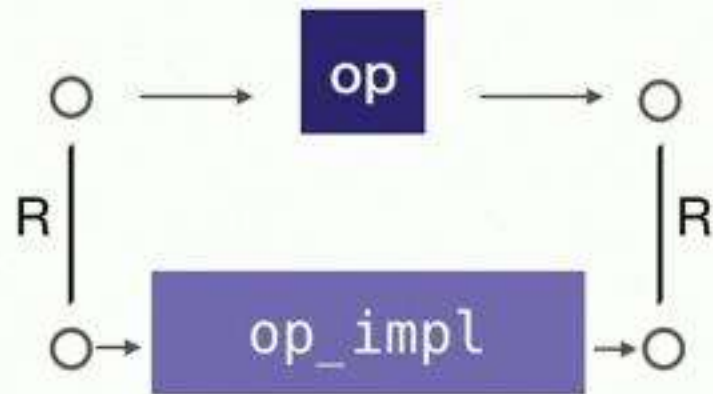
# After rewrite both proofs apply





# Kleene algebra also helps with refinement

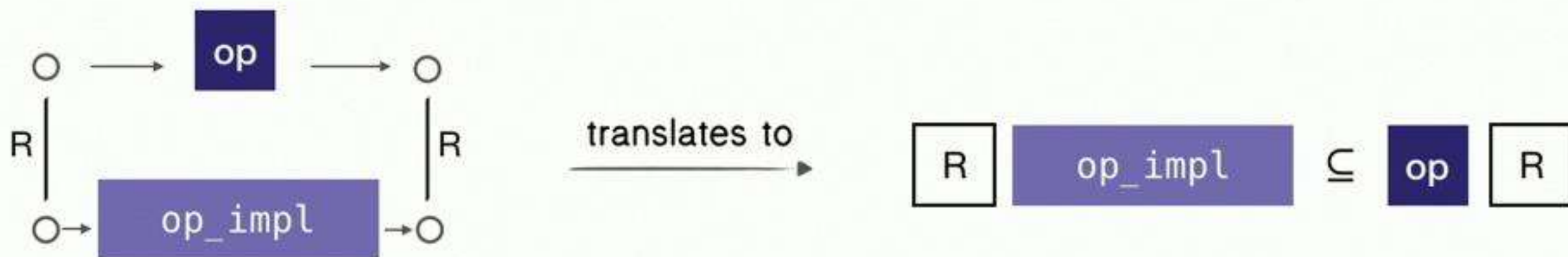
See paper





# Kleene algebra also helps with refinement

See paper



# An anecdote about modularity

# An anecdote about modularity

## Using Crash Hoare Logic for Certifying the FSCQ File System

Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M. Frans Kaashoek, and Nickolai Zeldovich  
*MIT CSAIL*

### Abstract

FSCQ is the first file system v  
(using the Coq proof assistant  
its specification and whose  
FSCQ provably avoids bugs  
systems, such as performing  
barriers or forgetting to zero  
happens at an inopportune tim

## Verifying a high-performance crash-safe file system using a tree specification

Haogang Chen,<sup>†</sup> Tej Chajed, Alex Konradi,<sup>‡</sup> Stephanie Wang,<sup>§</sup> Atalay İleri,  
Adam Chlipala, M. Frans Kaashoek, Nickolai Zeldovich  
*MIT CSAIL*

### ABSTRACT

DFSCQ is the first file system that (1) provides a precise specification for `fsync` and `fdatsync`, which allow applications to achieve high performance and crash safety, and (2) provides a machine-checked proof that its implementation meets this specification. DFSCQ's specification captures the behavior of sophisticated optimizations, including log-

### 1 INTRODUCTION

File systems achieve high I/O performance and crash safety by implementing sophisticated optimizations to increase disk throughput. These optimizations include deferring writing buffered data to persistent storage, grouping many transactions into a single I/O operation, checksumming journal entries, and bypassing the write-ahead log when writing to



# An anecdote about modularity

## Using Crash Hoare Logic for Certifying the FSCQ File System

Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M. Frans Kaashoek, and Nickolai Zeldovich  
*MIT CSAIL*

### Abstract

FSCQ is the first file system verified in Coq (using the Coq proof assistant) whose specification and whose implementation provably avoids bugs in file systems, such as performing operations without barriers or forgetting to zero out space that happens at an inopportune time.

## Verifying a high-performance crash-safe file system using a tree specification (my advisor)

Haogang Chen,<sup>†</sup> Tej Chajed,<sup>‡</sup> Alex Konradi,<sup>‡</sup> Stephanie Wang,<sup>§</sup> Atalay İleri,<sup>¶</sup>  
Adam Chlipala, M. Frans Kaashoek, Nickolai Zeldovich  
*MIT CSAIL*

### ABSTRACT

DFSCQ is the first file system that (1) provides a precise specification for `fsync` and `fdatasync`, which allow applications to achieve high performance and crash safety, and (2) provides a machine-checked proof that its implementation meets this specification. DFSCQ's specification captures the behavior of sophisticated optimizations, including log-

### 1 INTRODUCTION

File systems achieve high I/O performance and crash safety by implementing sophisticated optimizations to increase disk throughput. These optimizations include deferring writing buffered data to persistent storage, grouping many transactions into a single I/O operation, checksumming journal entries, and bypassing the write-ahead log when writing to

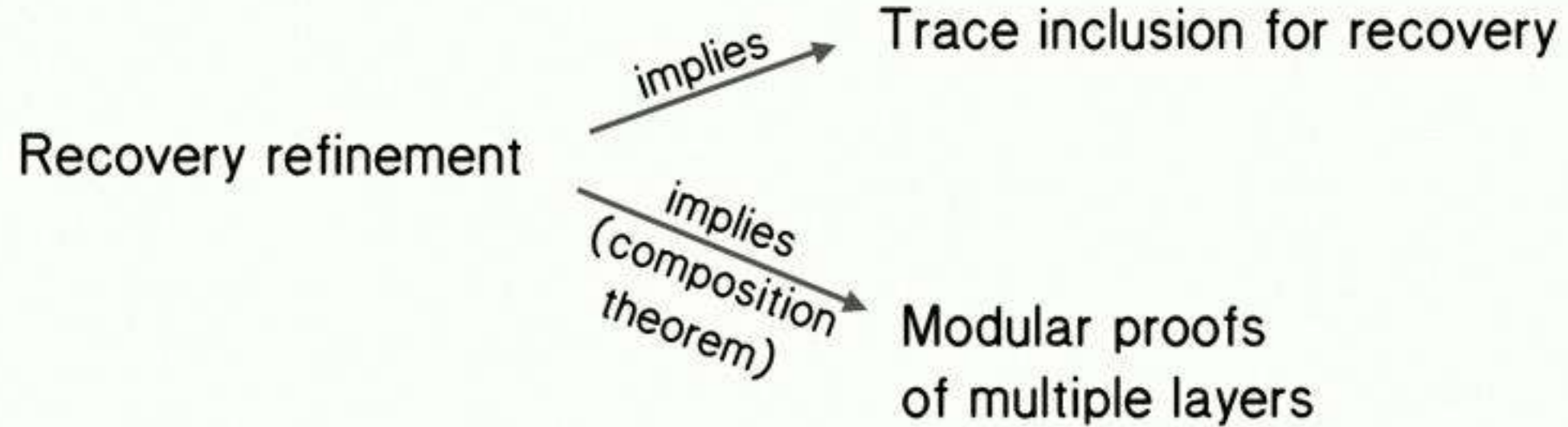
```
def atomic_save(data, path):  
    write_all(data, tmp)  
    rename(tmp, path)  
  
# runs on crash  
def recover():  
    fs_recover()  
    unlink(tmp)
```



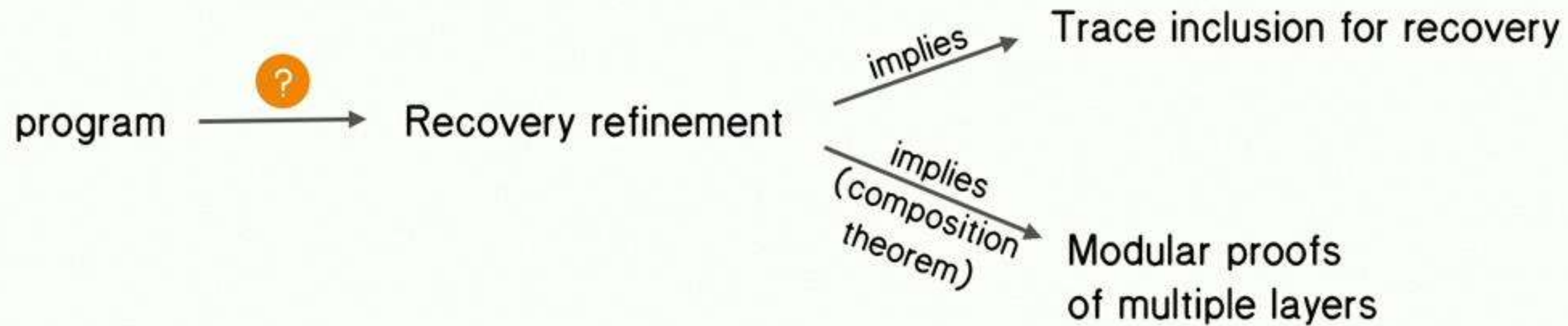
```
def atomic_save(data, path):  
    write_all(data, tmp)  
    rename(tmp, path)  
  
# runs on crash  
def recover():  
    fs_recover() ← this is non-modular and makes  
    unlink(tmp)    the proof much harder
```

Proving this code correct took 1500 lines of proof code!

# Argosy so far



# Argosy so far



# Crash Hoare Logic

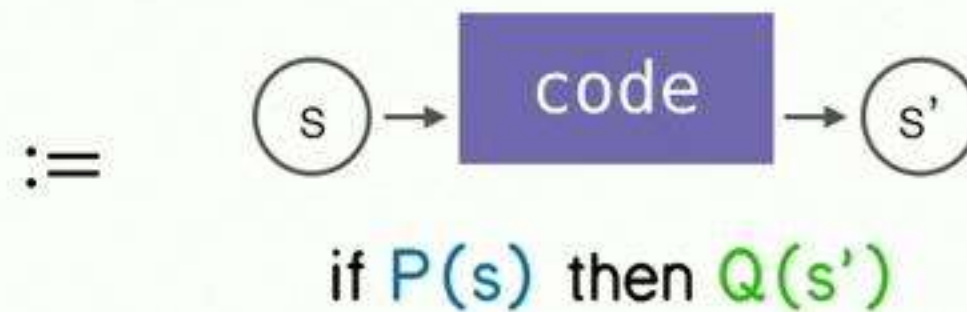
# Hoare Logic

“Hoare triple”     $\{P\}$     code     $\{Q\}$   
                    precondition                      postcondition



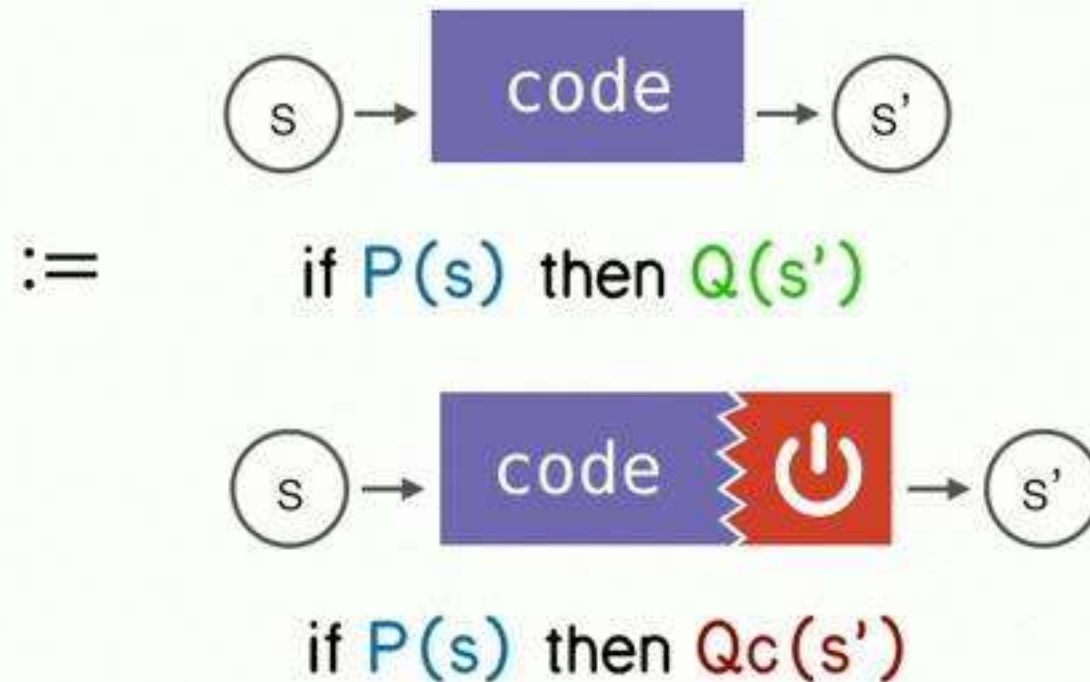
# Hoare Logic

“Hoare triple”       $\{P\}$       code       $\{Q\}$   
precondition      postcondition



# Crash Hoare Logic

“crash specification”     $\{P\}$     code     $\{Q\}$      $\{Q_c\}$   
precondition    postcondition    crash invariant



# Crash Hoare Logic

“recovery specification”

$\{P\}$

precondition

code ↻ recover

$\{Q\}$

postcondition

$\{Qr\}$

recovery postcondition

$:=$



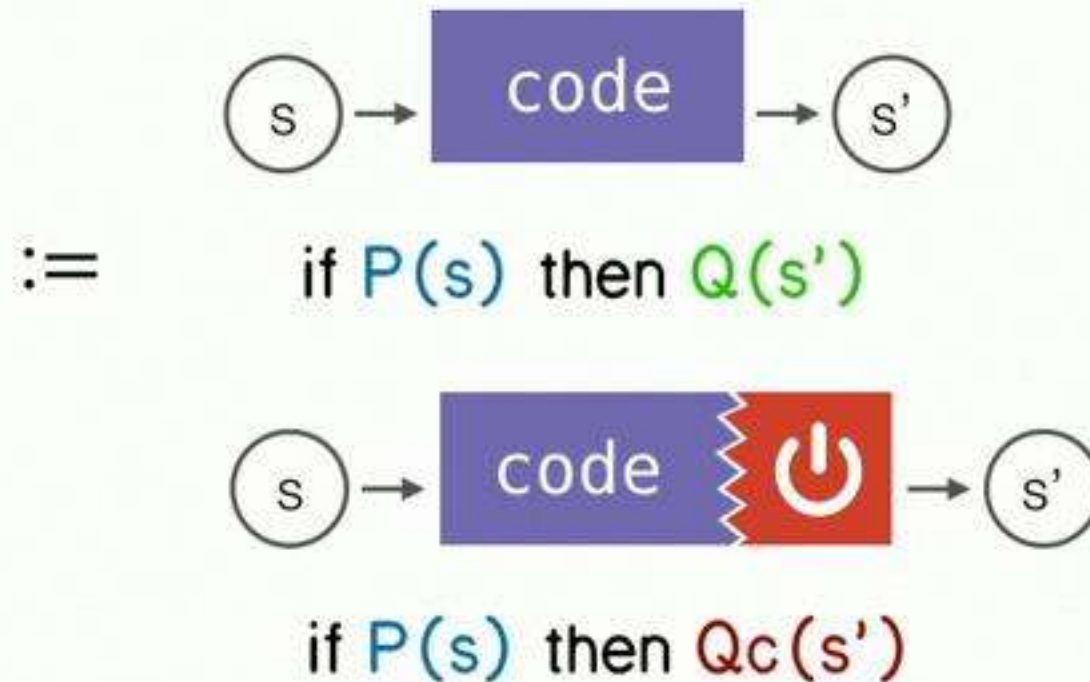
if  $P(s)$  then  $Q(s')$



if  $P(s)$  then  $Qr(s')$

# Crash Hoare Logic

“crash specification”     $\{P\}$     code     $\{Q\}$      $\{Q_c\}$   
precondition    postcondition    crash invariant





# Crash Hoare Logic

“recovery specification”

 $\{P\}$ 

code ↻ recover

 $\{Q\}$  $\{Qr\}$ 

precondition

postcondition

recovery postcondition



if  $P(s)$  then  $Q(s')$

$$\bullet =$$

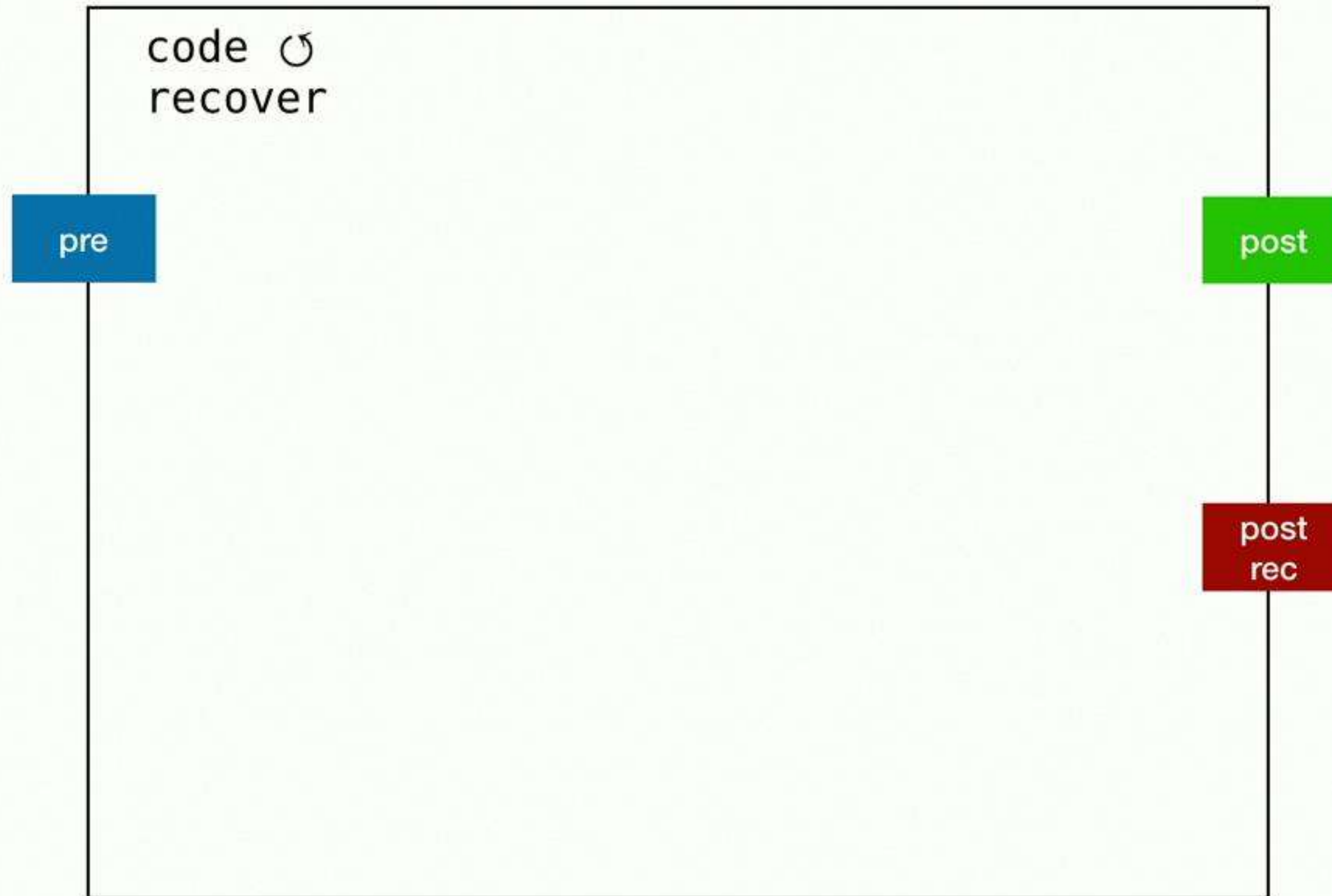

if  $P(s)$  then  $Qr(s')$



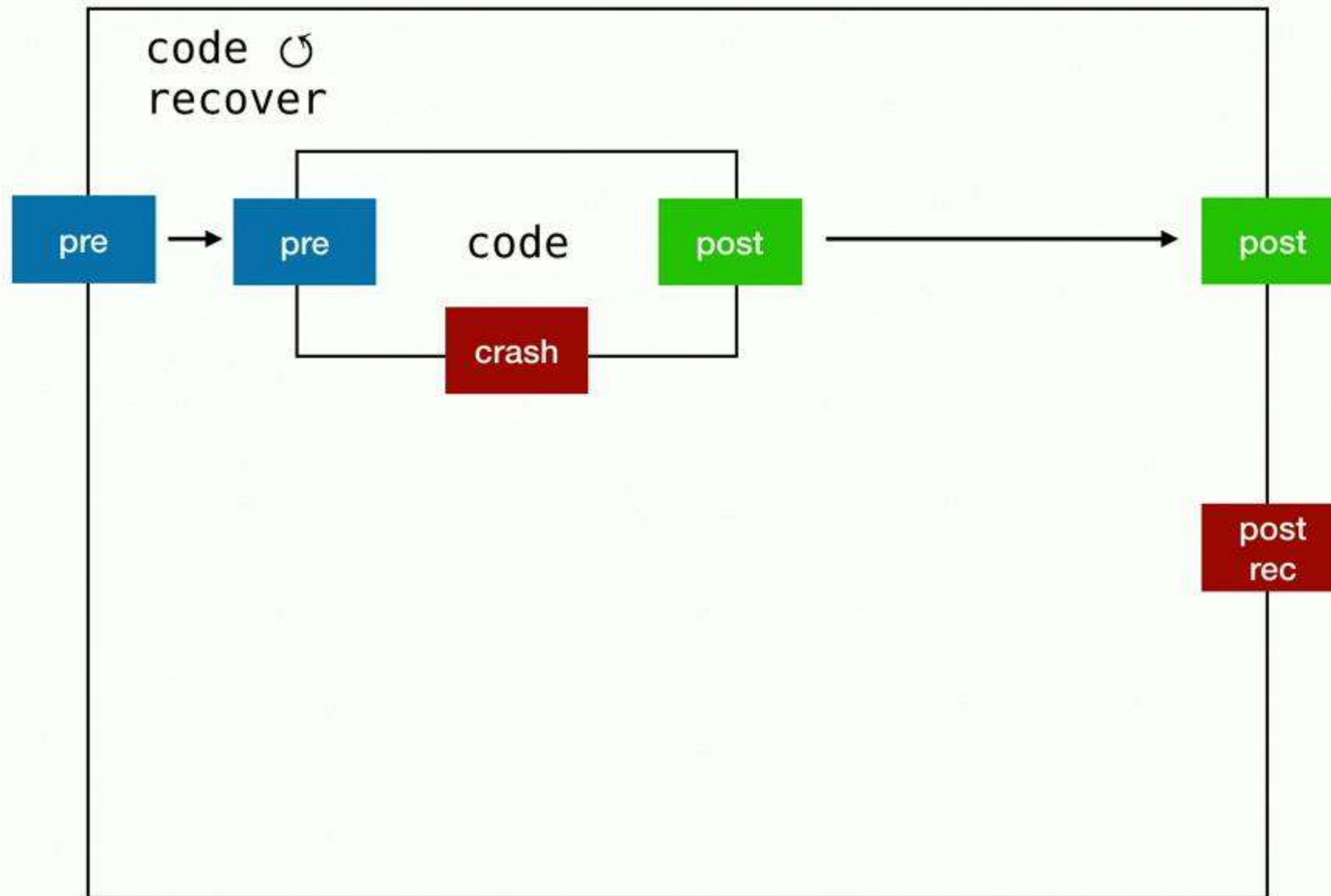
{pre} code ↻ recover {post} {post rec}



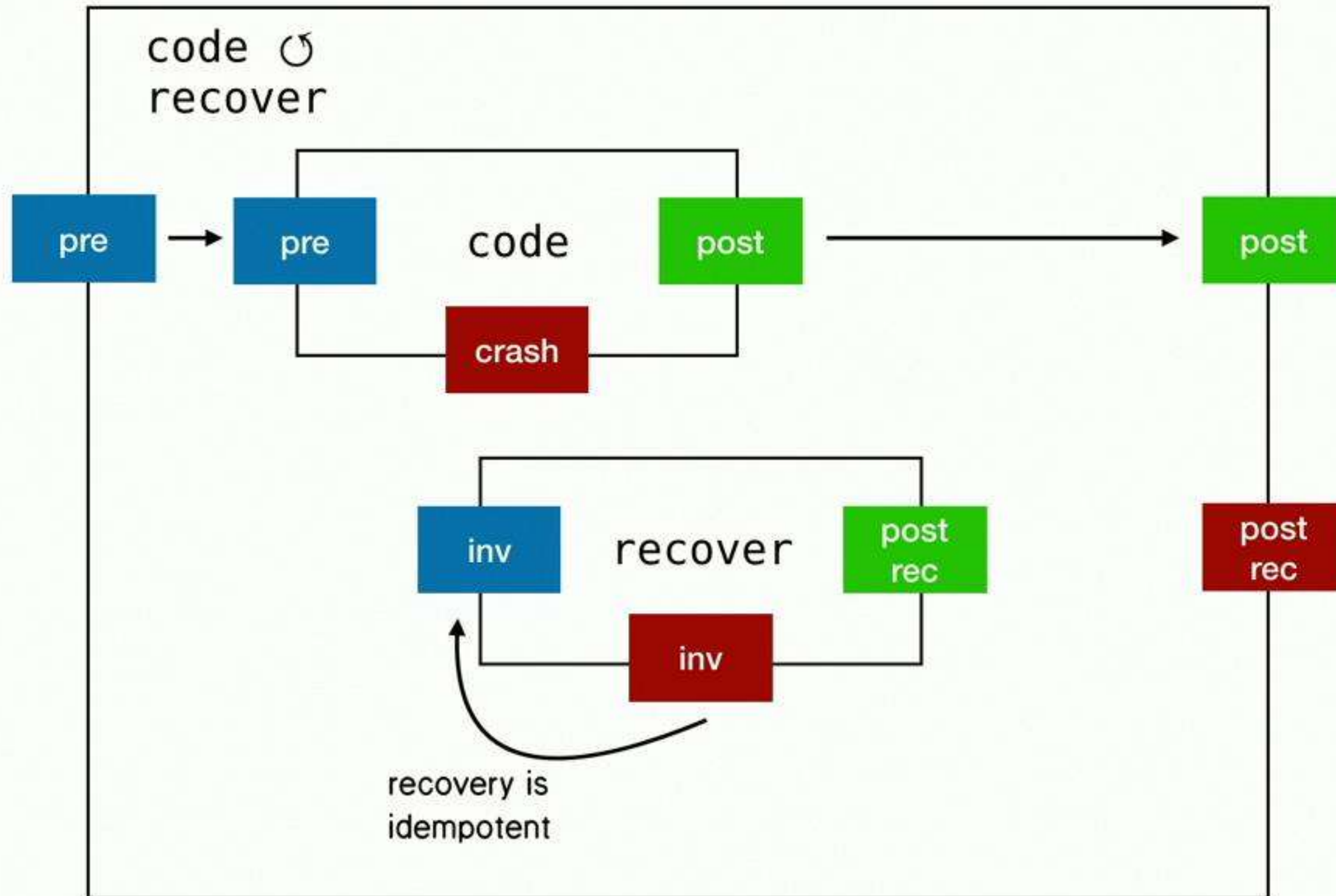
theorem in CHL to prove recovery specs from crash specs



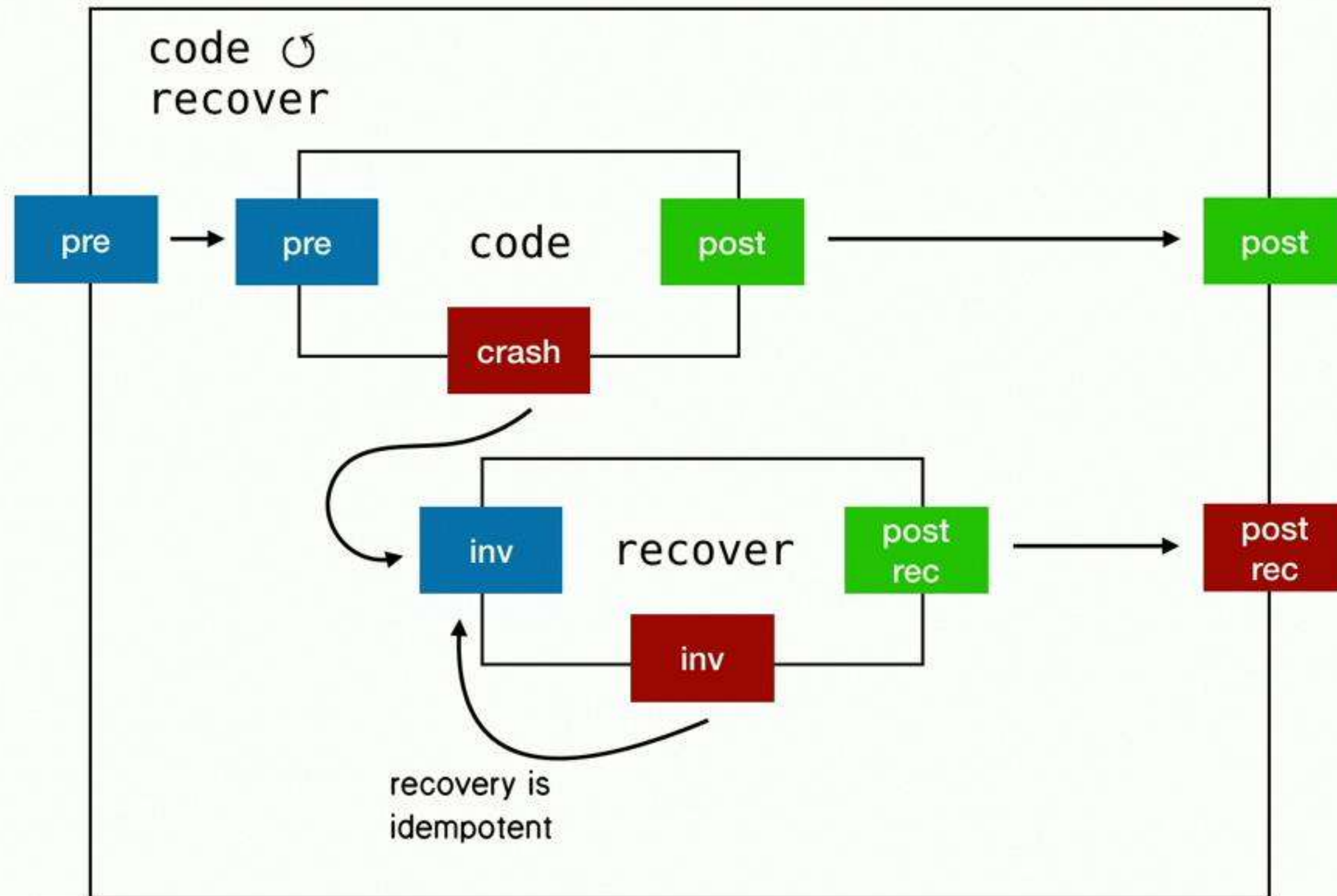
theorem in CHL to prove recovery specs from crash specs



theorem in CHL to prove recovery specs from crash specs



theorem in CHL to prove recovery specs from crash specs





# Argosy connects CHL to recovery refinement

Come up with *abstraction relation*

Prove a *refinement specification* for every operation

Gives recovery refinement for implementation

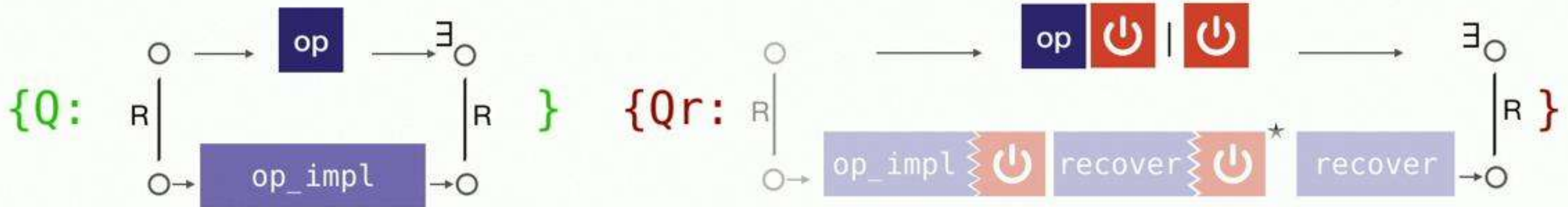
# Recovery refinement as a CHL spec

$\{P\}$  op\_impl  $\circ$  recover  $\{Q\}$   $\{Qr\}$

# Recovery refinement as a CHL spec



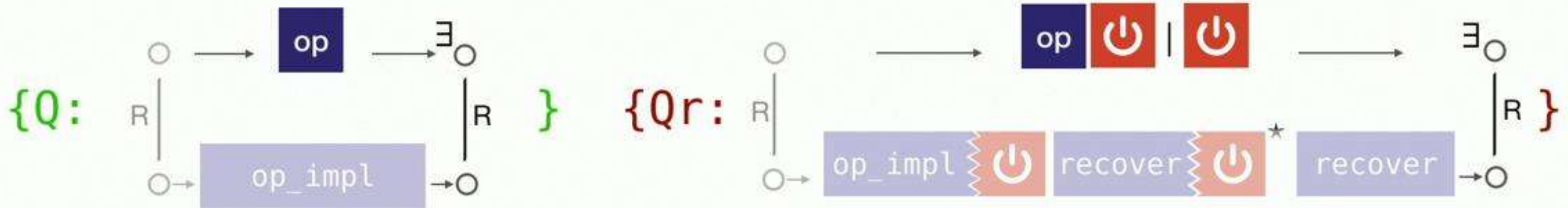
op\_impl  $\circlearrowright$  recover

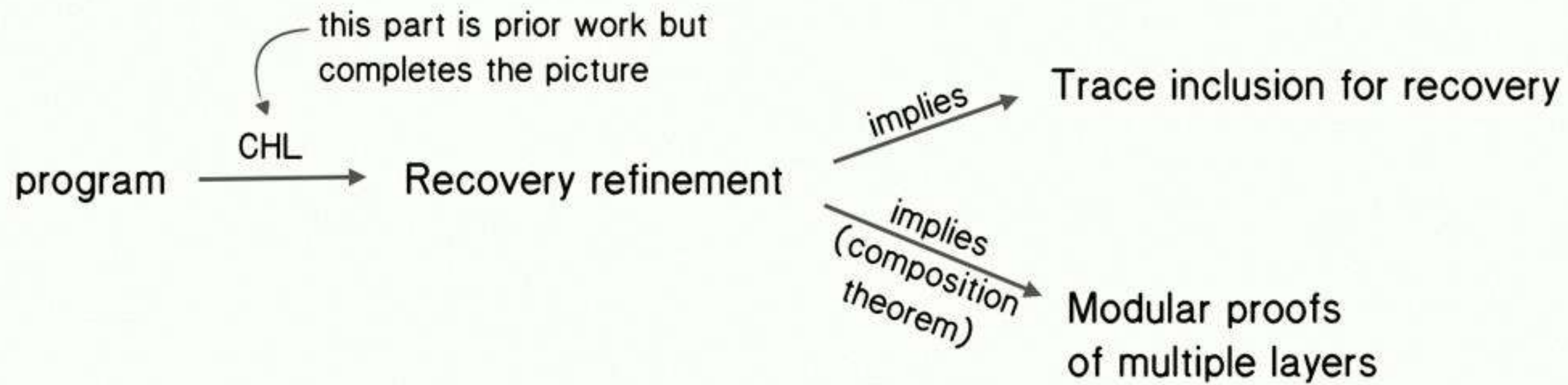


# Recovery refinement as a CHL spec



op\_impl  $\circlearrowright$  recover







# Argosy is implemented and verified in Coq



3,200 lines for framework

4,000 lines for verified example (logging + replication)

Example extracts to Haskell and runs

[github.com/mit-pdos/argosy](https://github.com/mit-pdos/argosy)

# Future work

# Concurrency

Extending Concurrent Separation Logic [originally 2007]

Implemented using Iris [originally POPL 2015]

# Better story for running code

Currently extract to Haskell

Performance problems (esp. for concurrency)

New plan: import Go into Coq

# Usability for students

Argosy spun off from course infrastructure

Now want to backport improvements



# Argosy: modular proofs of layered storage systems

Kleene algebra

$$\left( \text{rep} \begin{array}{|c|} \hline \text{power} \\ \hline \end{array} \mid \text{rep} \begin{array}{|c|} \hline \text{log} \\ \hline \end{array} \begin{array}{|c|} \hline \text{power} \\ \hline \end{array} \right)^*$$



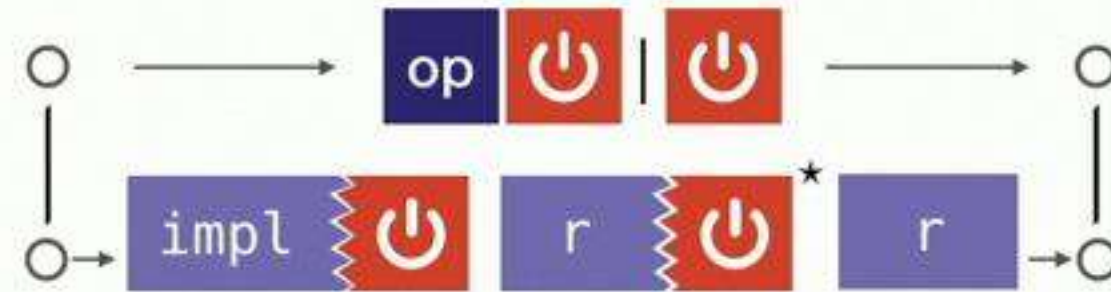
# Argosy: modular proofs of layered storage systems



Kleene algebra

$$(\text{rep} \text{ } \text{⏻} \mid \text{rep} \text{ log } \text{⏻})^*$$

recovery refinement



# Argosy: modular proofs of layered storage systems



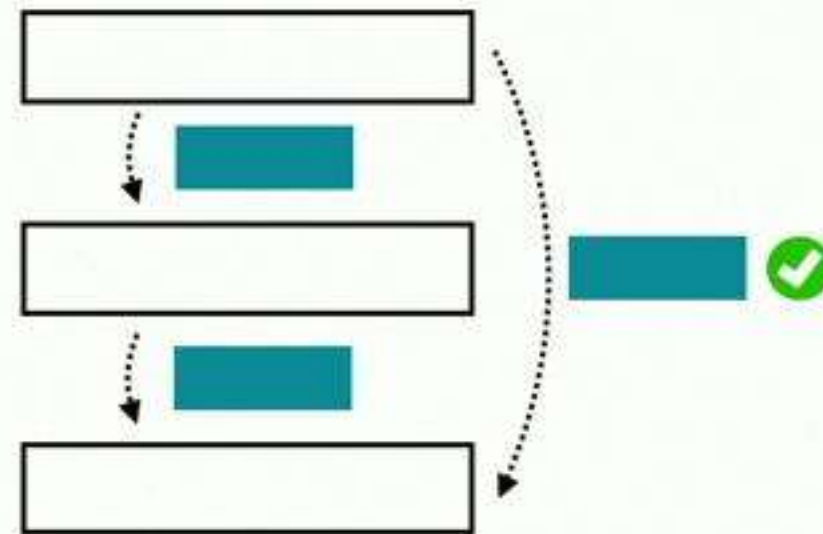
Kleene algebra

$$(\text{rep} \text{ } \text{⏻} \mid \text{rep} \text{ log } \text{⏻})^*$$

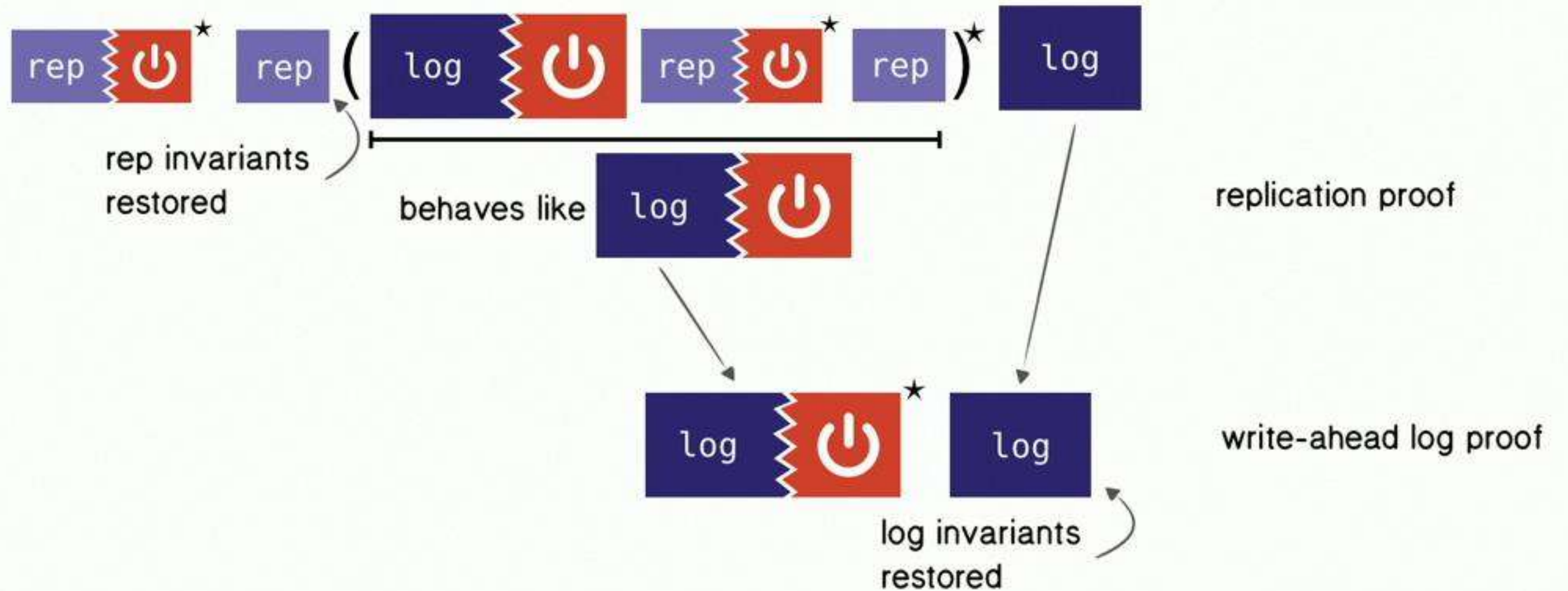
recovery refinement



modular proofs

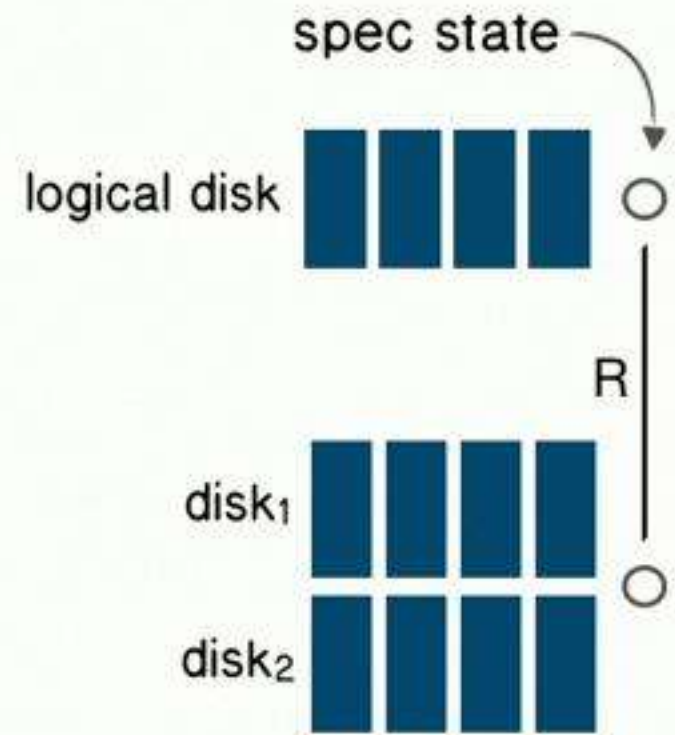


# After rewrite both proofs apply





# Proving correctness with an abstraction relation



1. developer provides abstraction relation  $R$



```
def atomic_save(data, path):  
    write_all(data, tmp)  
    rename(tmp, path)  
  
# runs on crash  
def recover():  
    fs_recover()  
    unlink(tmp)
```