

Conflict-Free Replicated Sets

Annette Bieniusa

together with Nuno Preguiça, Marc Shapiro, Marek Zawirski,
Carlos Baquero, Valter Balesgas

Replicated Sets

- Sets are standard container data type
- Building block for many other data types such as maps or graphs
- Previous replicated set designs seem flawed...

Amazon Dynamo Cart

Replica 1

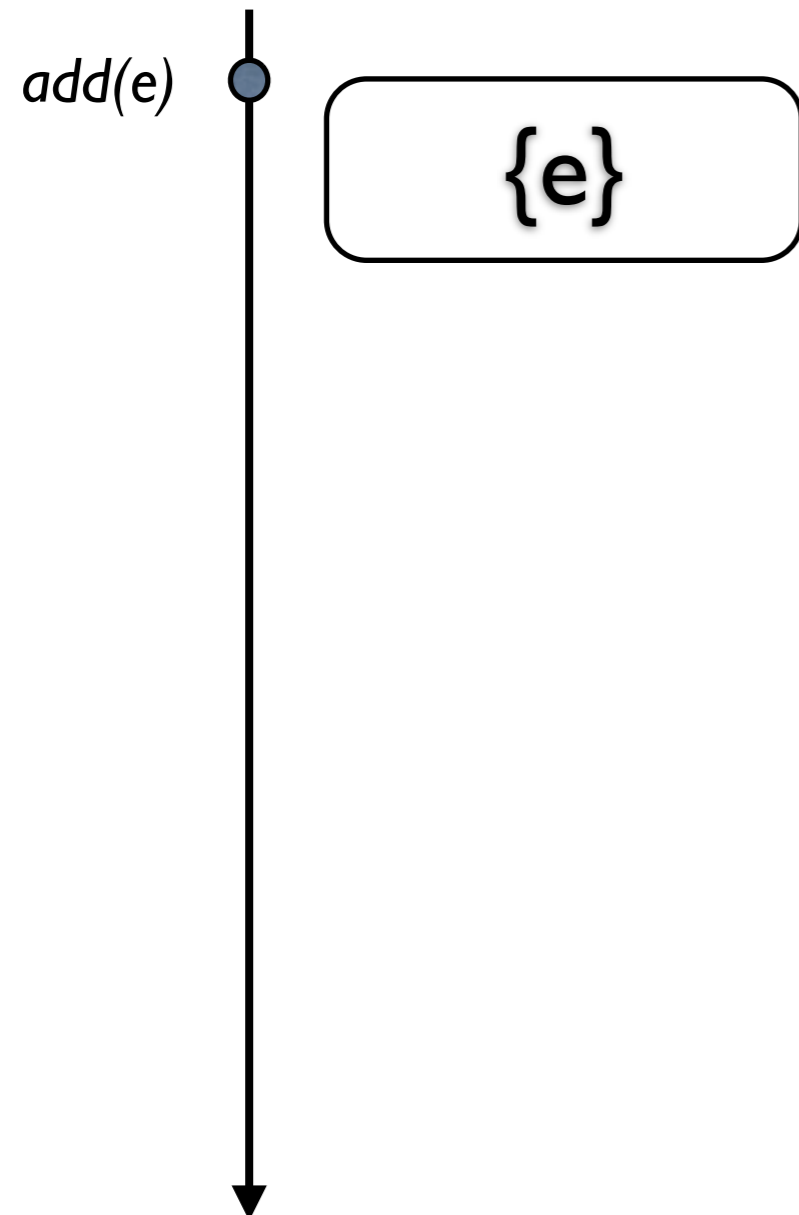


Replica 2



Amazon Dynamo Cart

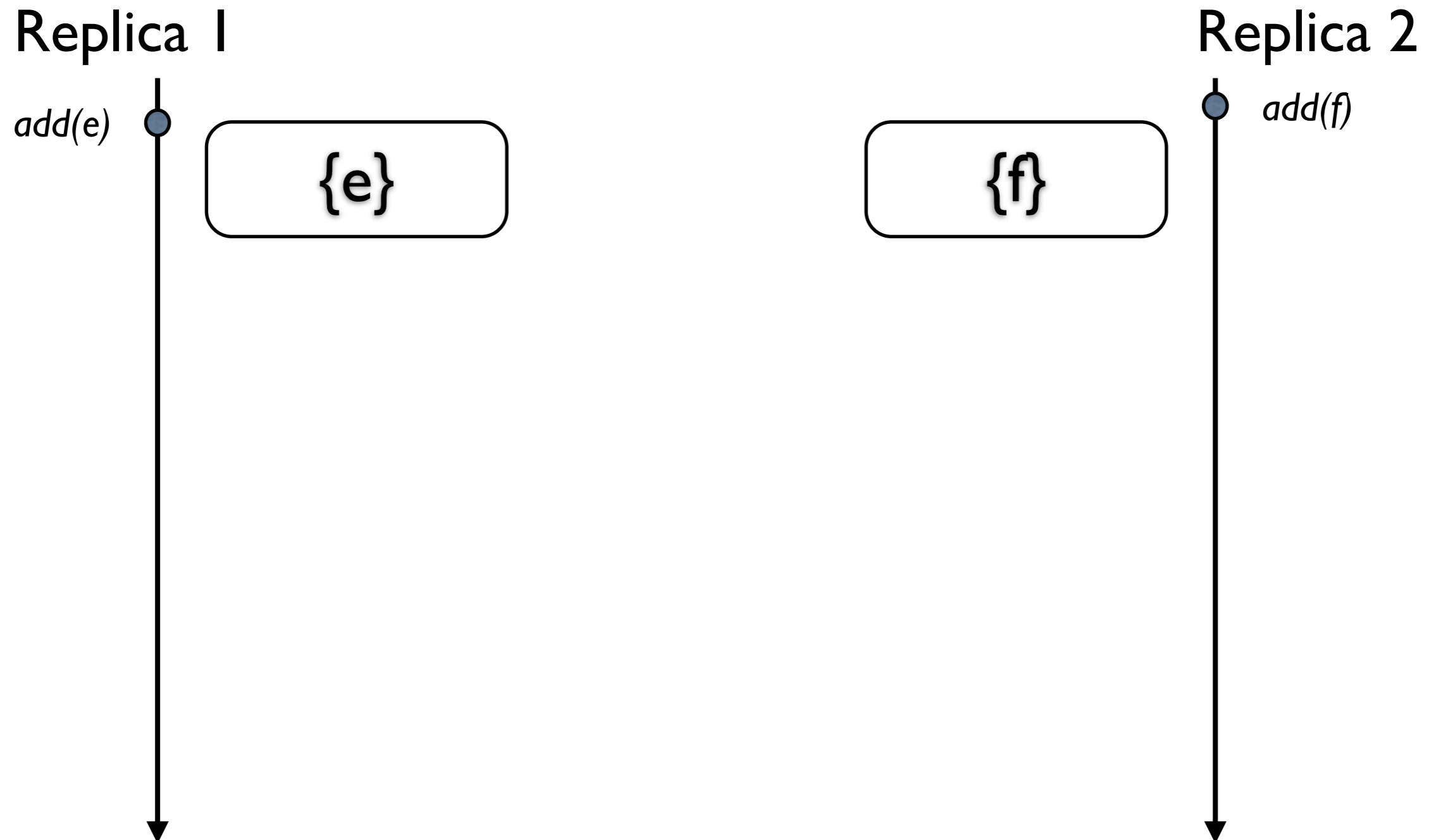
Replica 1



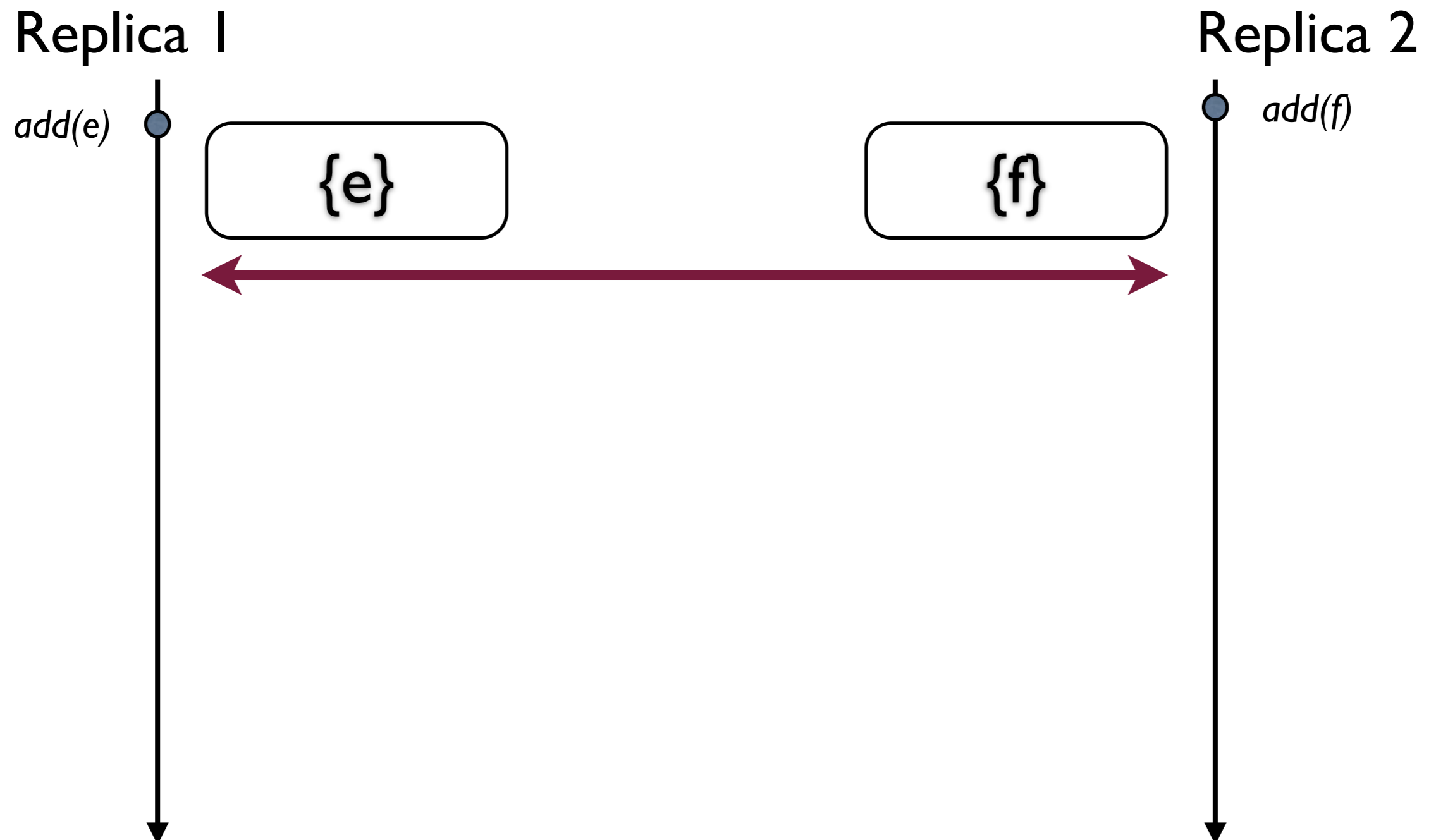
Replica 2



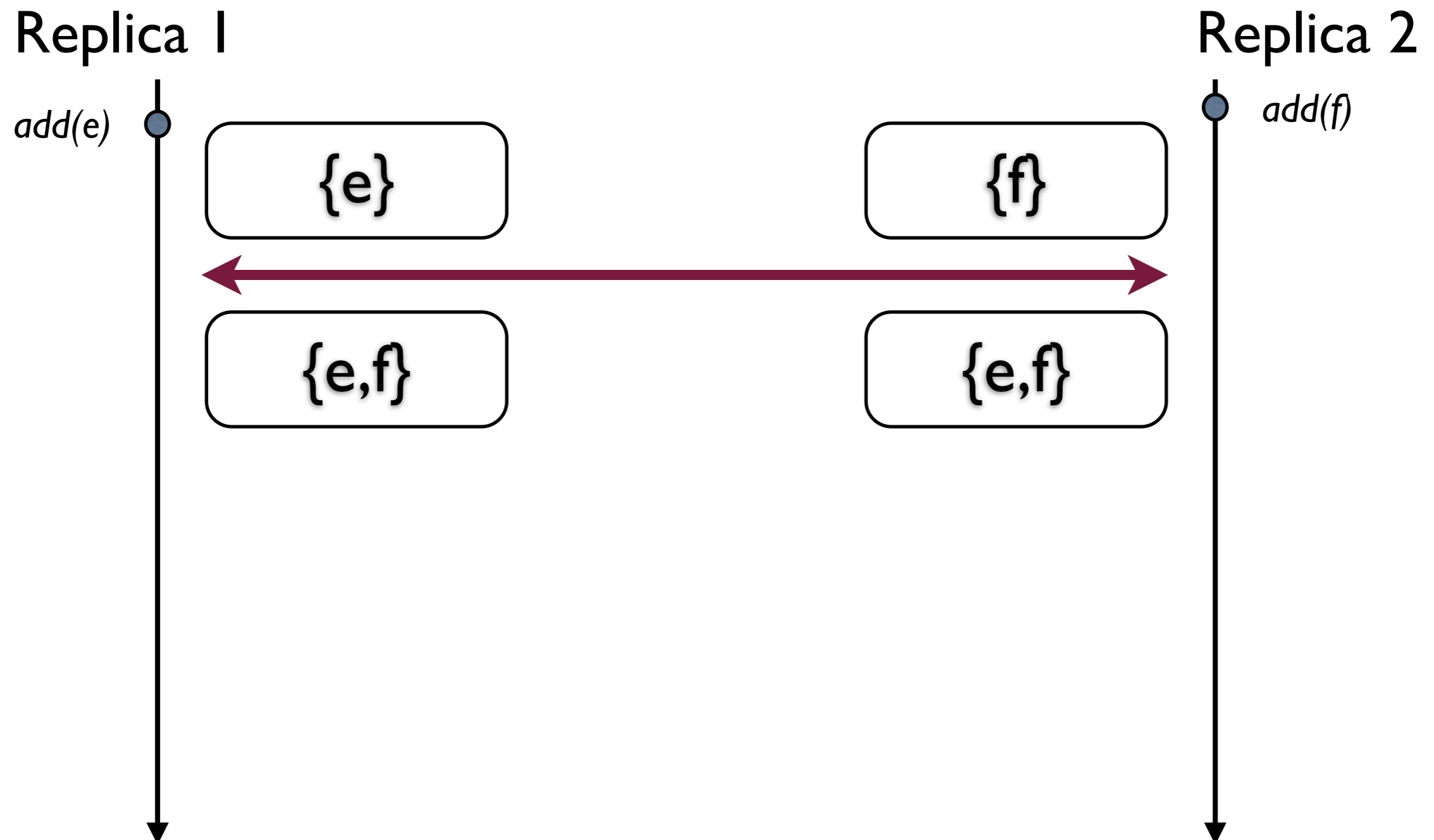
Amazon Dynamo Cart



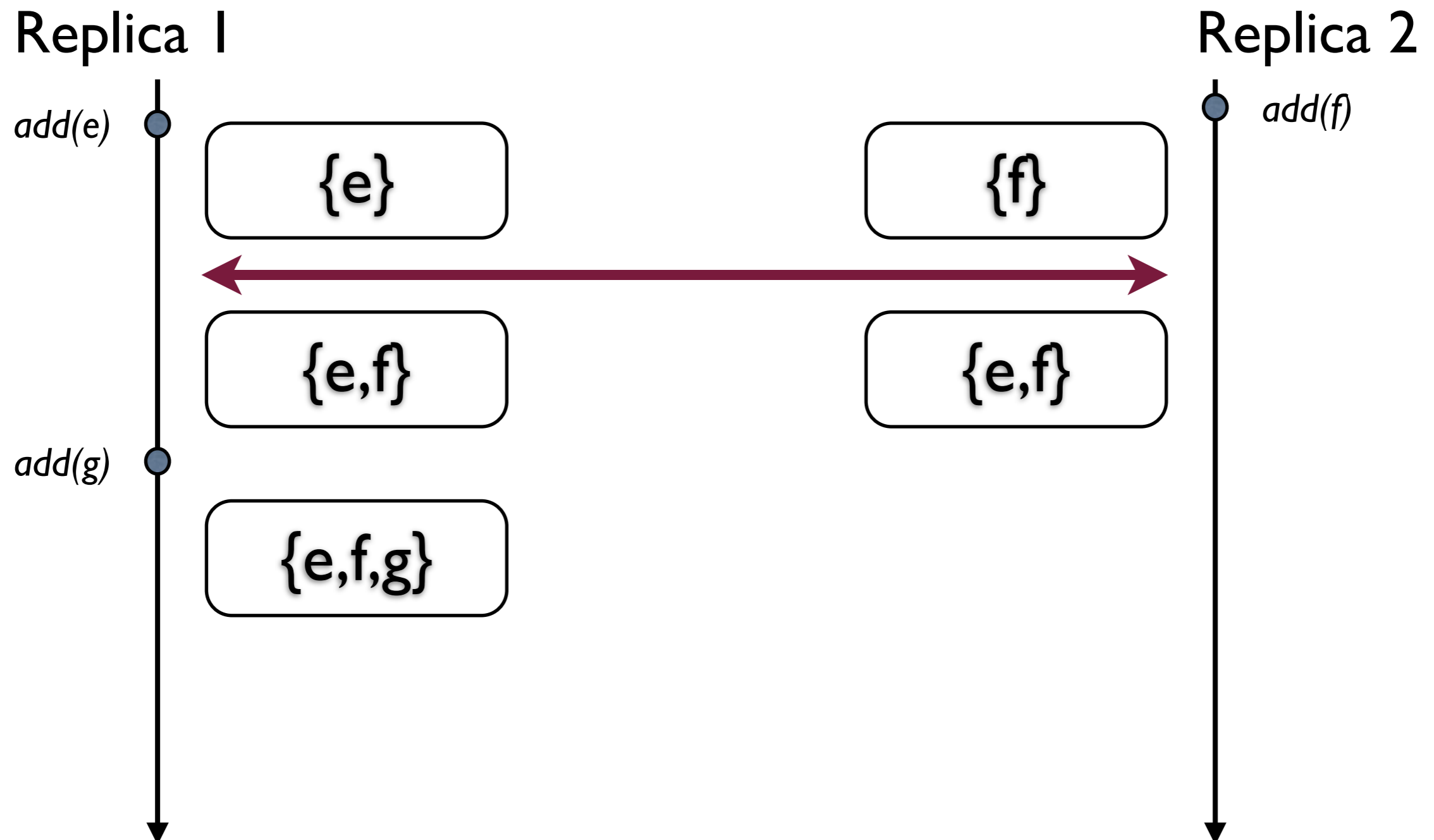
Amazon Dynamo Cart



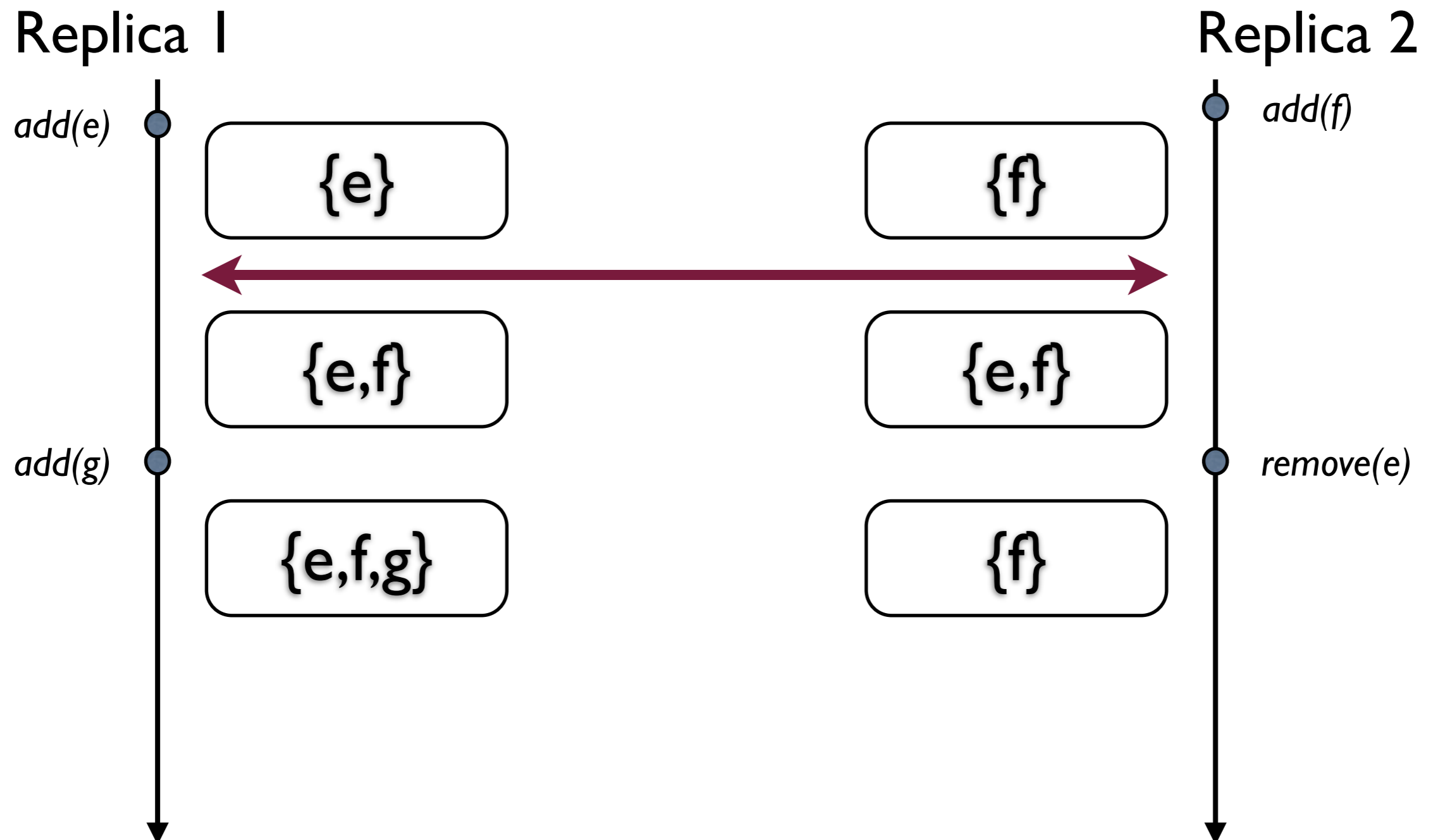
Amazon Dynamo Cart



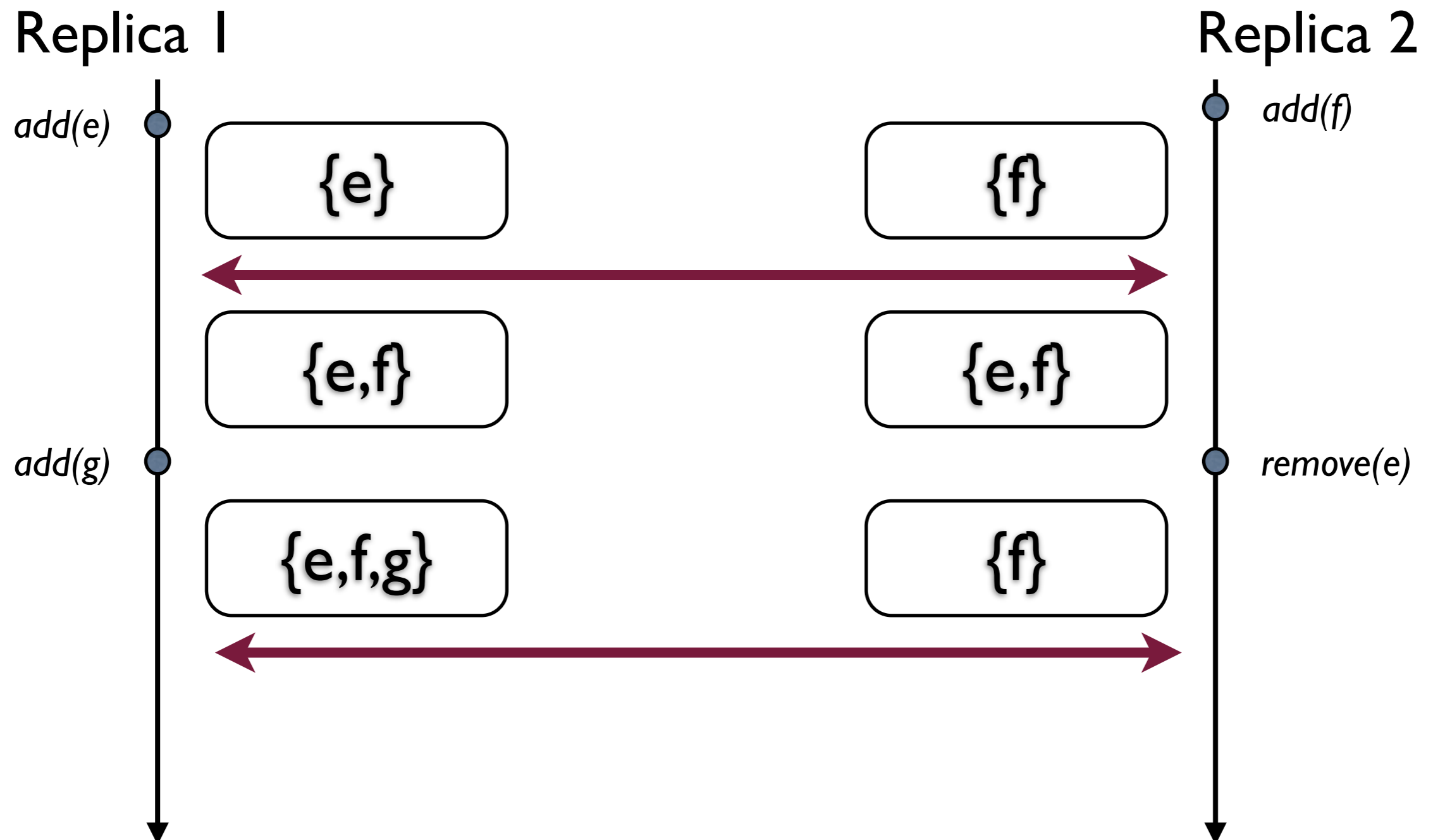
Amazon Dynamo Cart



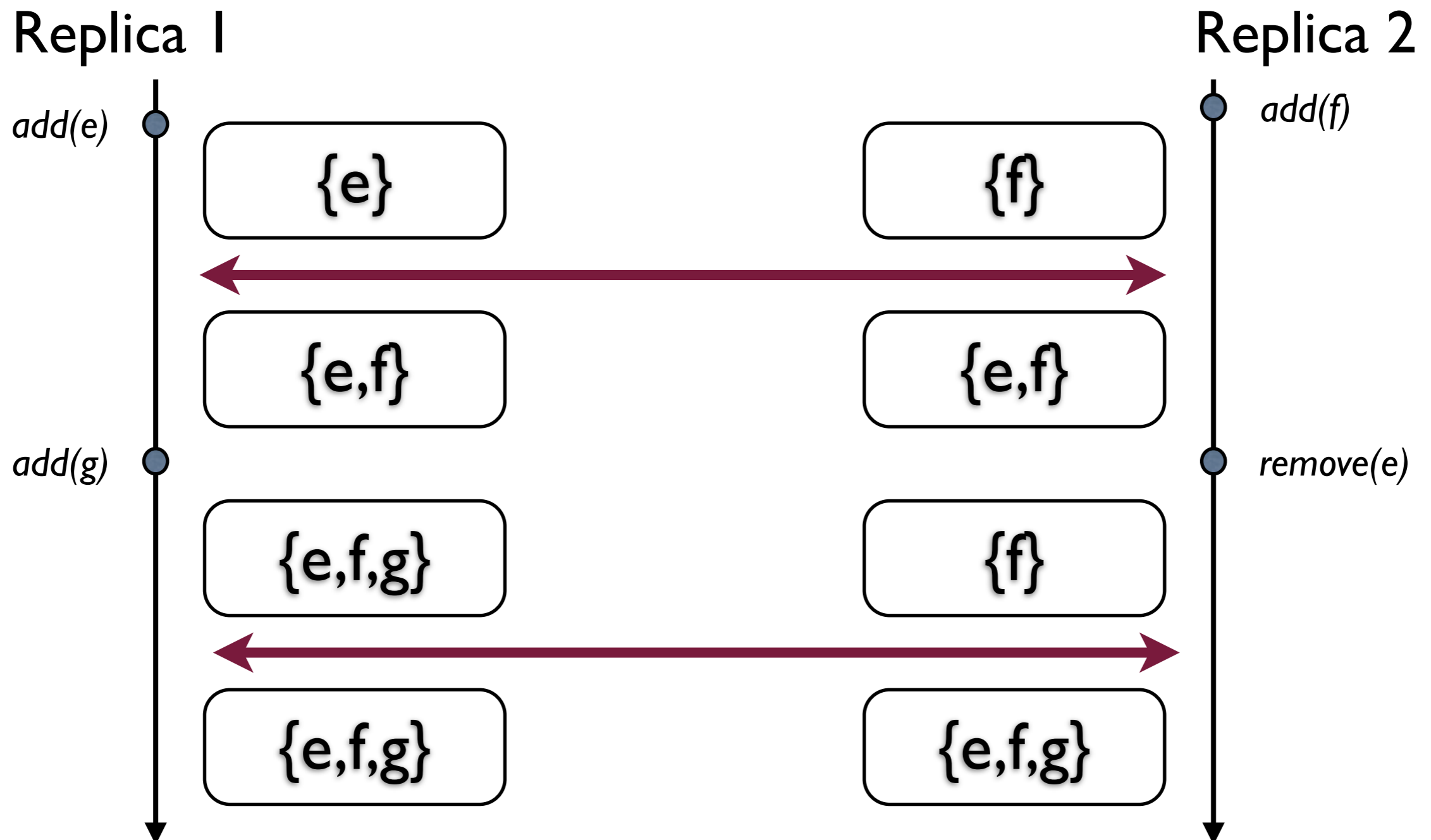
Amazon Dynamo Cart



Amazon Dynamo Cart



Amazon Dynamo Cart



C-Set

Replica 1



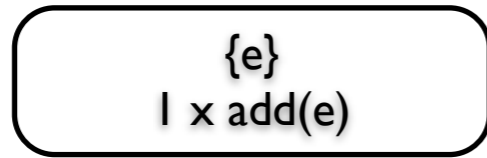
Replica 2



C-Set

Replica 1

$add(e)$



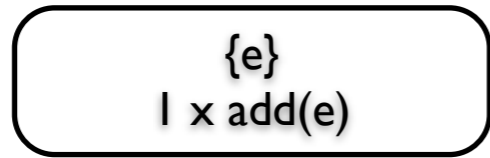
Replica 2



C-Set

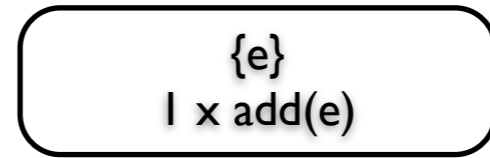
Replica 1

add(e)



Replica 2

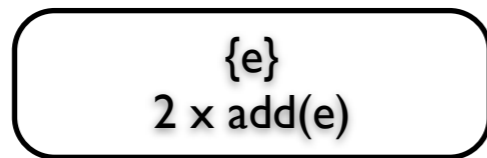
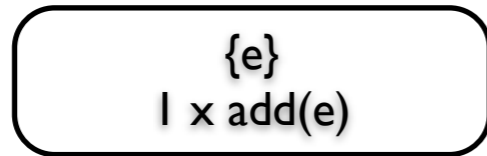
add(e)



C-Set

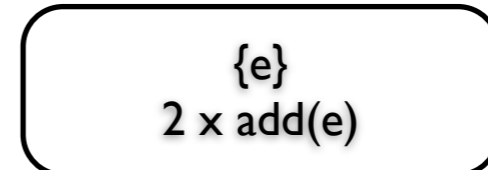
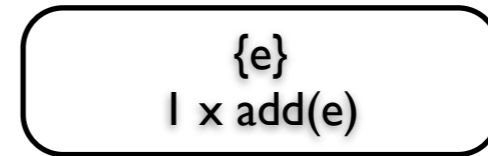
Replica 1

$add(e)$

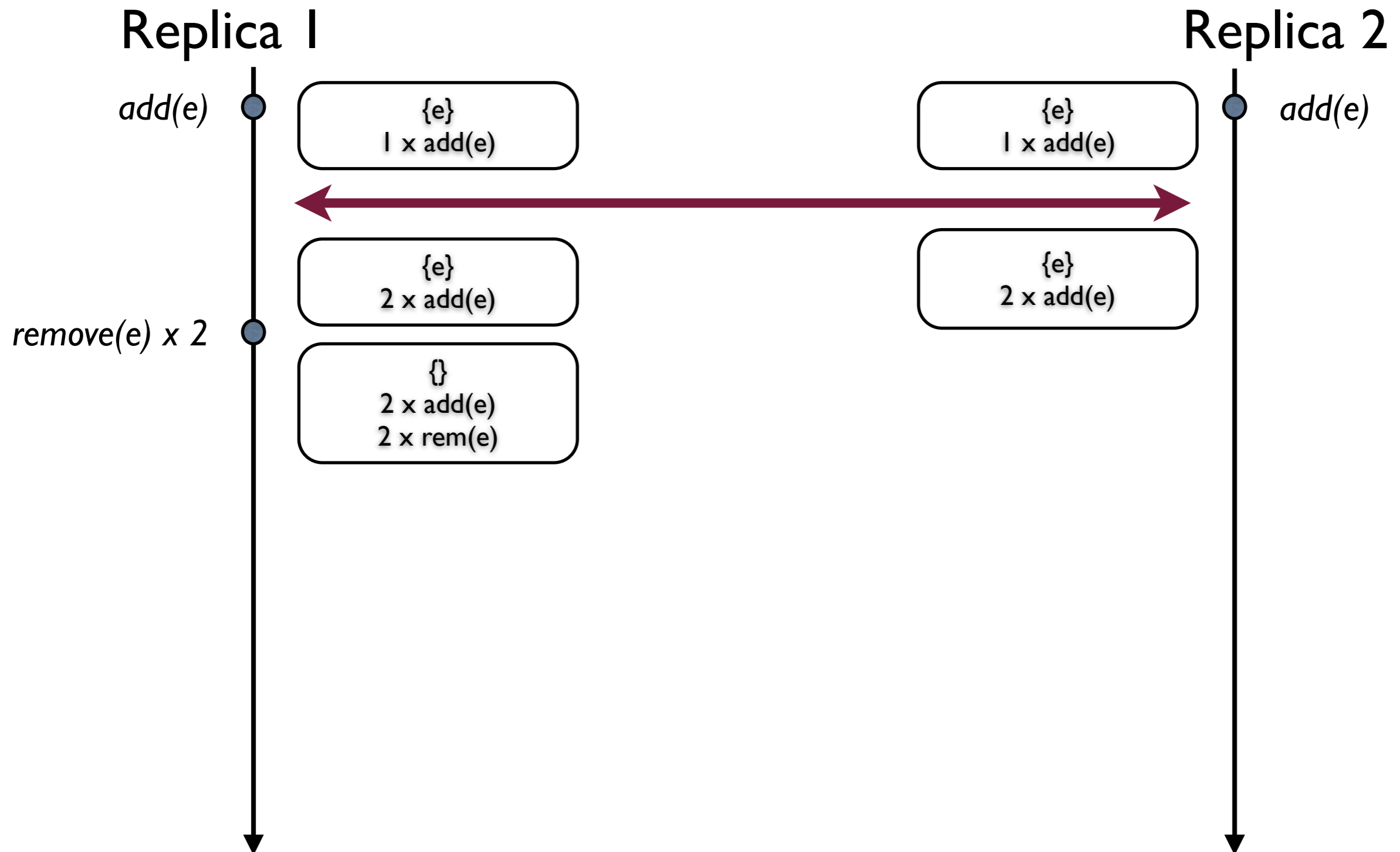


Replica 2

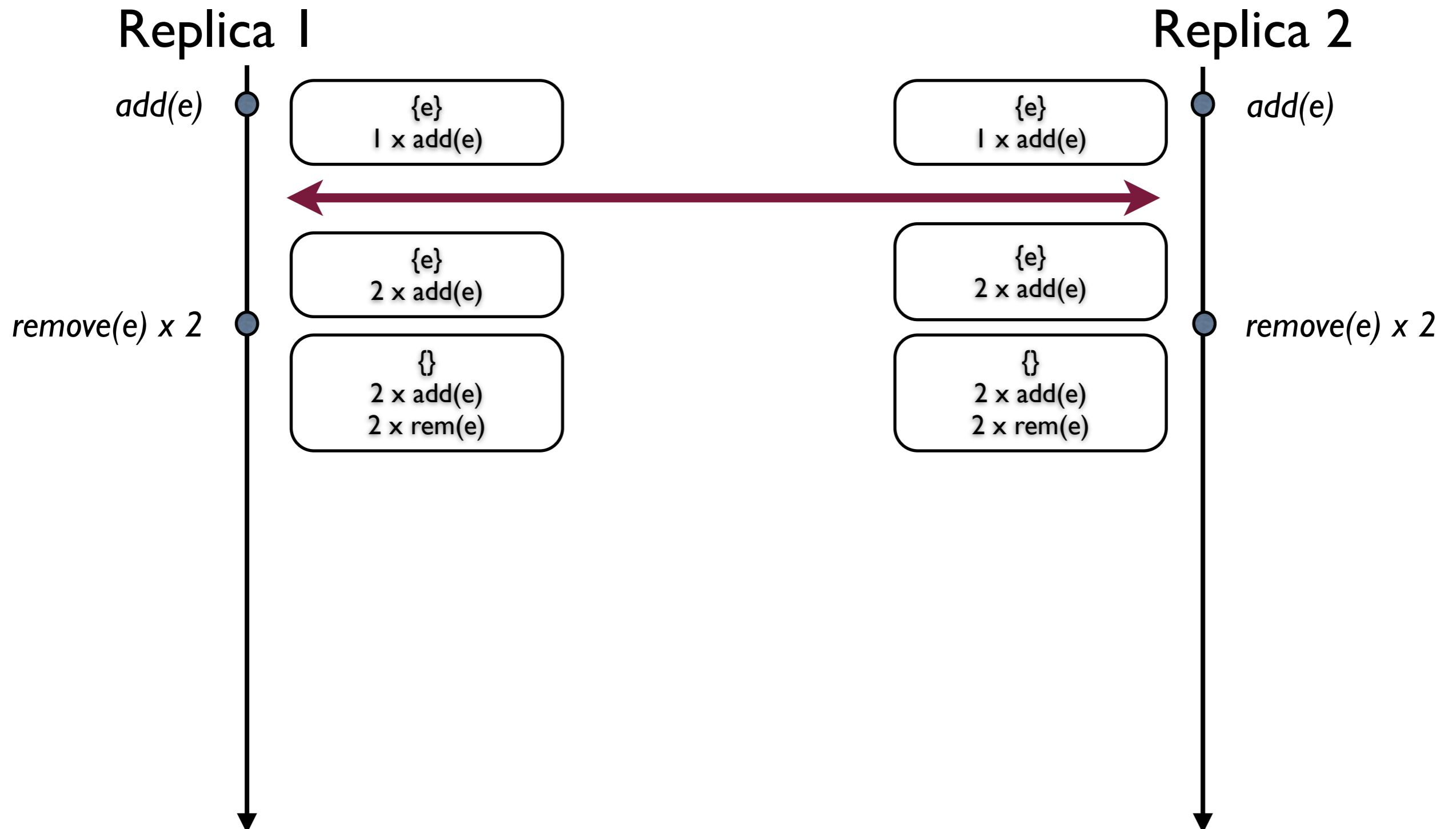
$add(e)$



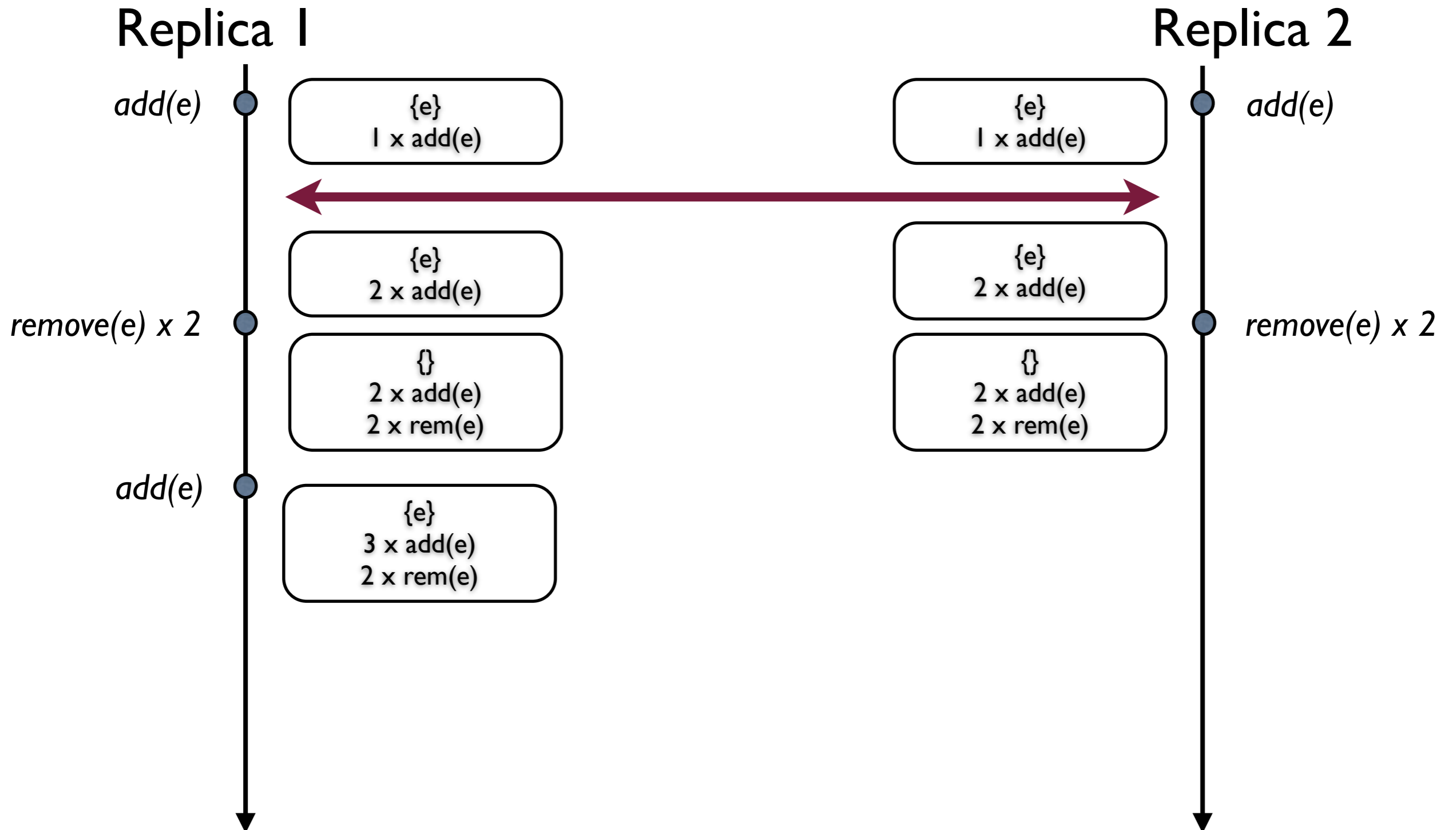
C-Set



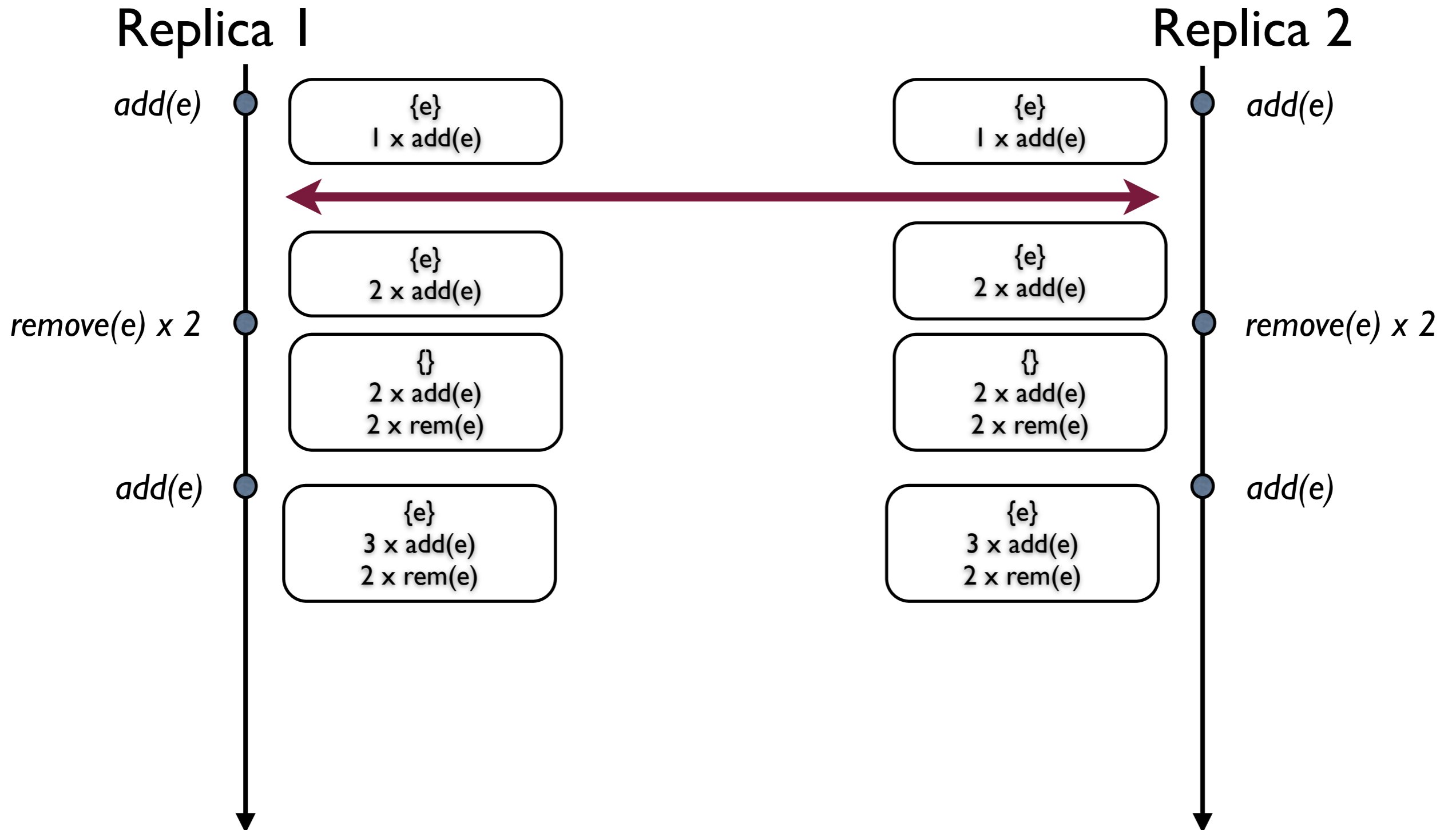
C-Set



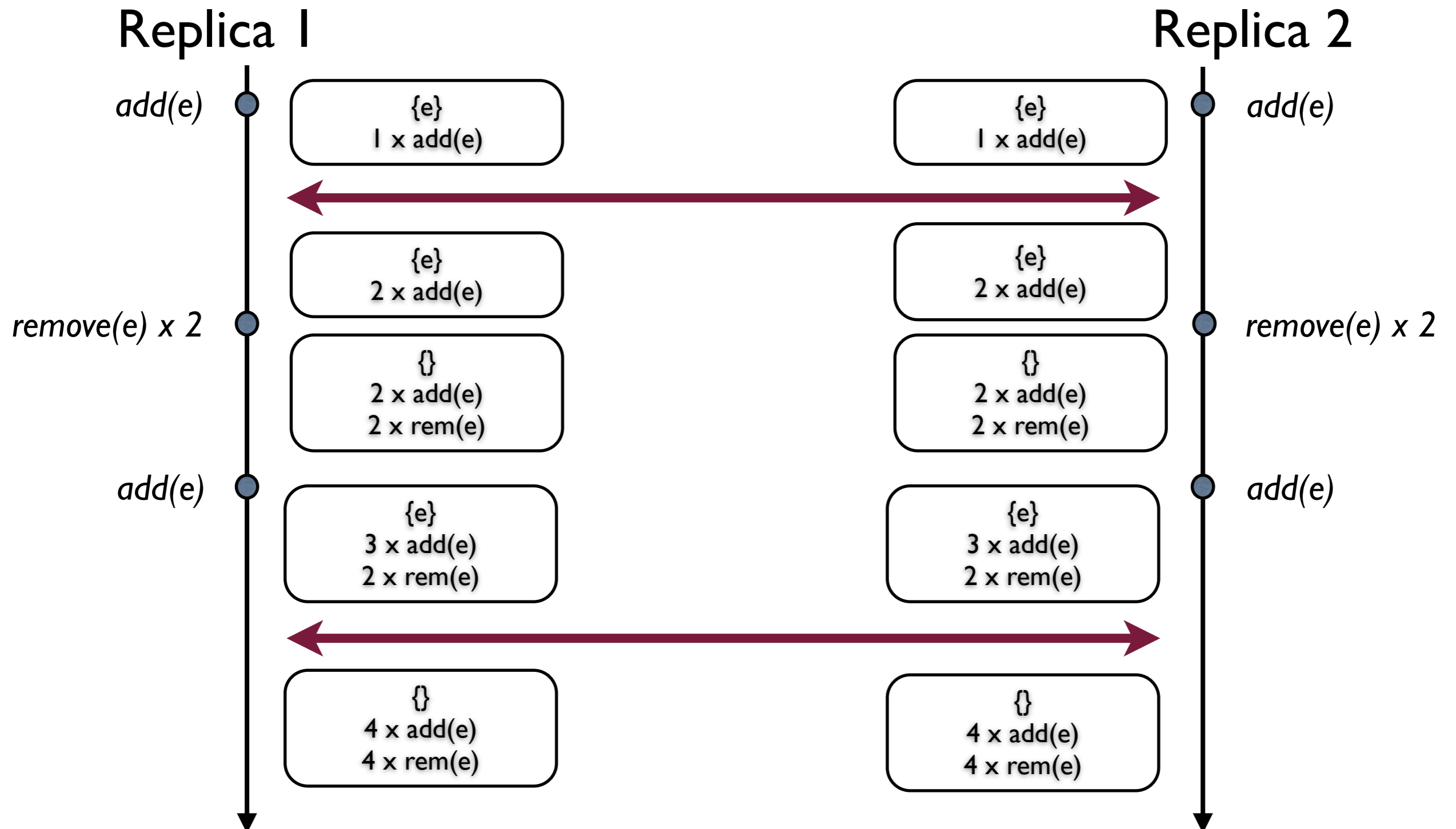
C-Set



C-Set



C-Set



Specification of concurrent operations

- For independent operations, concurrent execution should yield the same result as the sequential execution

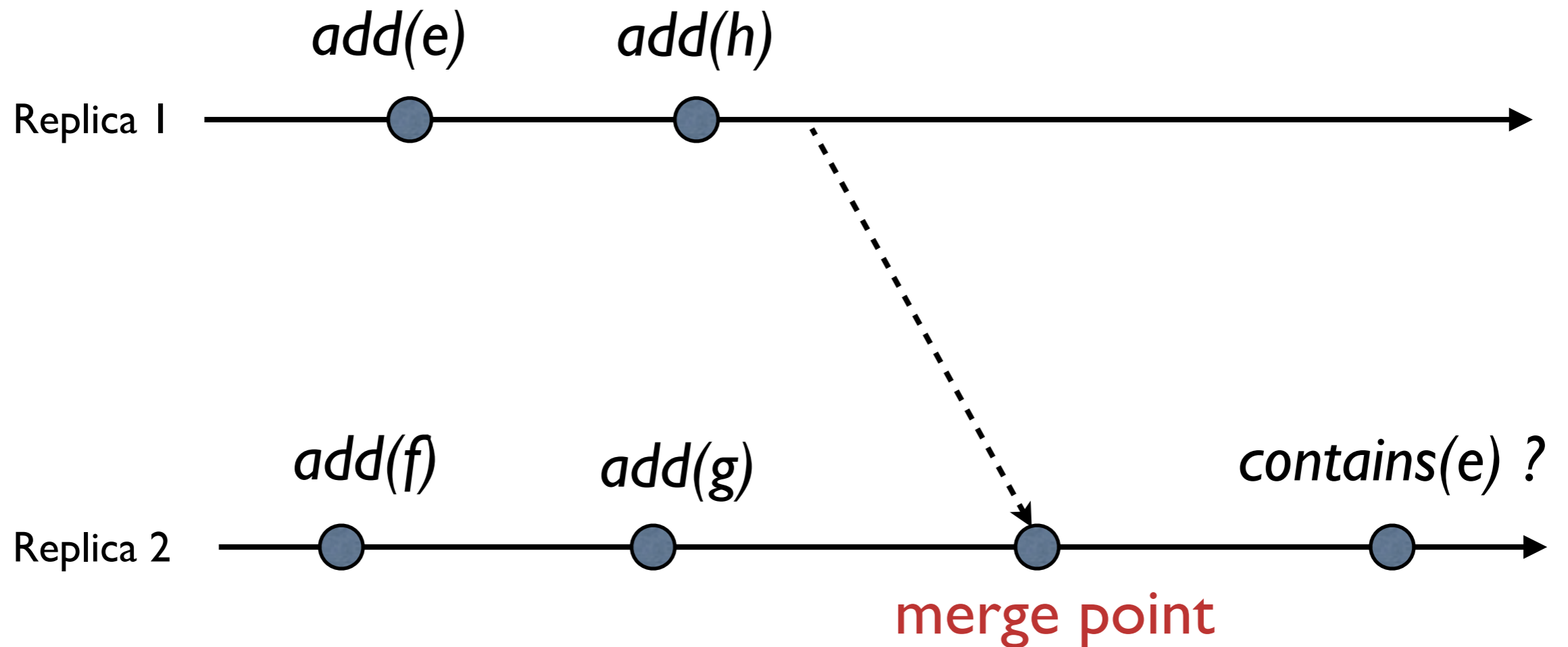
Specification of concurrent operations

- For independent operations, concurrent execution should yield the same result as the sequential execution

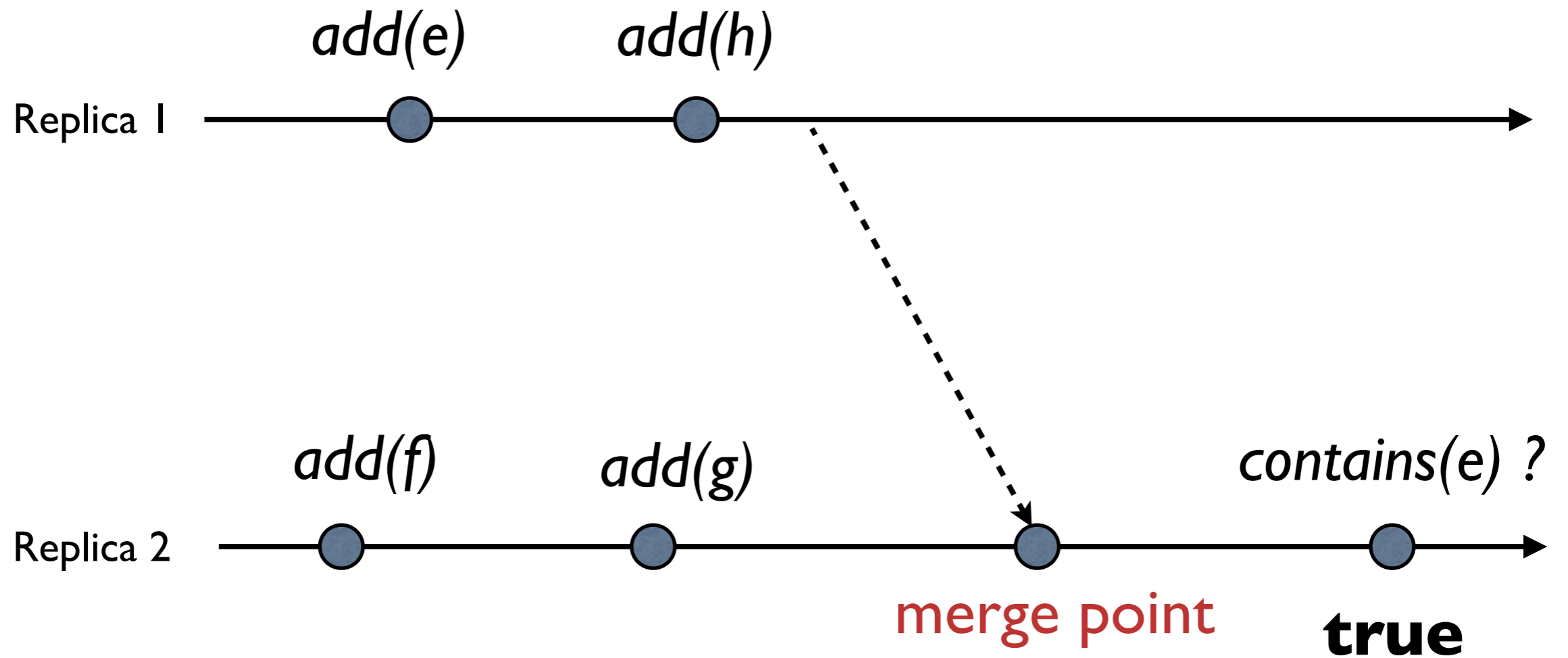
Principle of Permutation Equivalence

If $op;op'$ and $op';op$ have the same effect, then $op \parallel op'$ should also have the same effect.

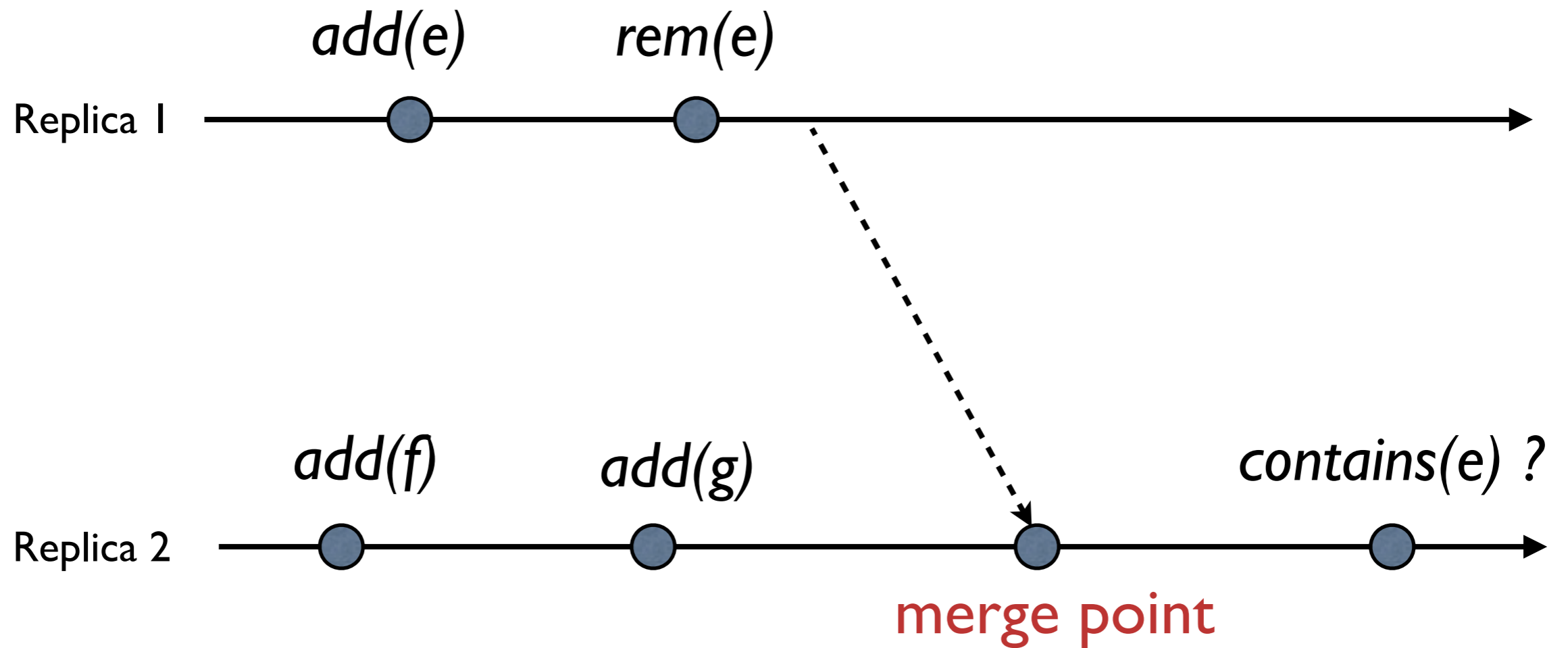
Applying the principle



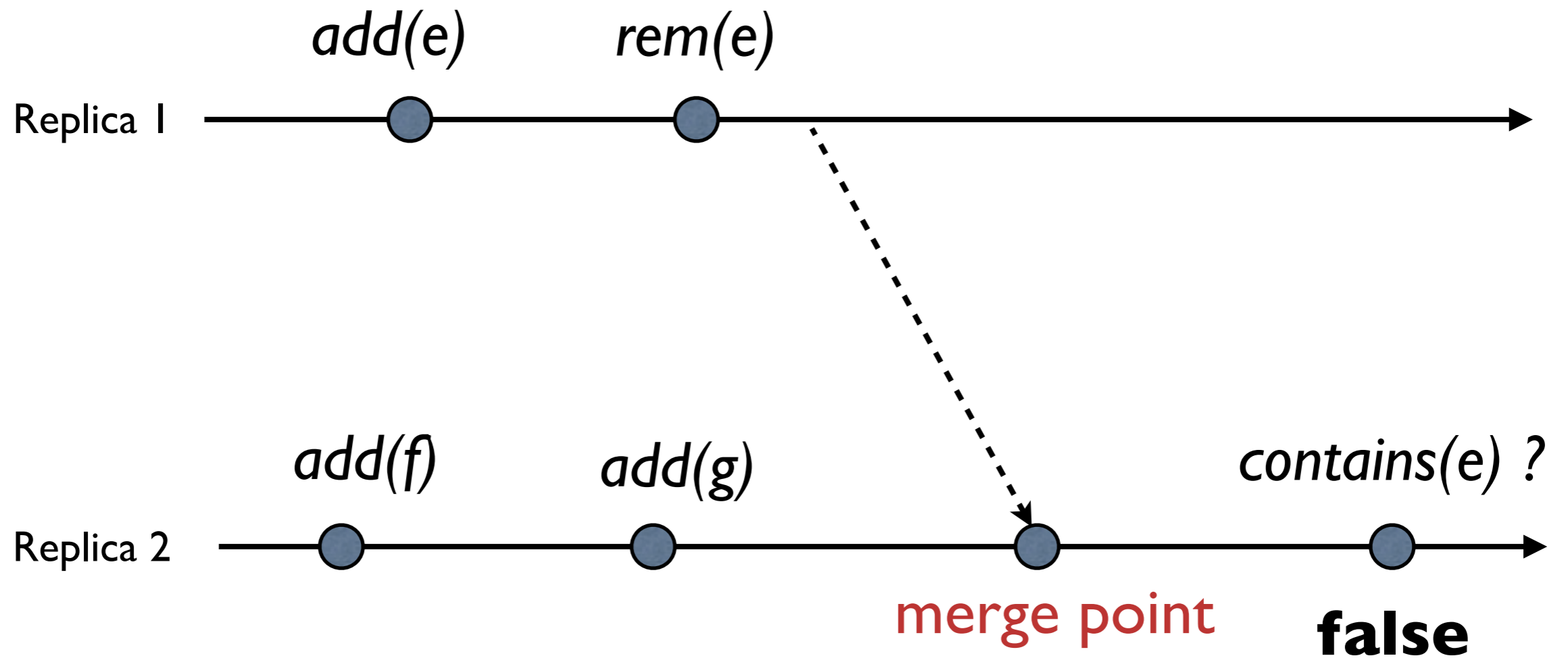
Applying the principle



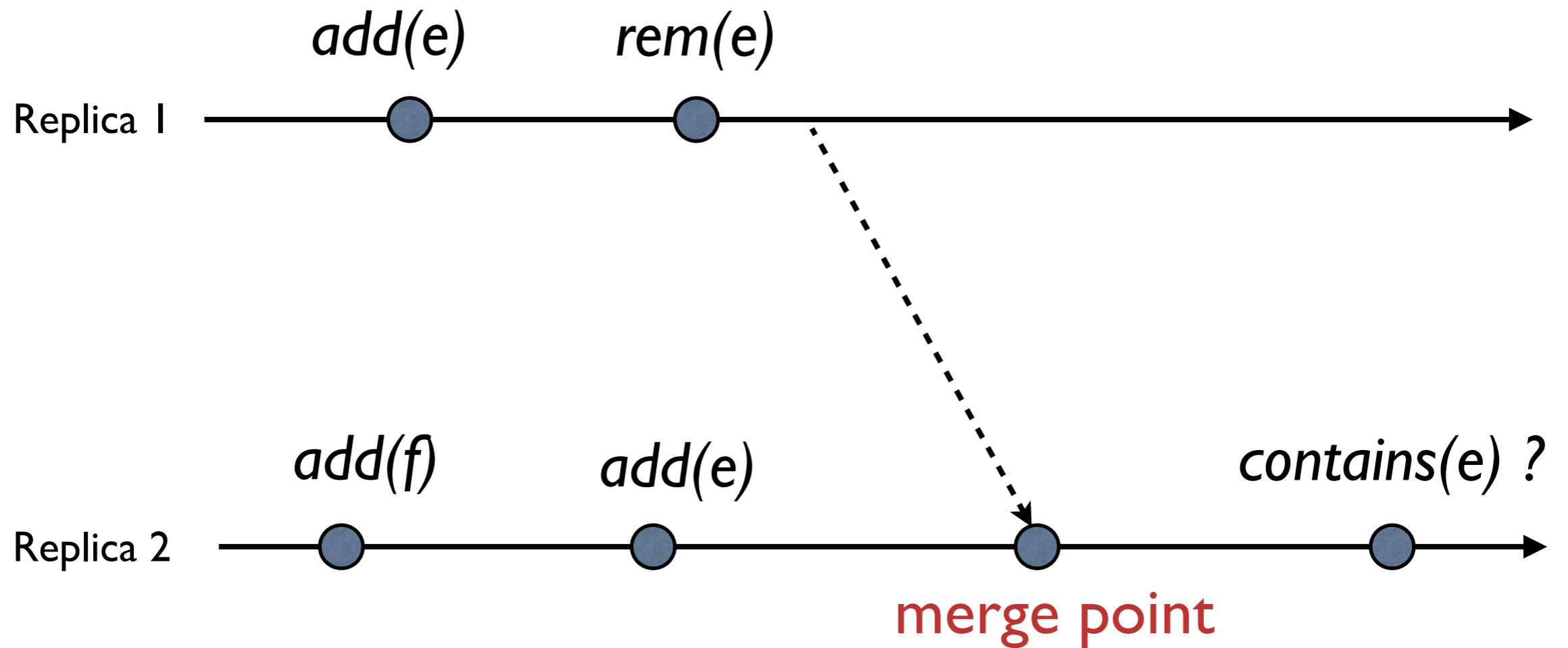
Applying the principle



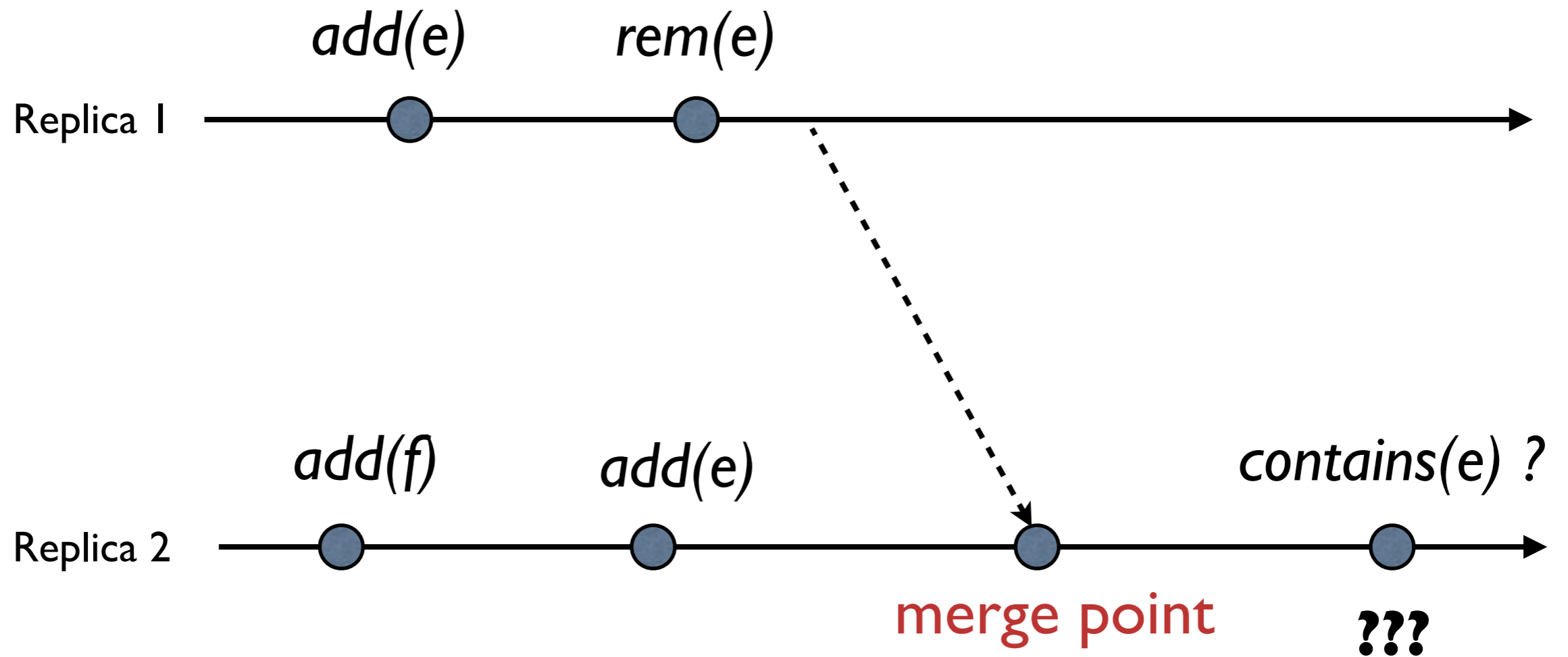
Applying the principle



Applying the principle



Applying the principle



Set semantics

The principle of permutation equivalence requires that

Set semantics

The principle of permutation equivalence requires that

$\{\text{true}\}$

$\text{add}(e) \parallel \text{add}(f)$

$\{e, f \in S\}$

Set semantics

The principle of permutation equivalence requires that

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(f)$ $\{e, f \in S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{rem}(f)$ $\{e \in S \wedge f \notin S\}$

Set semantics

The principle of permutation equivalence requires that

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(f)$ $\{e, f \in S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{rem}(f)$ $\{e \in S \wedge f \notin S\}$

$\{\text{true}\}$ $\text{rem}(e) \parallel \text{rem}(f)$ $\{e, f \notin S\}$

Set semantics

The principle of permutation equivalence requires that

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(f)$ $\{e, f \in S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{rem}(f)$ $\{e \in S \wedge f \notin S\}$

$\{\text{true}\}$ $\text{rem}(e) \parallel \text{rem}(f)$ $\{e, f \notin S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(e)$ $\{e \in S\}$

Set semantics

The principle of permutation equivalence requires that

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(f)$ $\{e, f \in S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{rem}(f)$ $\{e \in S \wedge f \notin S\}$

$\{\text{true}\}$ $\text{rem}(e) \parallel \text{rem}(f)$ $\{e, f \notin S\}$

$\{\text{true}\}$ $\text{add}(e) \parallel \text{add}(e)$ $\{e \in S\}$

$\{\text{true}\}$ $\text{rem}(e) \parallel \text{rem}(e)$ $\{e \notin S\}$

Semantics of $\text{add}(e) \parallel \text{rem}(e)$

Semantics of $\text{add}(e) \parallel \text{rem}(e)$

$\{\perp_e \in S\}$

Error indicator

Version control
systems

Semantics of $\text{add}(e) \parallel \text{rem}(e)$

$\{\perp_e \in S\}$

Error indicator

Version control
systems

$\{e \in S\}$

Add wins

Shopping cart

Semantics of $\text{add}(e) \parallel \text{rem}(e)$

$\{\perp_e \in S\}$

Error indicator

Version control
systems

$\{e \in S\}$

Add wins

Shopping cart

$\{e \notin S\}$

Remove wins

Shopping cart

Semantics of $\text{add}(e) \parallel \text{rem}(e)$

$\{\perp_e \in S\}$

Error indicator

Version control
systems

$\{e \in S\}$

Add wins

Shopping cart

$\{e \notin S\}$

Remove wins

Shopping cart

$\{\text{add}(e) > \text{rem}(e) \Rightarrow e \in S$
 $\wedge \text{rem}(e) > \text{add}(e) \Rightarrow e \notin S\}$

Last Writer wins

Read/write registers

OR-SET

payload set E, set T -- elements and tombstones: sets of pairs (element, id)

query contains (element e) : boolean b
let b = $(\exists n : (e, n) \in E)$

update **add** (element e)
prepare (e)
let n = unique() -- unique() returns a unique tag
effect (e, n)
 $E := E \cup \{(e, n)\} \setminus T$

update **remove** (element e)
prepare (e)
let R = $\{(e, n) \mid \exists n : (e, n) \in E\}$ -- collect all unique pairs containing e
effect (R)
 $E := E \setminus R$
 $T := T \cup R$ -- remove pairs observed at source

OR-SET

payload set E, set T

compare (ORSet B) : boolean b

let b = $(E \cup T) \subseteq (B.E \cup B.T) \wedge T \subseteq B.T$

merge (OrSet B)

$E := (E \setminus B.T) \cup (B.E \setminus T)$

$T := T \cup B.T$

- Specification supports mixing of state-based and operation-based updating!

Example

Replica 1



Replica 2



Example

Replica 1

$add(e)$



$\{e\}$
 $E = \{(e, n1)\}$
 $T = \{\}$



Replica 2



Example

Replica 1

$add(e)$



$\{e\}$
 $E = \{(e, n1)\}$
 $T = \{\}$



Replica 2

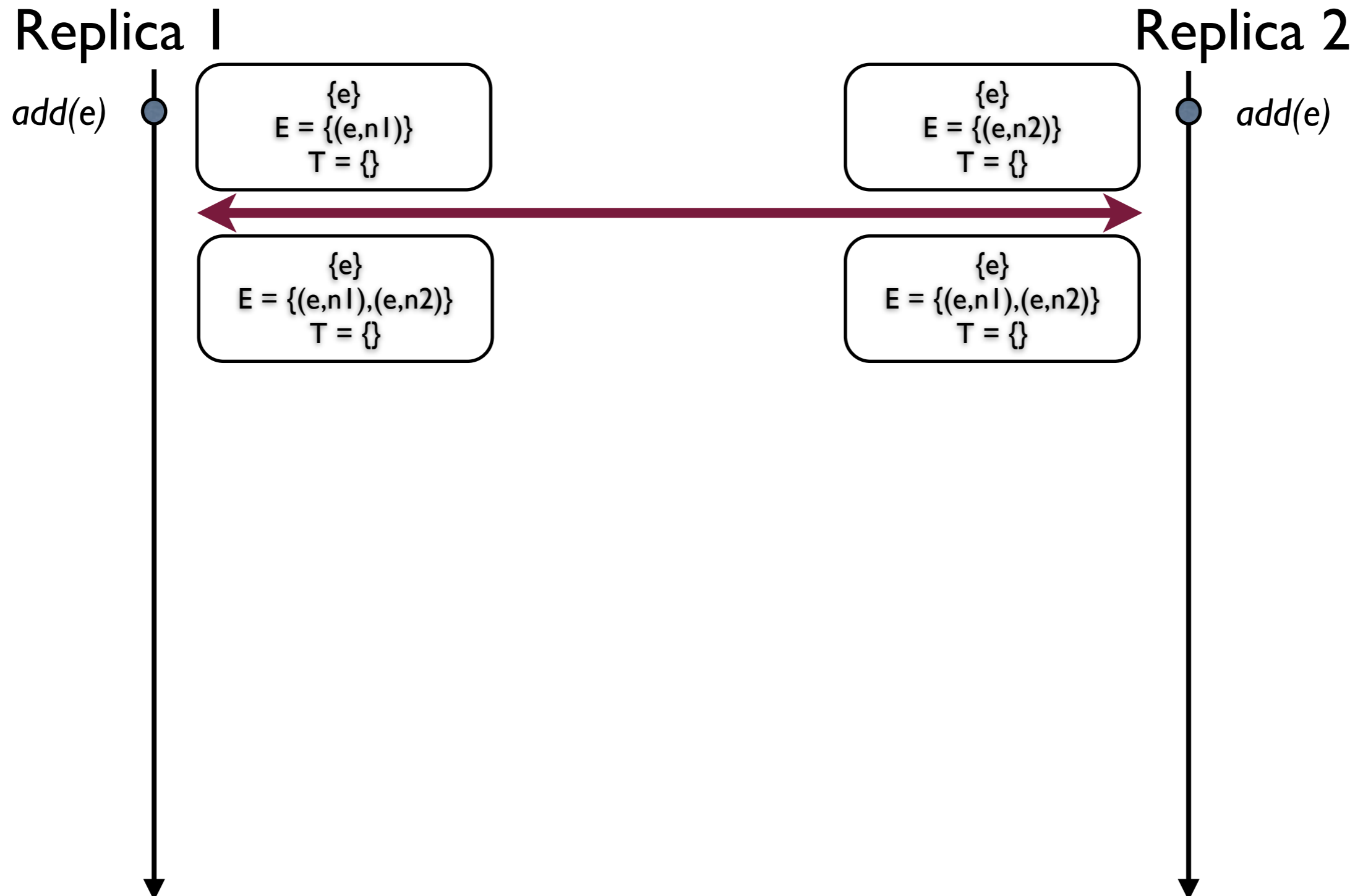
$add(e)$



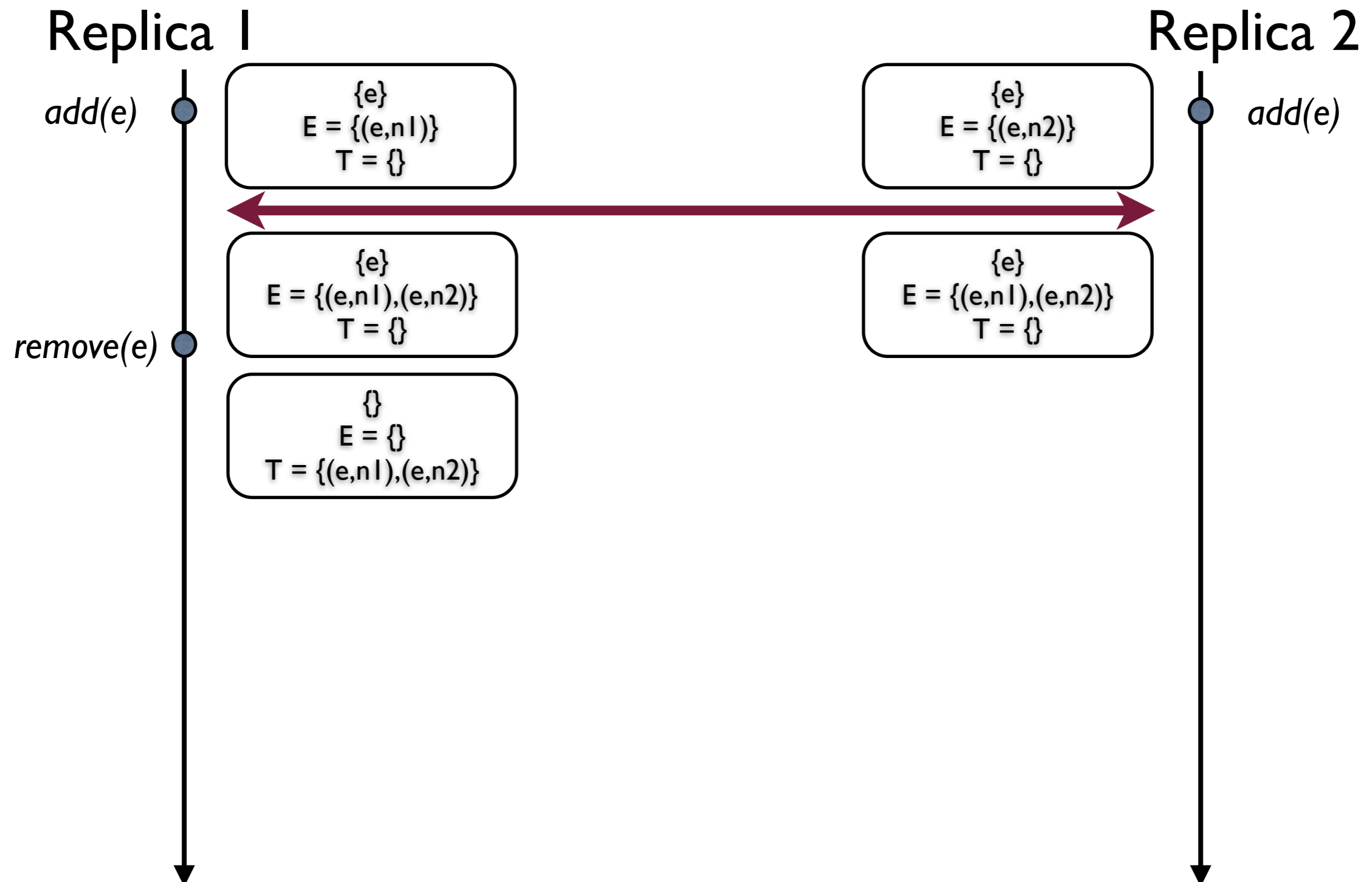
$\{e\}$
 $E = \{(e, n2)\}$
 $T = \{\}$



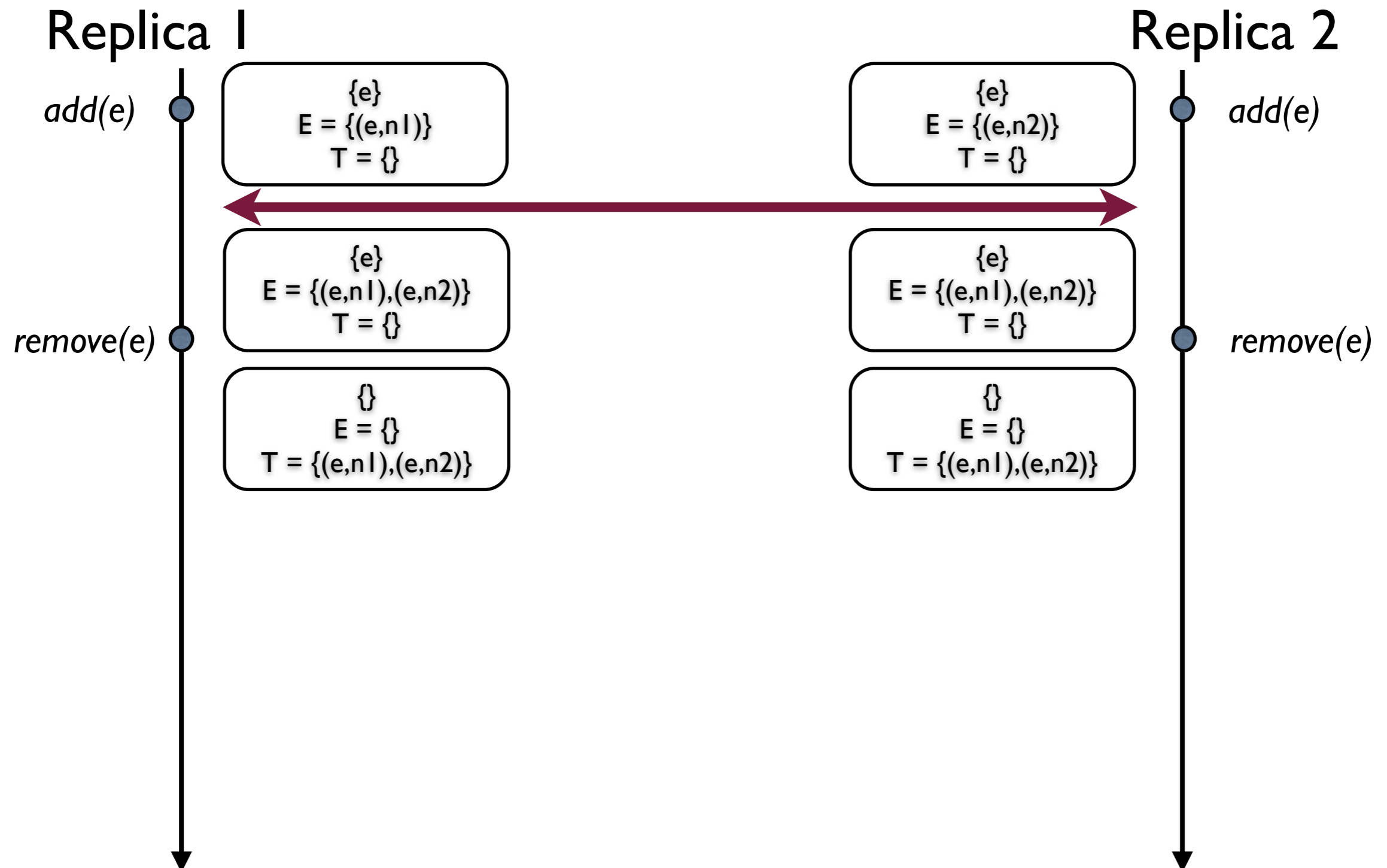
Example



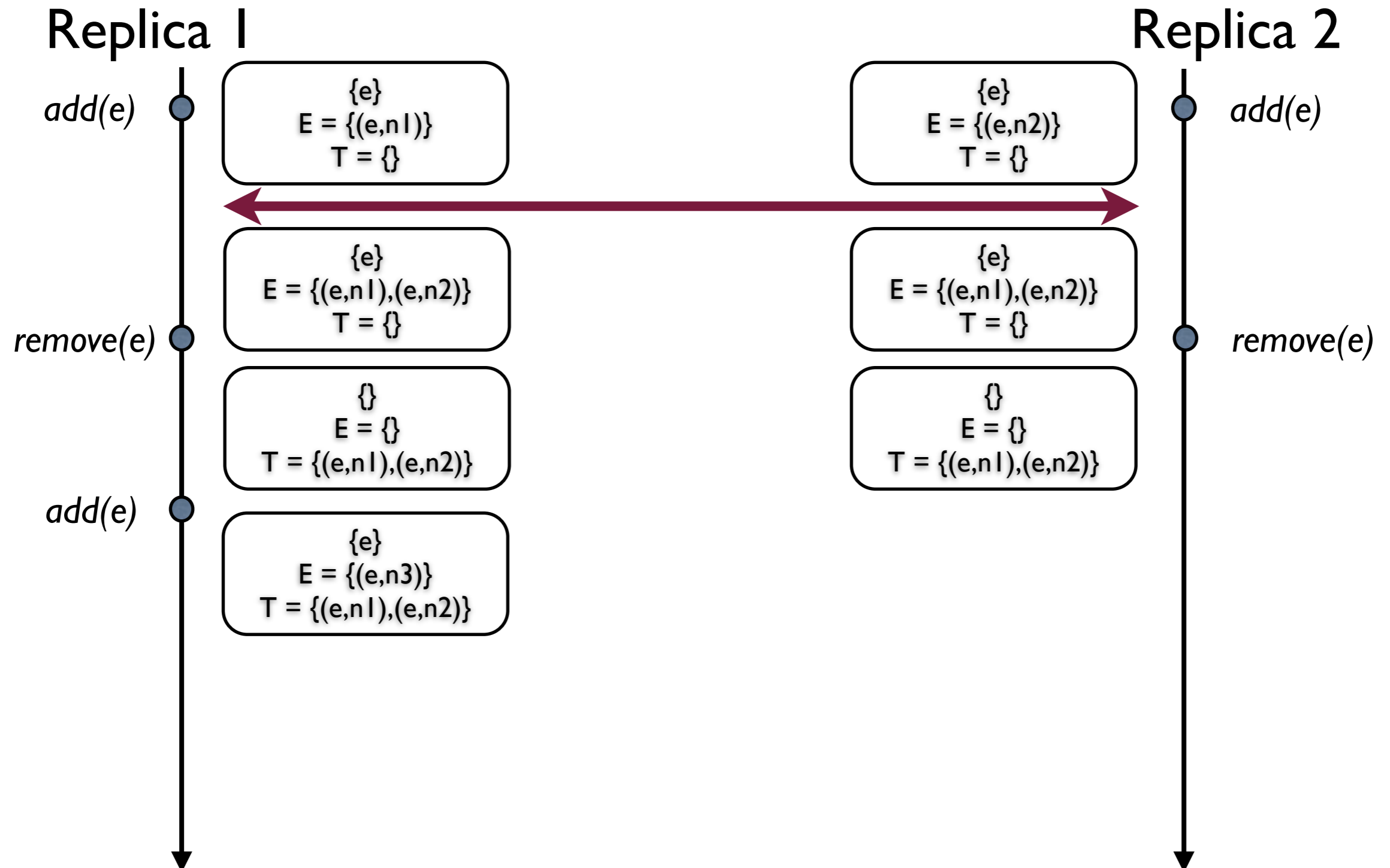
Example



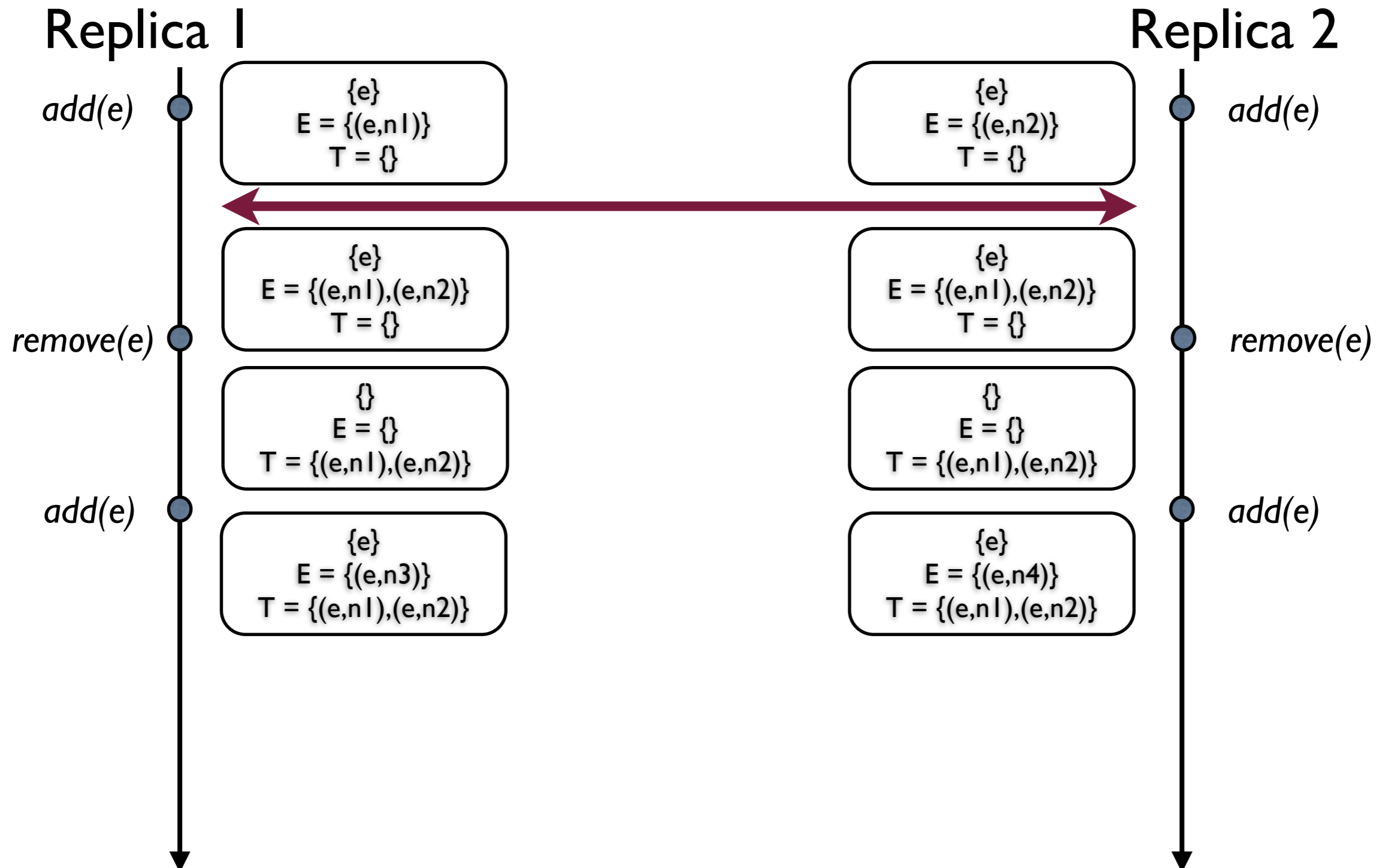
Example



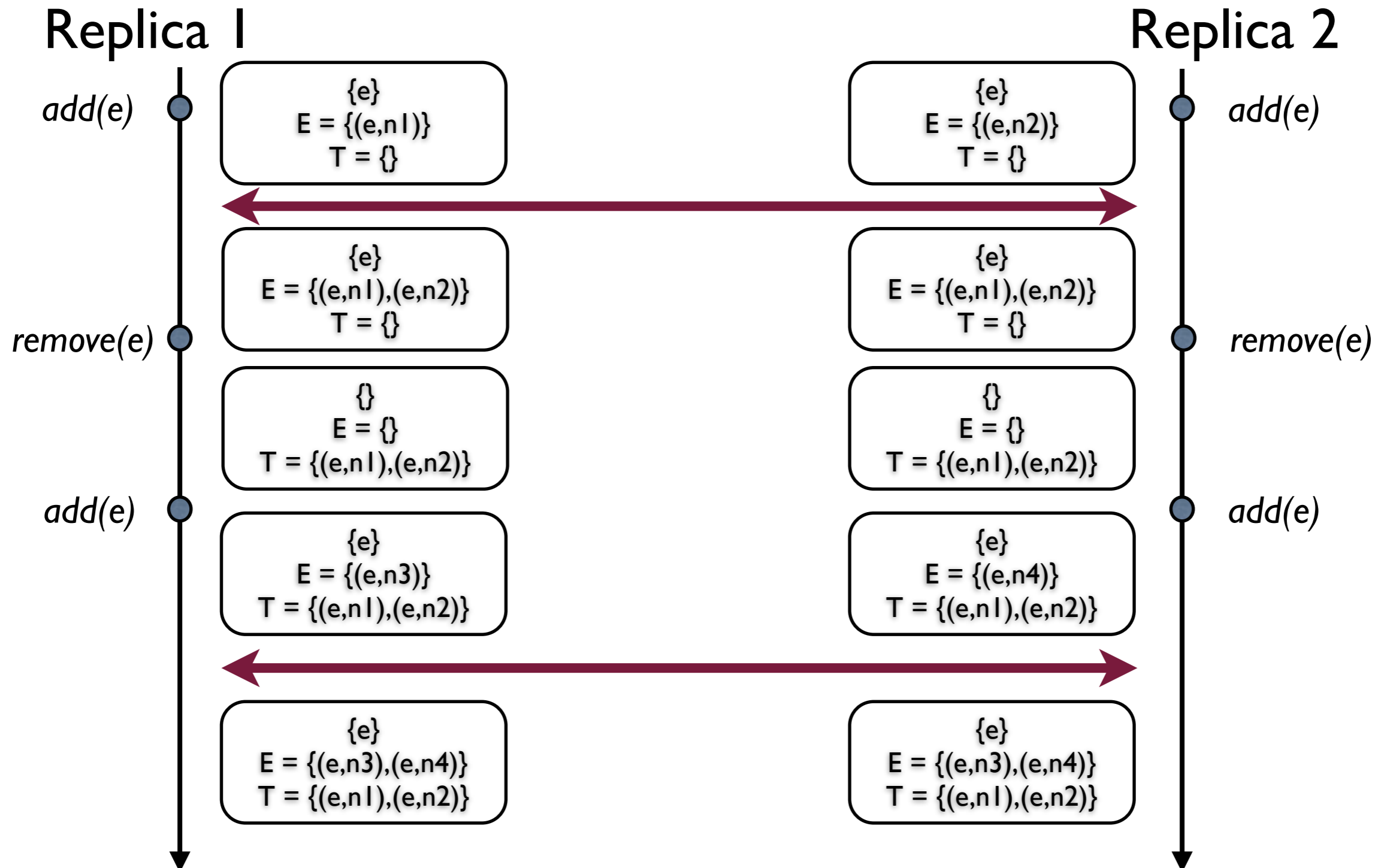
Example



Example



Example



Tombstone management

- Idea: Garbage-collect tombstones that have been delivered everywhere
 - ▶ Requires reliable membership service to deliver acknowledgements
 - ▶ Increases distributed processing
- Can we do better?

Tombstone management

- Insight: Adding an element (e,n) always happens-before removing (e,n)
 - ▶ For operation-based ORSet, no tombstones required!
 - ▶ To support merge, we need to record the happens-before information using version vectors (requires causal delivery)
- Trick can also be applied to other CRDTs

ORSWOT

payload set E , vector v
initial $\emptyset, [0, \dots, 0]$

-- E : elements, set of triples (element e , timestamp c , replica i)
-- v : vector (summary) of received process identifiers, indexed by replica

query contains (element e) : boolean b
let $b = (\exists c, i : (e, c, i) \in E)$

update add (element e)

prepare (e)

let $r = \text{myID}()$

-- r = source replica

let $c = v[r] + 1$

effect (e, c, r)

-- requires causal delivery

if $c > v[r]$ then

$v[r] := c$

$E := E \cup \{(e, c, r)\}$

update remove (element e)

prepare (e)

let $R = \{(e, c, i) \in E\}$

-- Collect all unique triples containing e

effect (R)

$E := E \setminus R$

ORSWOT

compare (A, B) : boolean b

let R = $\{(c,i) \mid 0 < c \leq A.v[i] \wedge \nexists e : (e,c,i) \in A.E\}$

let R' = $\{(c,i) \mid 0 < c \leq B.v[i] \wedge \nexists e : (e,c,i) \in B.E\}$

let b = $A.v \leq B.v \wedge R \subseteq R'$

-- Compare sets of removed ids

merge (B)

let M = $E \cap B.E$

let M' = $\{(e,c,i) \in E \setminus B.E \mid c > B.v[i]\}$

let M'' = $\{(e,c,i) \in B.E \setminus E \mid c > v[i]\}$

E := $M \cup M' \cup M''$

v := $[\max(v[0], B.v[0]), \dots, \max(v[n], B.v[n])]$

Preliminary evaluation

	<i>Write probability</i>					
	0	0,2	0,4	0,6	0,8	1
HashSet	1228	1251	1266	1239	1249	1235
C-Set	1228	980	849	818	697	672
ORSWOT	1241	1223	1167	1110	1074	1036

Throughput (K operations per second) on single machine

Summary

- Principle of Permutation Equivalence
- Concurrency semantics for set operations
- Several specifications for sets unifying state- and operation-based approach
- Efficient tombstone management
- Results are available as Technical Report

Remove-Wins Set

payload set E, set T (*initial* \emptyset, \emptyset) -- sets of pairs { (element e, unique-tag n), . . . }

query contains (element e) : boolean b
let b = $(\exists (e,n) \in E \wedge \nexists n' : (e,n') \in T)$

update add (element e)

prepare (e)

-- Collect all unique pairs of tombstones containing e

let $R' = \{(e, n) \mid \exists n : (e, n) \in T\}$

if $R' = \emptyset$ then $R = \{(e, \text{unique}())\}$

else $R = R'$

effect (R)

-- Remove pairs of tombstones observed at source

$E := E \cup R$

$T := T \setminus R$

update remove (element e)

prepare (e)

let $n = \text{unique}()$

-- unique() returns a unique tag

effect (e,n)

$T := T \cup \{(e,n)\} \setminus E$

Last-Writer-Wins Set

payload set E (*initial* \emptyset) -- sets of tuples { (element e , time t , boolean inSet), . . . }

query **contains** (element e) : boolean b
let $b = (\exists (e, t, \text{true}) \in E)$

update **add** (element e)
prepare (e)
let $t = \text{unique}()$ -- $\text{unique}()$ returns unique time
effect (e, t)
 $T := \{(e, t', v) \in E\}$
if $\forall (e, t', v) \in T : t' < t$ then $E := E \setminus T \cup \{(e, t, \text{true})\}$

update **remove** (element e)
prepare (e)
let $t = \text{unique}()$ -- $\text{unique}()$ returns unique time
effect (e, t)
 $T := \{(e, t', v) \in E\}$
if $\forall (e, t', v) \in T : t' < t$ then $E := E \setminus T \cup \{(e, t, \text{false})\}$