# Verifying Implementations of CRDTs

# Recommended Reading

INRIA
ROCQUENCOURT

## A comprehensive study of Convergent and Commutative Replicated Data Types [*]

Marc Shapiro, INRIA & LIP6, Paris, France
Nuno Preguiça, CITI, Universidade Nova de Lisboa, Portugal
Carlos Baquero, Universidade do Minho, Portugal
Marek Zawirski, INRIA & UPMC, Paris, France

**Abstract:**    Eventual consistency aims to ensure that replicas of some mutable shared object converge without foreground synchronisation. Previous approaches to eventual consistency are ad-hoc and error-prone. We study a principled approach: to base the design of shared data types on some simple formal conditions that are sufficient to guarantee eventual consistency. We call these types Convergent or Commutative Replicated Data Types either state based or op-... This paper formalises asynchronous object replication, ...condition appropriate for each case. It describes ...g both *add* and *remove* op-... s graphs, montonic DAGs, ...t non-trivial CRDTs. and sequences. ...

http://bit.ly/1PBC4zc

# Recommended Reading

## A comprehensive study of Convergent and Commutative Replicated Data Types [*]

Marc Shapiro, INRIA & LIP6, Paris, France
Nuno Preguiça, CITI, Universidade Nova de Lisboa, Portugal
Carlos Baquero, Universidade do Minho, Portugal
Marek Zawirski, INRIA & UPMC, Paris, France

Thème COM — Systèmes communicants

Projet Regal

**Abstract:** Eventual consistency aims to ensure that replicas of some mutable shared object converge without foreground synchronisation. Previous approaches to eventual consistency are ad-hoc and error-prone. We study a principled approach: to base the design of shared data types on some simple formal conditions that are sufficient to guarantee eventual consistency. We call these types Convergent or Commutative Replicated Data Types, either state based or optual consistency. This paper formalises asynchronous object replication, condition appropriate for each case. It describes both *add* and *remove* operations graphs, montonic DAGs, non-trivial CRDTs. and sequences.

**Key-words:** Data replication, optimistic replication, commutative operations

http://bit.ly/1PBC4zc

## 'Cause I'm Strong Enough: Reasoning about Consistency Choices in Distributed Systems

Alexey Gotsman
IMDEA Software Institute, Spain

Hongseok Yang
University of Oxford, UK

Carla Ferreira
NOVA LINCS, DI, FCT,
Universidade NOVA de Lisboa, Portugal

Mahsa Najafzadeh
Sorbonne Universités, Inria,
UPMC Univ Paris 06, France

Marc Shapiro
Sorbonne Universités, Inria,
UPMC Univ Paris 06, France

**Abstract**
Large-scale distributed systems often rely on replicated databases that allow a programmer to request different data consistency guarantees for different operations, and thereby control their performance. Using such databases is far from trivial: requesting stronger consistency in too many places may hurt performance, and requesting it in too few places may violate correctness. To help programmers in this task, we propose the first proof rule for establishing that a particular choice of consistency guarantees for various operations on a replicated database is enough to ensure the preservation of ations. Our rule is modular: it allows tool. We present a use on several examples.

use. Ideally, we would like replicated databases to provide *strong consistency*, i.e., to behave as if a single centralised node handles all operations. However, achieving this ideal usually requires synchronisation among replicas, which slows down the database and even makes it unavailable if network connections between replicas fail [2, 24].

For this reason, modern replicated databases often eschew synchronisation completely; such databases are commonly dubbed *eventually consistent* [47]. In these databases, a replica performs an operation requested by a client locally without any synchronisation with other replicas and immediately returns to the client; the *effect* of t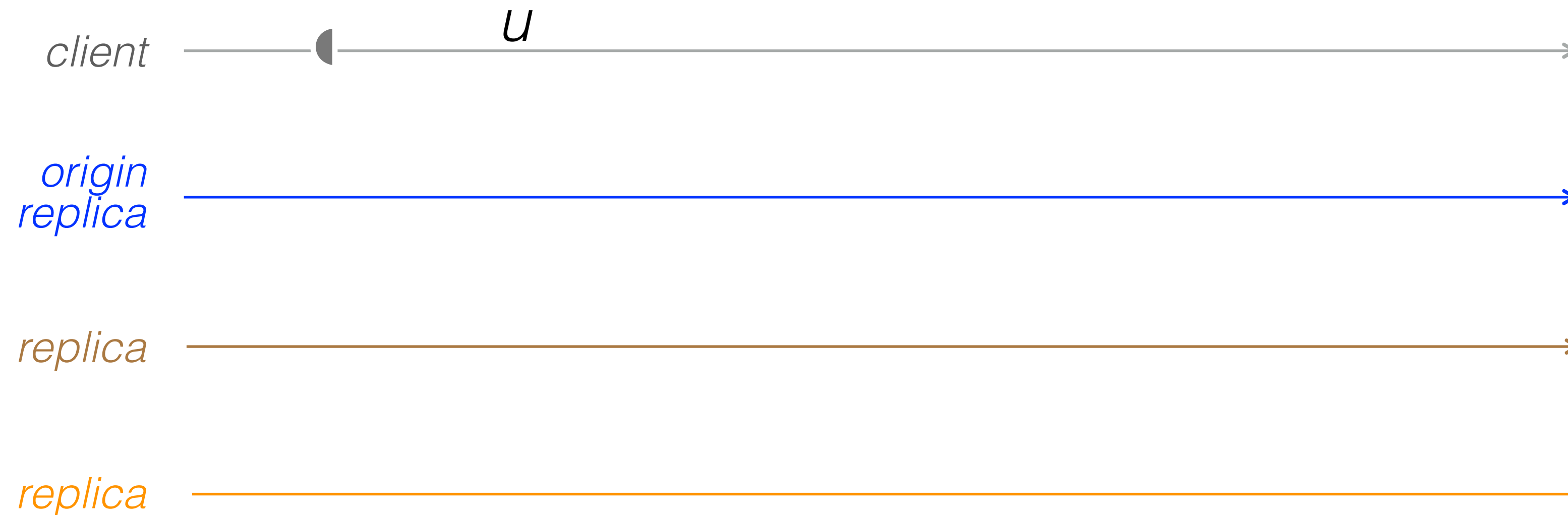he operation is propagated to the other replicas only *even* lead to *anomalies*—behaviours deviating from. One of them is illustrated in Figure 1(a). Here t while connected to a replica $r_1$, and Bob, also sees the post and comments on it. After each of s, $r_1$ sends a message to the other replicas in the system with the update performed by the user. If the messages with the updates by Alice and Bob arrive to a replica $r_2$ out of order, sees Bob's comment,

http://bit.ly/2nM96mT

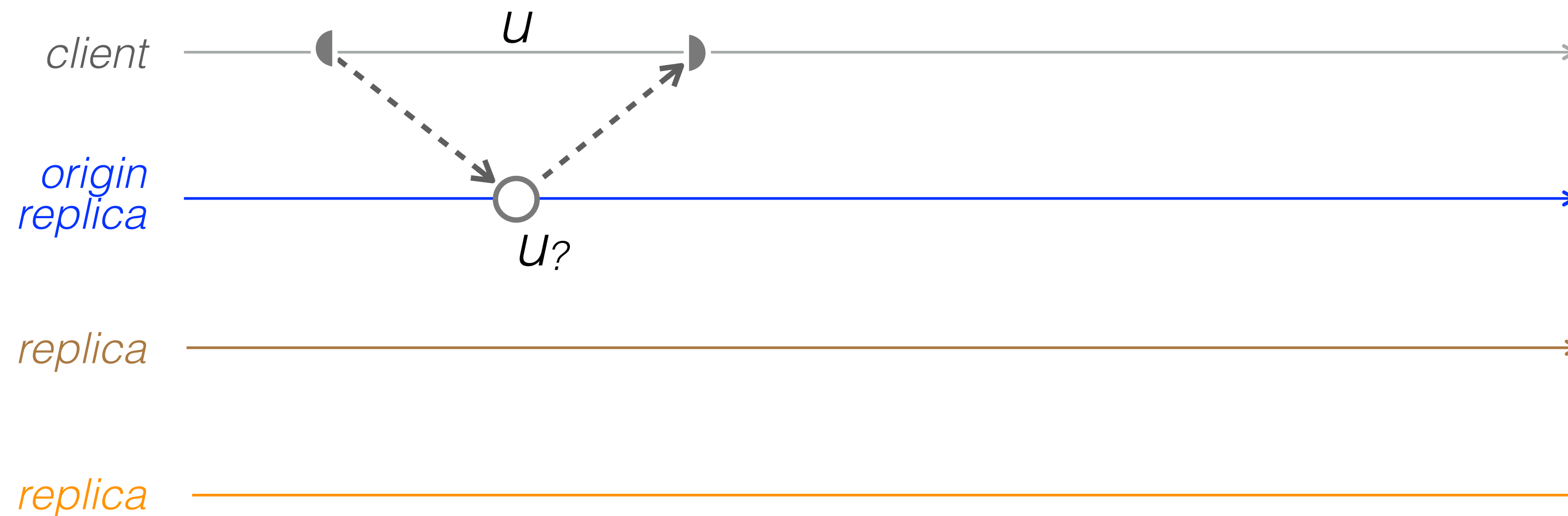<u>Disclaimer</u>:

Slides kindly provided by Marc Shapiro

(all errors are mine)

# Operation



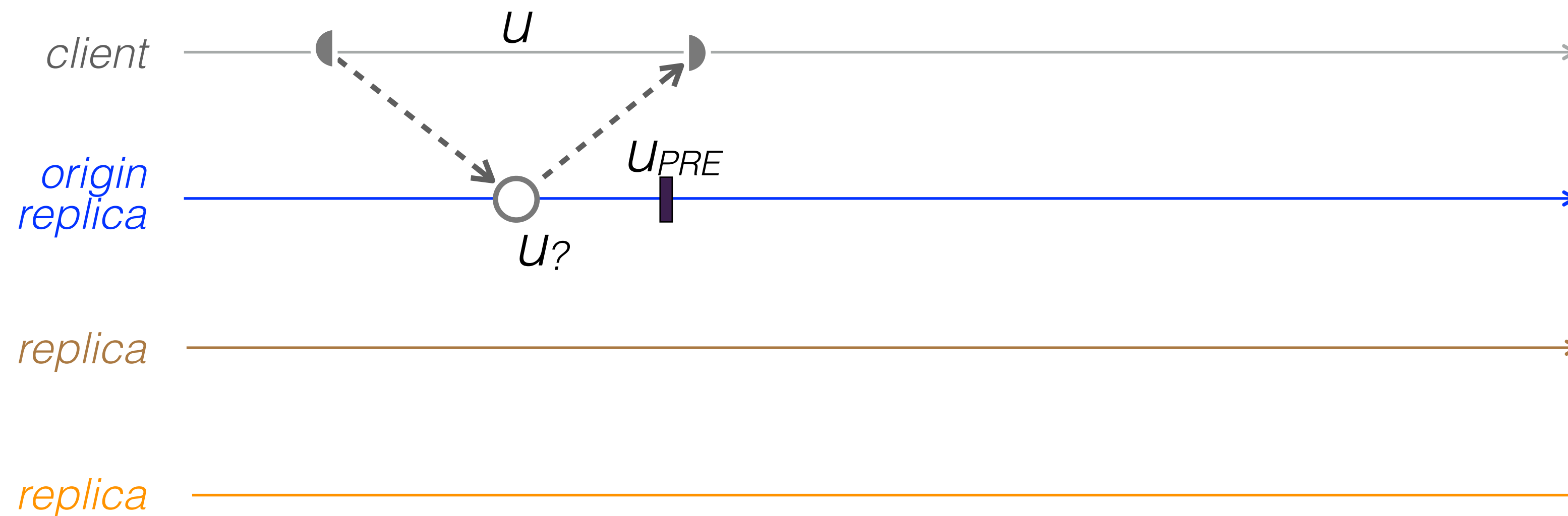- *u: state $\hookrightarrow$ (retval, (state $\hookrightarrow$ state))*
- Prepare (@origin) $u_?$; deliver $u_!$
- Read one, write all (ROWA)
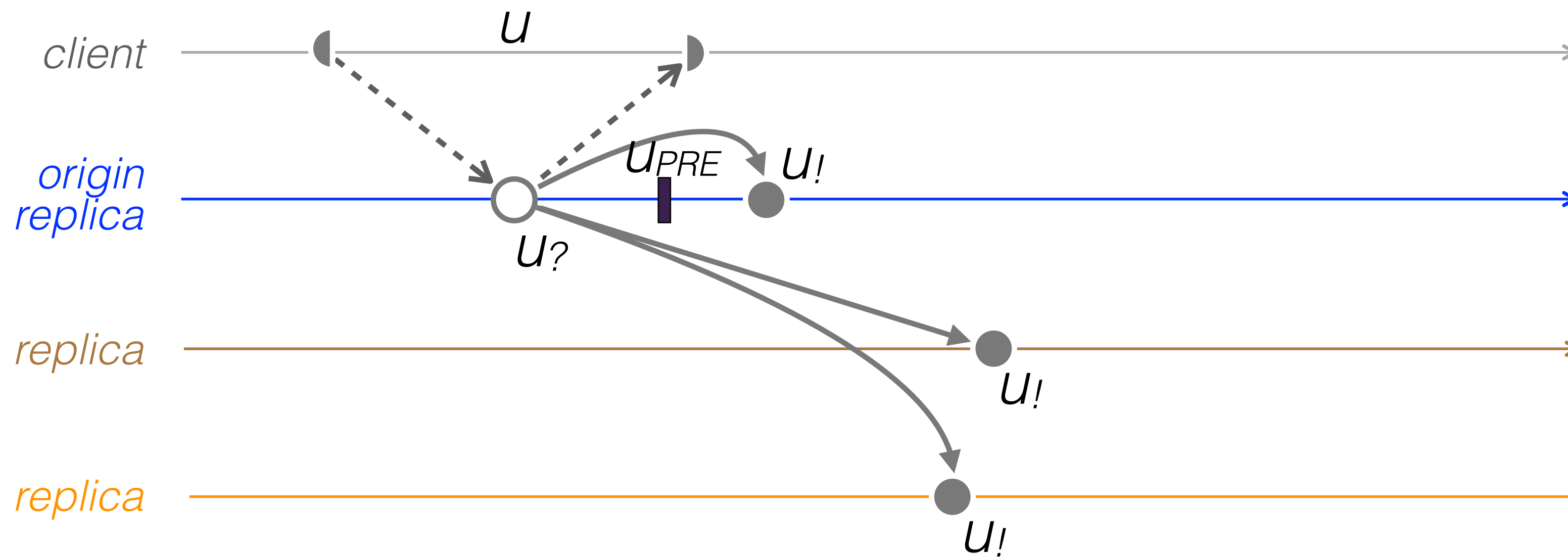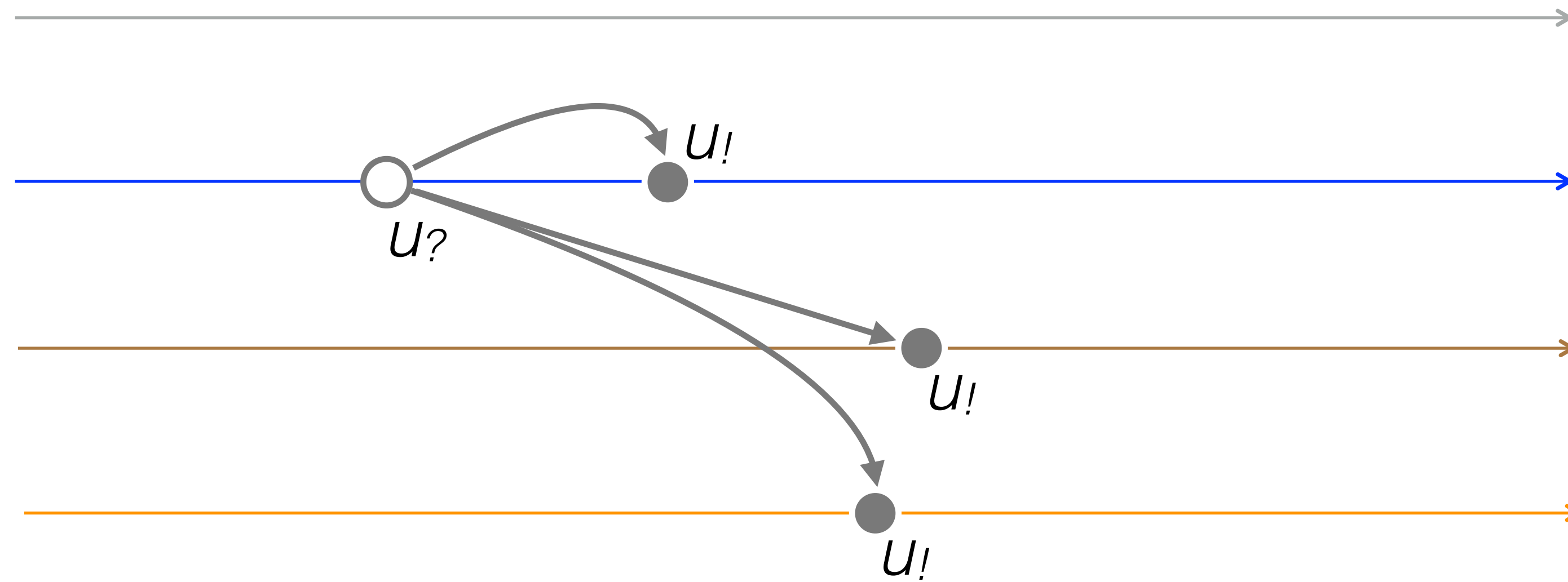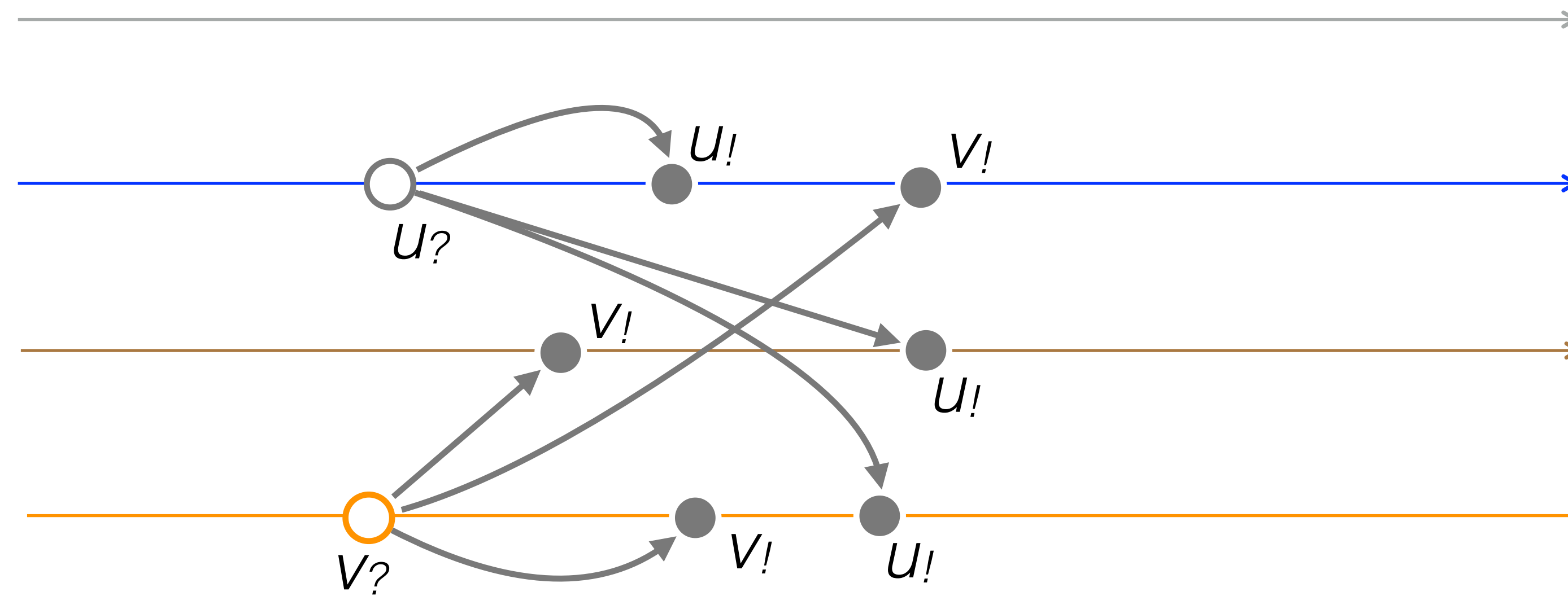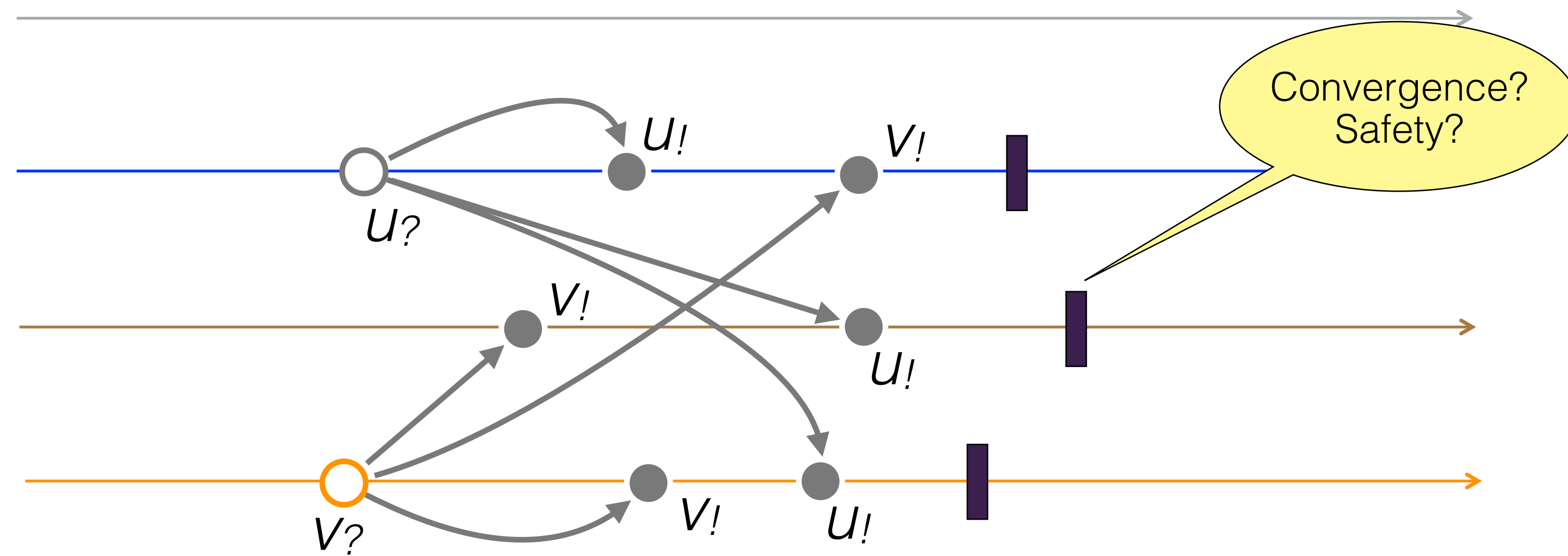- Deferred-update replication (DUR)

# Operation



- *u: state ↪ (retval, (state ↪ state))*
- Prepare (@origin) *u?*; deliver *u!*
- Read one, write all (ROWA)
- Deferred-update replication (DUR)

# Operation



- *u: state $\hookrightarrow$ (retval, (state $\hookrightarrow$ state))*
- Prepare (@origin) $u_?$; deliver $u_!$
- Read one, write all (ROWA)
- Deferred-update replication (DUR)

# Operation



- ▸ *u: state ↪ (retval, (state ↪ state))*
- ▸ Prepare (@origin) *u$_?$*; deliver *u$_!$*
- ▸ Read one, write all (ROWA)
- ▸ Deferred-update replication (DUR)

# Concurrency



▶ Concurrent, Multi-master

▶ Strong: total order, identical state

▶ Weak: concurrent, interleaving, no global state

# Concurrency



▸ Concurrent, Multi-master

▸ Strong: total order, identical state

▸ Weak: concurrent, interleaving, no global state

# Concurrency



- ▶ Concurrent, Multi-master
- ▶ Strong: total order, identical state
- ▶ Weak: concurrent, interleaving, no global state

# Anomalies of concurrent updates

- Bank:
  - $\sigma_{init}$ = 100€
  - Alice: *credit(20)* = { σ ≔ 120 }
  - Bob: *debit (60)* = { σ ≔ 40 }
  - σ = ???

# Anomalies of concurrent updates

- Bank:
  - $\sigma_{init} = 100€$
  - Alice: *credit(20)* = { σ ≔ 120 }
  - Bob: *debit (60)* = { σ ≔ 40 }
  - σ = ???

- File system:
  - $\sigma_{init} = $ "/"
  - Alice: *mkdir ("/foo"); mkdir ("/foo/bar")*
  - Bob: receives *mkdir ("/foo/bar")*
  - σ = ???

# Eventual Consistency

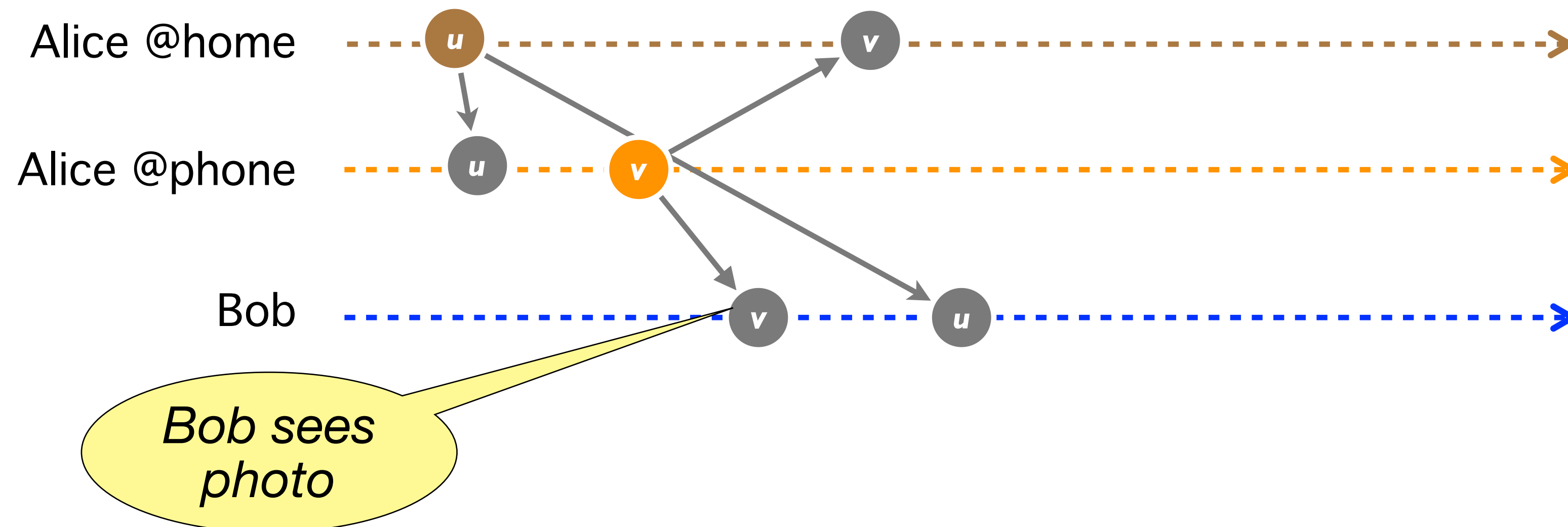Don't show photos to Bob

Alice @home    *u*

Alice @phone

Bob

- *access (Bob, photo)* $\implies$ *ACL (Bob, photo)*
- *v observed effects of u $\implies$ v should be delivered after u*
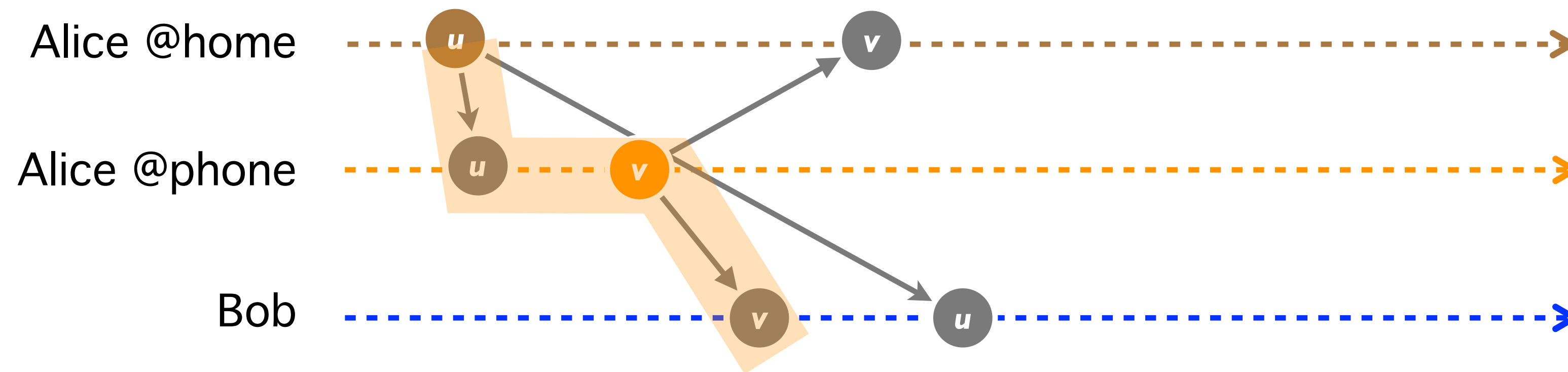- *Available: doesn't slow down sender*

# Eventual Consistency



- *access (Bob, photo) $\implies$ ACL (Bob, photo)*

- *v observed effects of u $\implies$ v should be delivered after u*

- *Available: doesn't slow down sender*
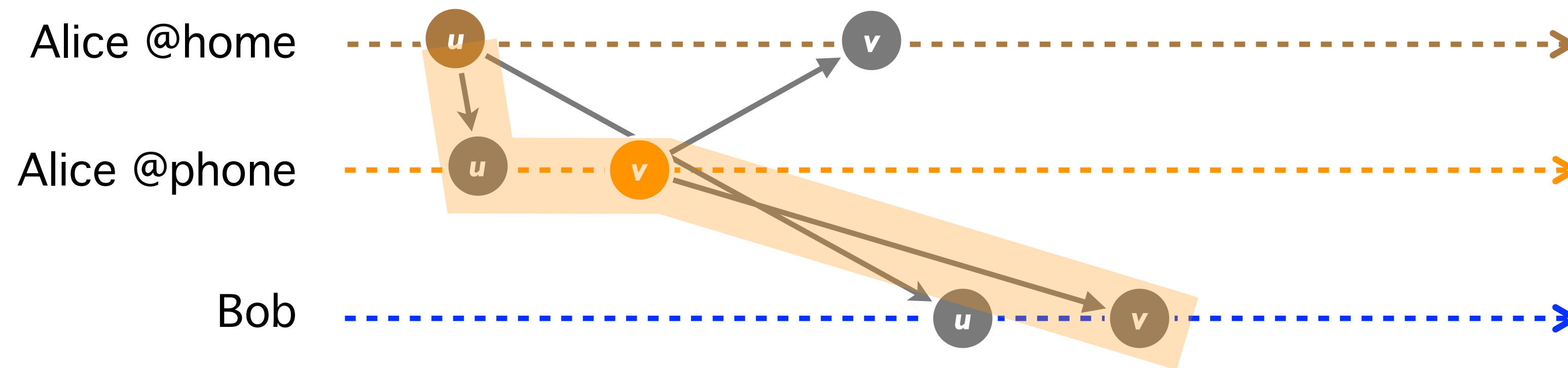
# Eventual Consistency



- *access (Bob, photo) $\Longrightarrow$ ACL (Bob, photo)*

- *v observed effects of $u \Longrightarrow v$ should be delivered after $u$*

- *Available: doesn't slow down sender*

# Eventual Consistency

Alice @home

Alice @phone

Bob

*post photo*

- *access (Bob, photo) $\implies$ ACL (Bob, photo)*
- *v observed effects of $u \implies v$ should be delivered after u*
- *Available: doesn't slow down sender*

# Eventual Consistency



- *access (Bob, photo) $\implies$ ACL (Bob, photo)*

- *v observed effects of u $\implies$ v should be delivered after u*

- *Available: doesn't slow down sender*

# Eventual Consistency



▶ *access (Bob, photo)* $\implies$ *ACL (Bob, photo)*

▶ *v observed effects of u $\implies$ v should be delivered after u*

▶ *Available: doesn't slow down sender*

# (1) Causal consistency



- *access (Bob, photo) $\implies$ ACL (Bob, photo)*

- *v observed effects of u $\implies$ v should be delivered after u*

- *Available: doesn't slow down sender*

# (2) Conflict-free Replicated Data Types (CRDTs)

- ▸ Data type
  - ▸ Encapsulates state
- ▸ Replicated
  - ▸ At multiple nodes
- ▸ Available
  - ▸ Update my replica without coordination
  - ▸ Convergence guaranteed by design
  - ▸ Decentralized, peer-to-peer

# Commute $\implies$ Converge

- Bank account:
  - **credit(amt)**$_l$ = { *local_balance += amt* }
  - **debit(amt)**$_l$ = { *local_balance −= amt* }
  - **interest()**$_l$ = { *local_balance += origin_balance*.05* }
- File system:
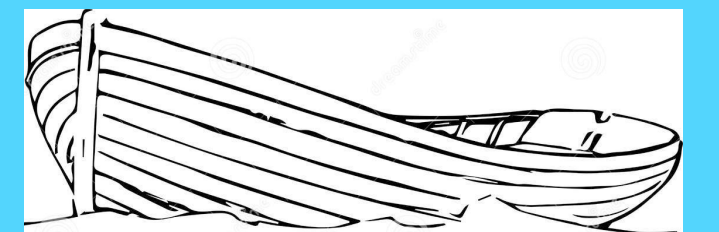  - *write(f)*$_l$ = { *local_f ⊔ f* }

# CRDT design concept

▸ Backward-compatible with sequential datatype

▸ Commute $\implies$ concurrent is same

    ▸ *add(e); rm(f) = rm(f); add(e)* $\triangleq$ *add(e) || rm (f)*

▸ Otherwise, *concurrency semantics*

    ▸ Example: *add(e) || rm (e)*

    ▸ Deterministic, similar to sequential

        ▸ $\approx$ *rm(e);add(e)* or $\approx$ *add(e); rm(e)*

    ▸ Merge, don't lose updates

    ▸ Result doesn't depend on order received

    ▸ Stable preconditions

# CRDT design concept

▸ Backward-compatible with sequential datatype

▸ Commute $\implies$ concurrent is same

  ▸ *add(e); rm(f) = rm(f); add(e) $\triangleq$ add(e) || rm (f)*

▸ Otherwise, *concurrency semantics*

▸ Example: *add(e) || rm (e)*

▸ Deterministic, similar to sequential

  ▸ *$\approx$ rm(e);add(e)* or *$\approx$ add(e); rm(e)*

▸ Merge, don't lose updates

▸ Result doesn't depend on order received

▸ Stable preconditions

# CRDT design concept

▸ Backward-compatible with sequential datatype

▸ Commute $\implies$ concurrent is same

  ▸ *add(e); rm(f) = rm(f); add(e) $\triangleq$ add(e) || rm (f)*

▸ Otherwise, *concurrency semantics*

  ▸ Example: *add(e) || rm (e)*

  ▸ Deterministic, similar to sequential

    ▸ $\approx$ *rm(e);add(e)* or $\approx$ *add(e); rm(e)*

  ▸ Merge, don't lose updates

  ▸ Result doesn't depend on order received

  ▸ Stable preconditions

# Application invariants

- *South ⊎ Boat ⊎ North = { sheep, dog, wolf }*

- *carryNorth(S) $\implies$ 1 ≤ |S| ≤ 2*

- *carrySouth(S) $\implies$ 1 ≤ |S| ≤ 2*

- *∀S ∈ {South, Boat, North} : sheep ∈ S ∧ wolf ∈ S $\implies$ dog ∈ S*

- *Hard to tease invariants out*
  - Silent invariants

# Seq. consistency examples

- Bank account
  - *deposit(amt), withdraw(amt), accrueInterest(amt)*
  - Invariant: *"balance ≥ 0"*
  - *{ amt ≤ balance ∧ Inv } withdraw(amt) { Inv }*

# Seq. consistency examples

- ▸ Bank account
  - ▸ *deposit(amt), withdraw(amt), accrueInterest(amt)*
  - ▸ Invariant: *"balance ≥ 0"*
  - ▸ { *amt ≤ balance ∧ Inv* } *withdraw(amt)* { *Inv* }

- ▸ File system
  - ▸ *mkdir, rmdir, create, write, rm, ls,* etc.
  - ▸ Invariant: Tree
  - ▸ { Tree ∧ ¬ *x/…/y* } *mv(x,y)* { Tree }

# Just-Right Consistency

- ▸ CRDT geo-replicated database
  - ▸ Lots of internal parallelism
  - ▸ Transactional, causal consistency by default
- ▸ Specification of application updates, invariant
  - ▸ **CISE: do all state transitions preserve the invariant?**
  - ▸ If not, fix: adjust
    - ▸ either specification
    - ▸ or synchronisation
  - ▸ Repeat until safe
- ▸ App / synch co-design: Minimal synchronisation

$\sigma{:}I$ ————————————————————→

$\sigma{:}I$ ————————————————————→

Asynchronous, replicated updates
  ‣ State $\sigma$
  ‣ Invariant $I$
  ‣ Prepare: read one, generate effector
  ‣ Update all, deferred: deliver effector
Converge?  Invariant OK?

Asynchronous, replicated updates

‣ State $\sigma$

‣ Invariant $I$

‣ Prepare: read one, generate effector
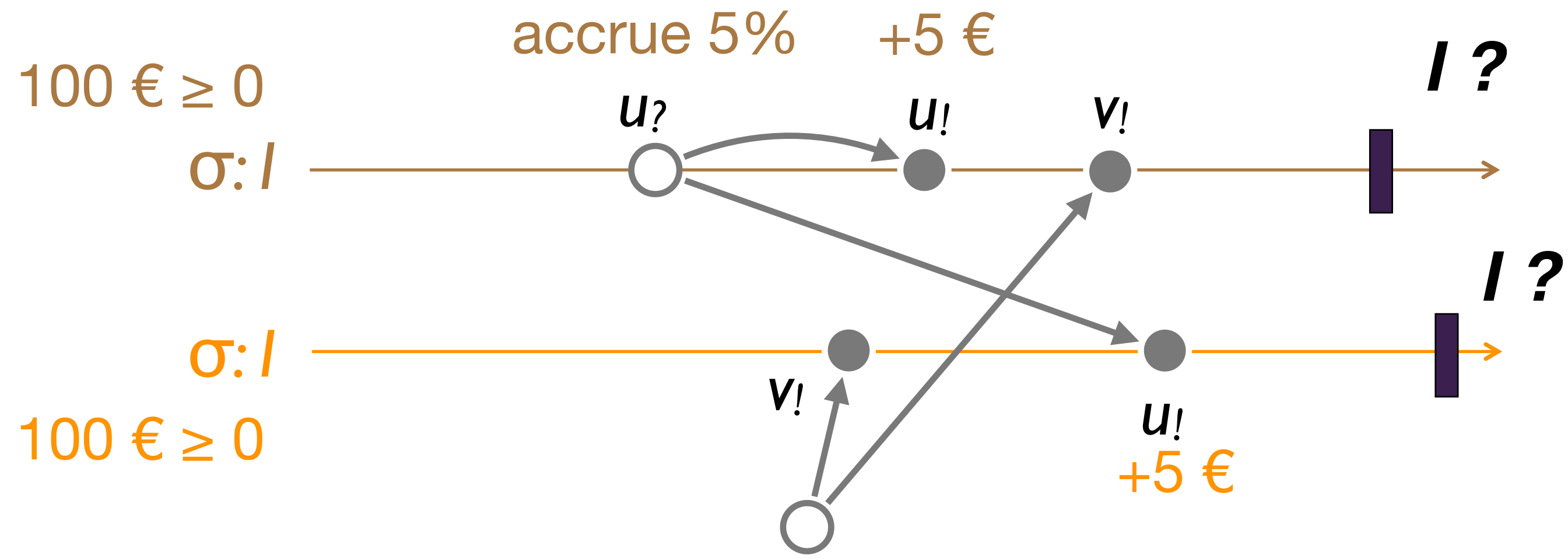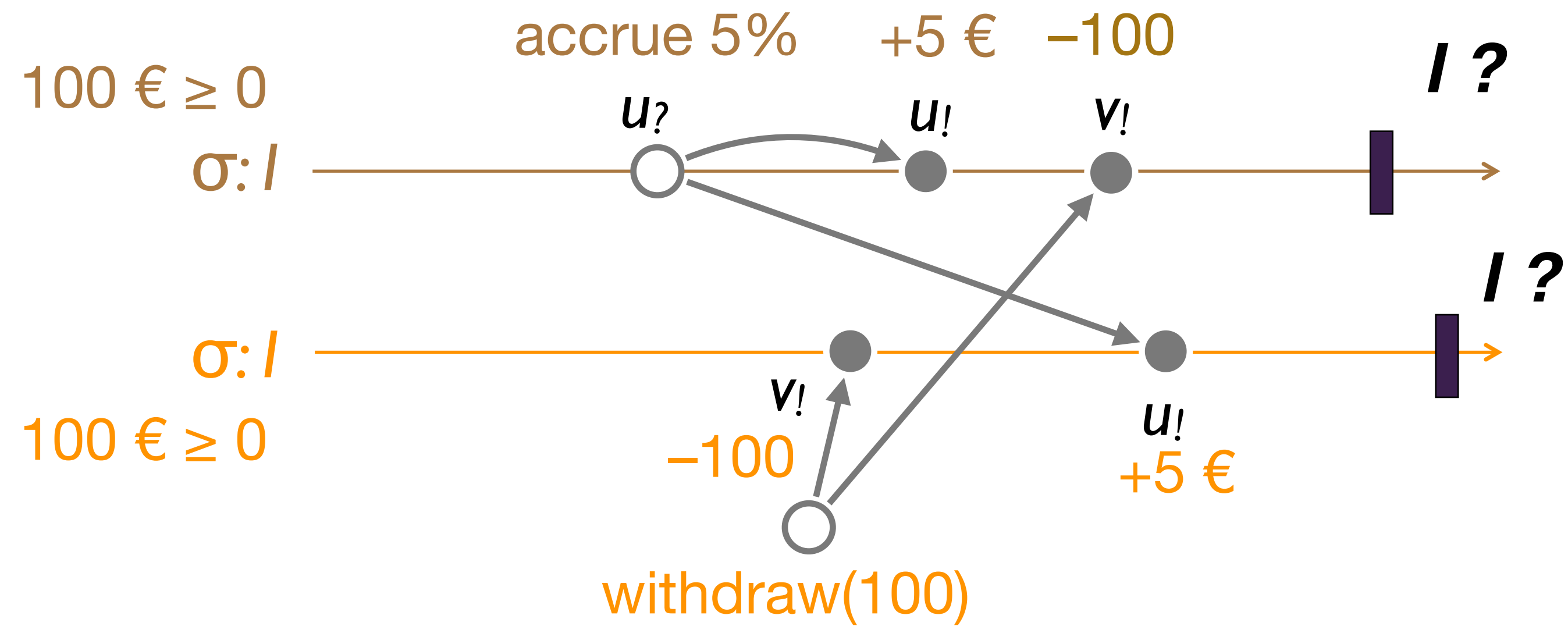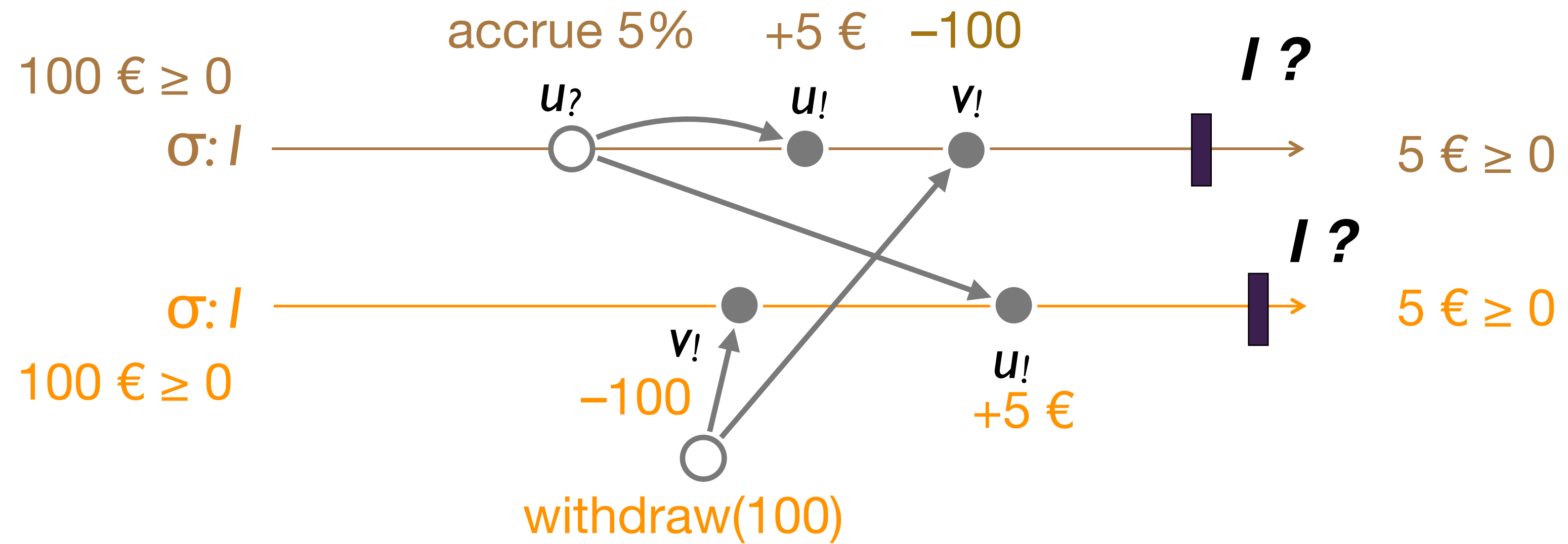
‣ Update all, deferred: deliver effector
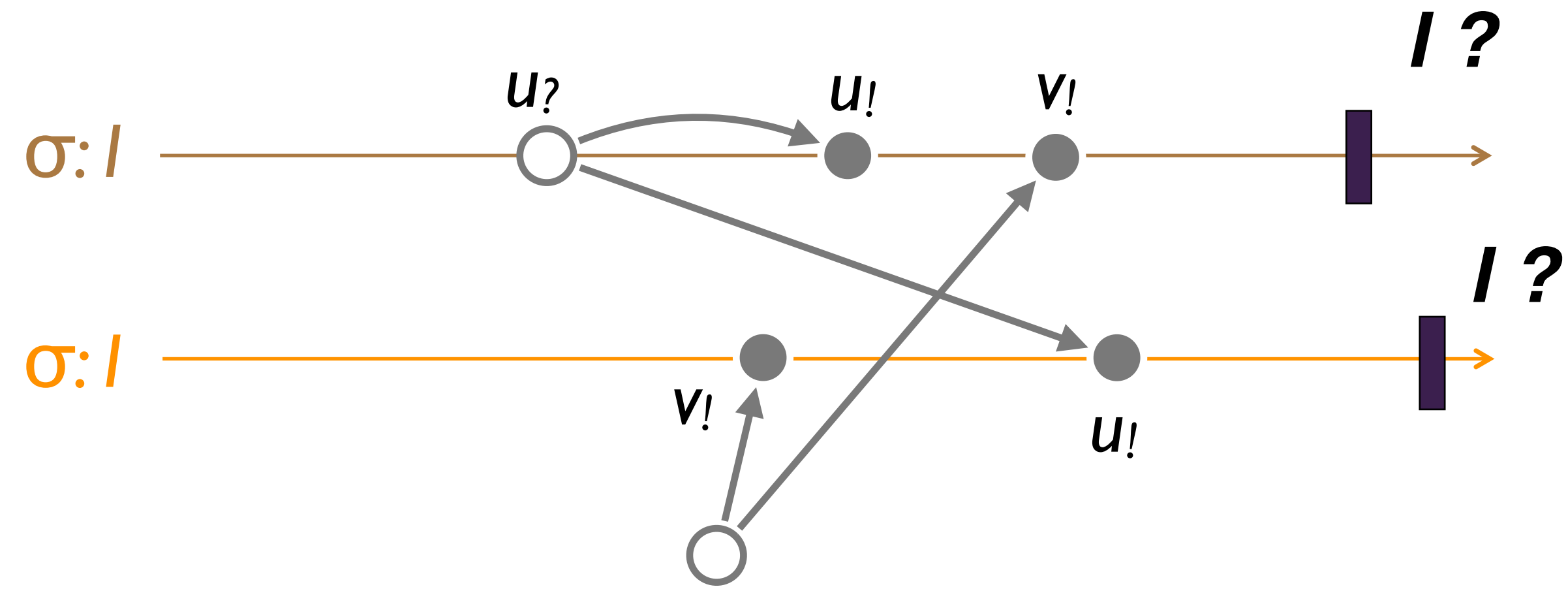
Converge?  Invariant OK?

Asynchronous, replicated updates
- ‣ State $\sigma$
- ‣ Invariant $I$
- ‣ Prepare: read one, generate effector
- ‣ Update all, deferred: deliver effector

Converge?  Invariant OK?

Asynchronous, replicated updates
‣ State $\sigma$
‣ Invariant $I$
‣ Prepare: read one, generate effector
‣ Update all, deferred: deliver effector
Converge?  Invariant OK?

Asynchronous, replicated updates
‣ State $\sigma$
‣ Invariant $I$
‣ Prepare: read one, generate effector
‣ Update all, deferred: deliver effector
Converge? Invariant OK?

Asynchronous, replicated updates
‣ State $\sigma$
‣ Invariant $I$
‣ Prepare: read one, generate effector
‣ Update all, deferred: deliver effector
Converge?  Invariant OK?

Asynchronous, replicated updates

‣ State $\sigma$

‣ Invariant $I$

‣ Prepare: read one, generate effector

‣ Update all, deferred: deliver effector

Converge?  Invariant OK?

Asynchronous, replicated updates

‣ State $\sigma$

‣ Invariant $I$

‣ Prepare: read one, generate effector

‣ Update all, deferred: deliver effector

Converge?  Invariant OK?

Asynchronous, replicated updates

‣ State $\sigma$

‣ Invariant $I$

‣ Prepare: read one, generate effector

‣ Update all, deferred: deliver effector

Converge?  Invariant OK?

CISE Rules

1: Sequential correctness
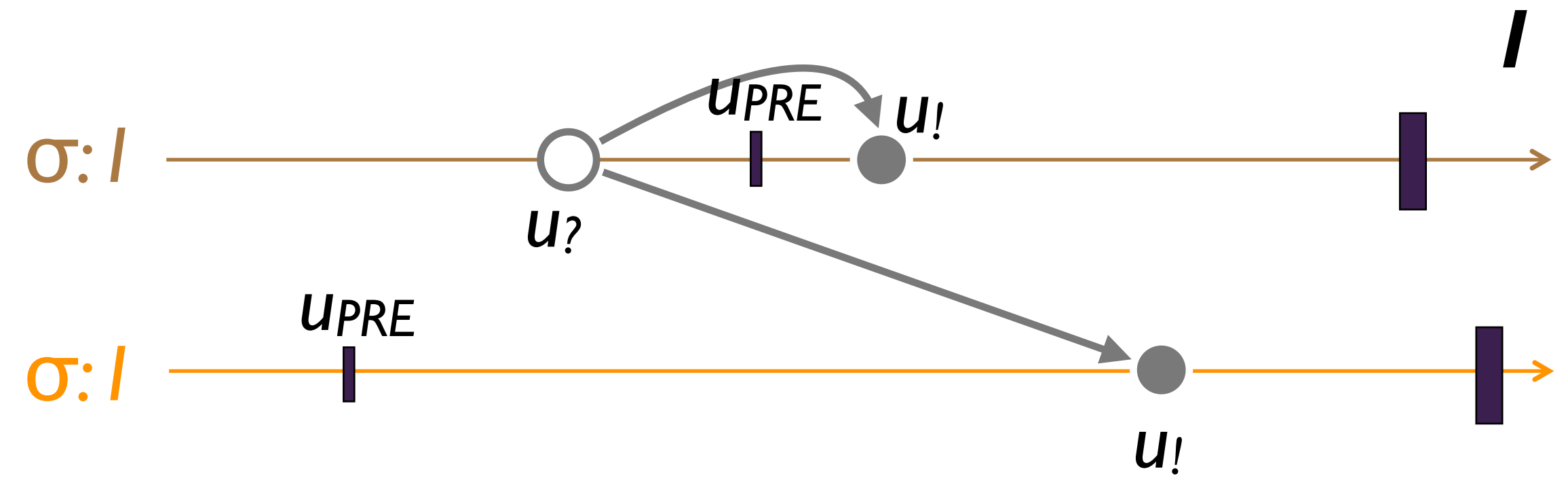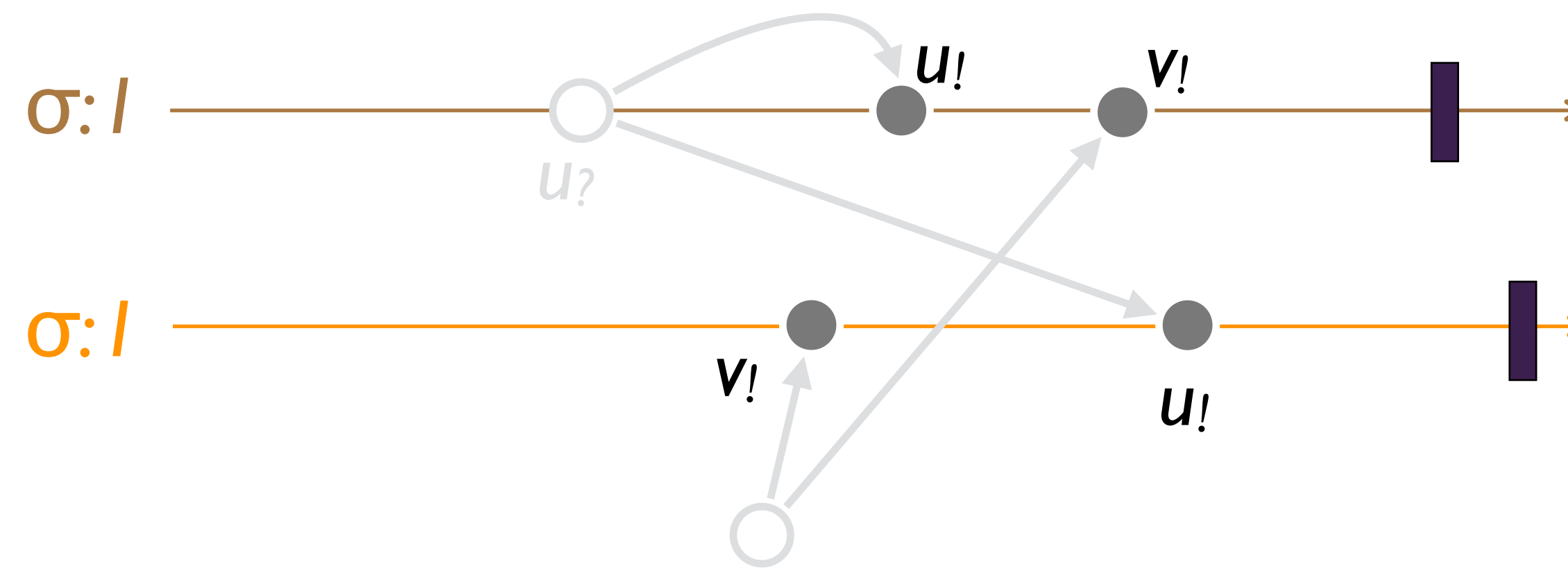‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

# Simple example: bank account

‣ Operations: *deposit(amount), withdraw(amount)*

‣ Invariant: *balance ≥ 0*

   ‣ Start with weak specification

   ‣ Rule 1 ⟶ strengthen precondition for withdraw

   ‣ Rule 2: OK

   ‣ Rule 3 ⟶ *withdraw || withdraw* unsafe

      ‣ fixed with concurrency control

## CISE Rules

# 1: Sequential correctness
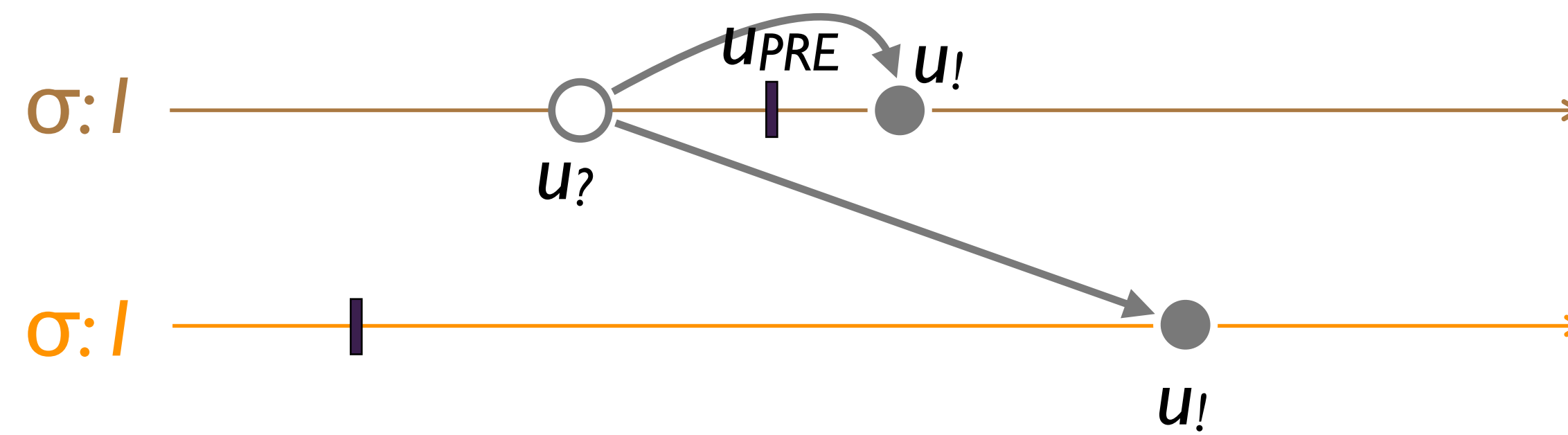‣ Any single operation maintains the invariant

## 2: Convergence
‣ Concurrent effectors commute

## 3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

## CISE Rules

**1: Sequential correctness**
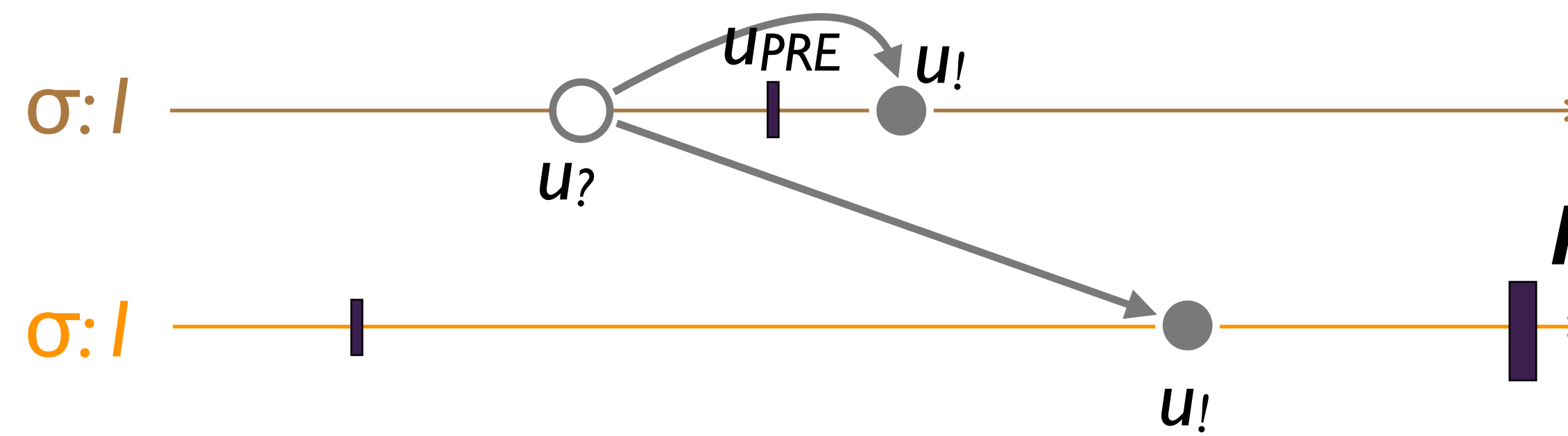‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

# 1: Sequential correctness
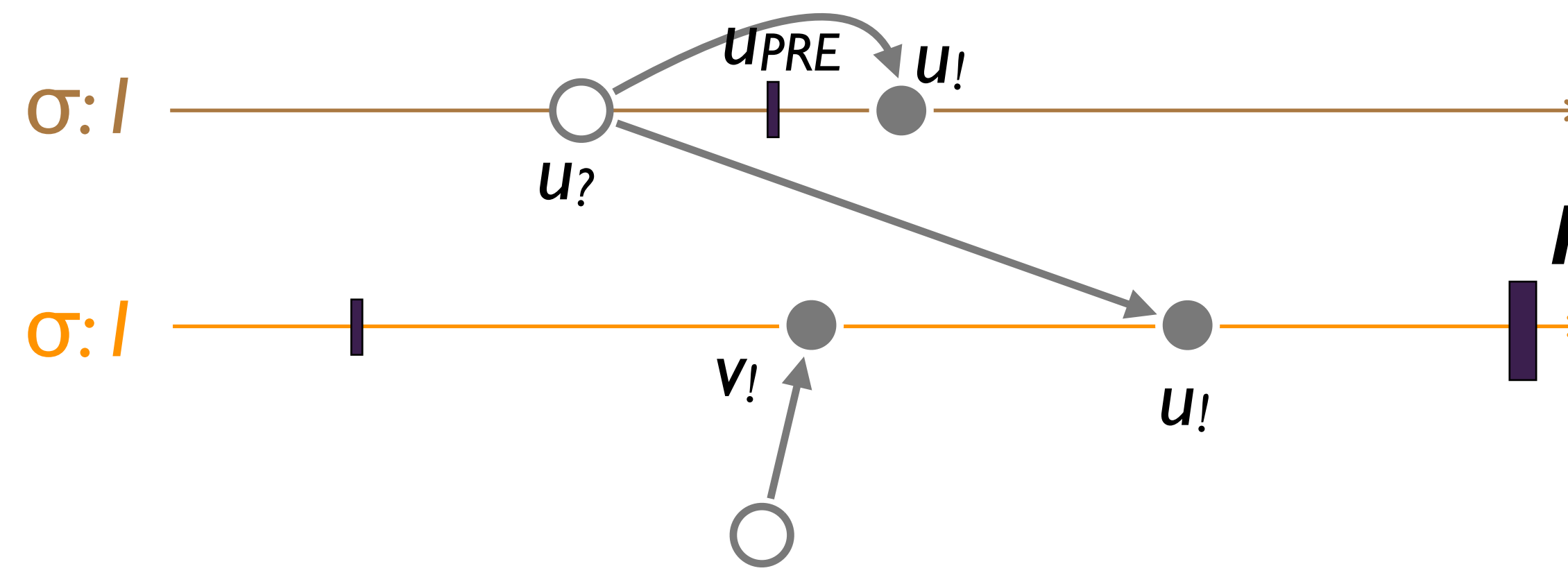- ‣ Any single operation maintains the invariant

## 2: Convergence
- ‣ Concurrent effectors commute

## 3: Precondition Stability
- ‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
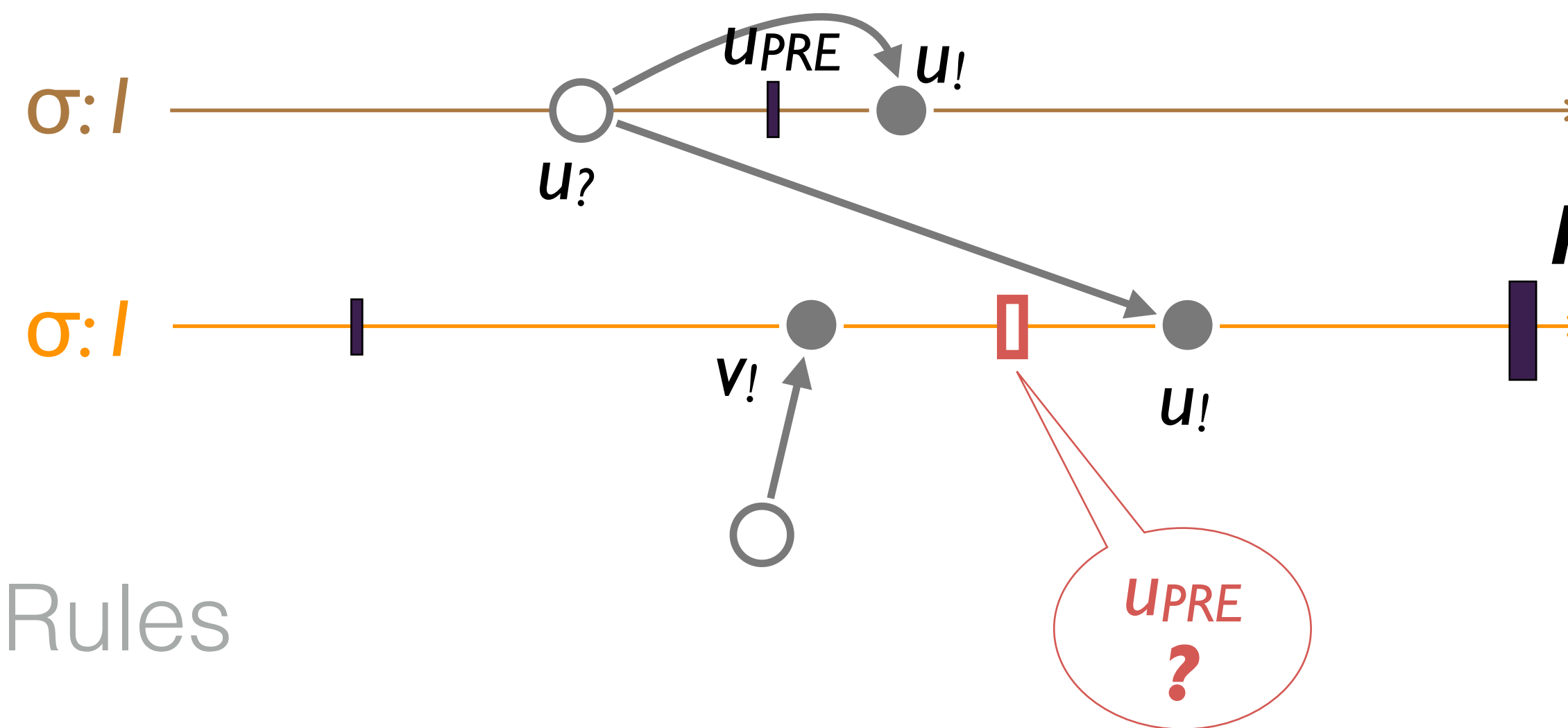‣ Any single operation maintains the invariant

**2: Convergence**
‣ **Concurrent effectors commute**

3: Precondition Stability
‣ Every precondition is stable under every
  concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
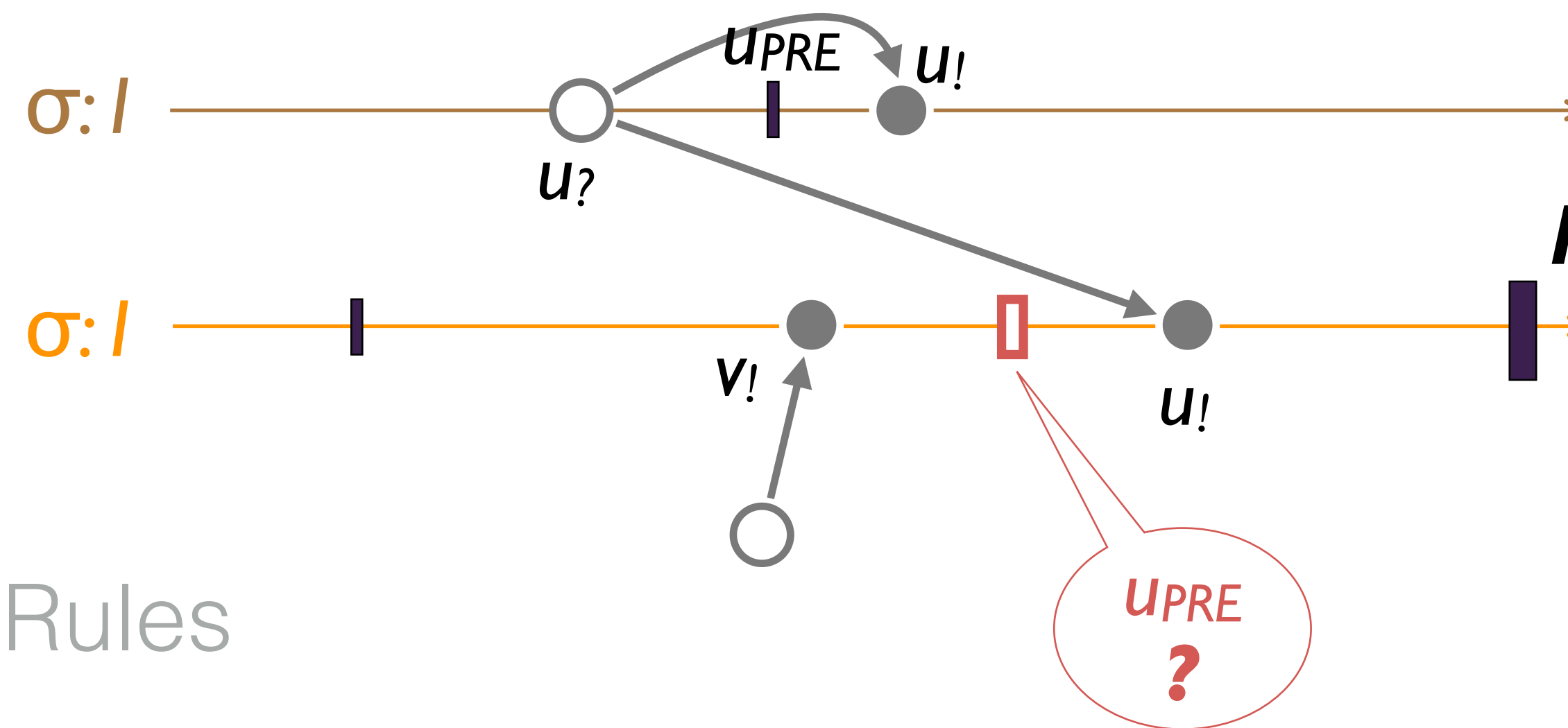‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

**3: Precondition Stability**
‣ **Every precondition is stable under every concurrent operation**

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
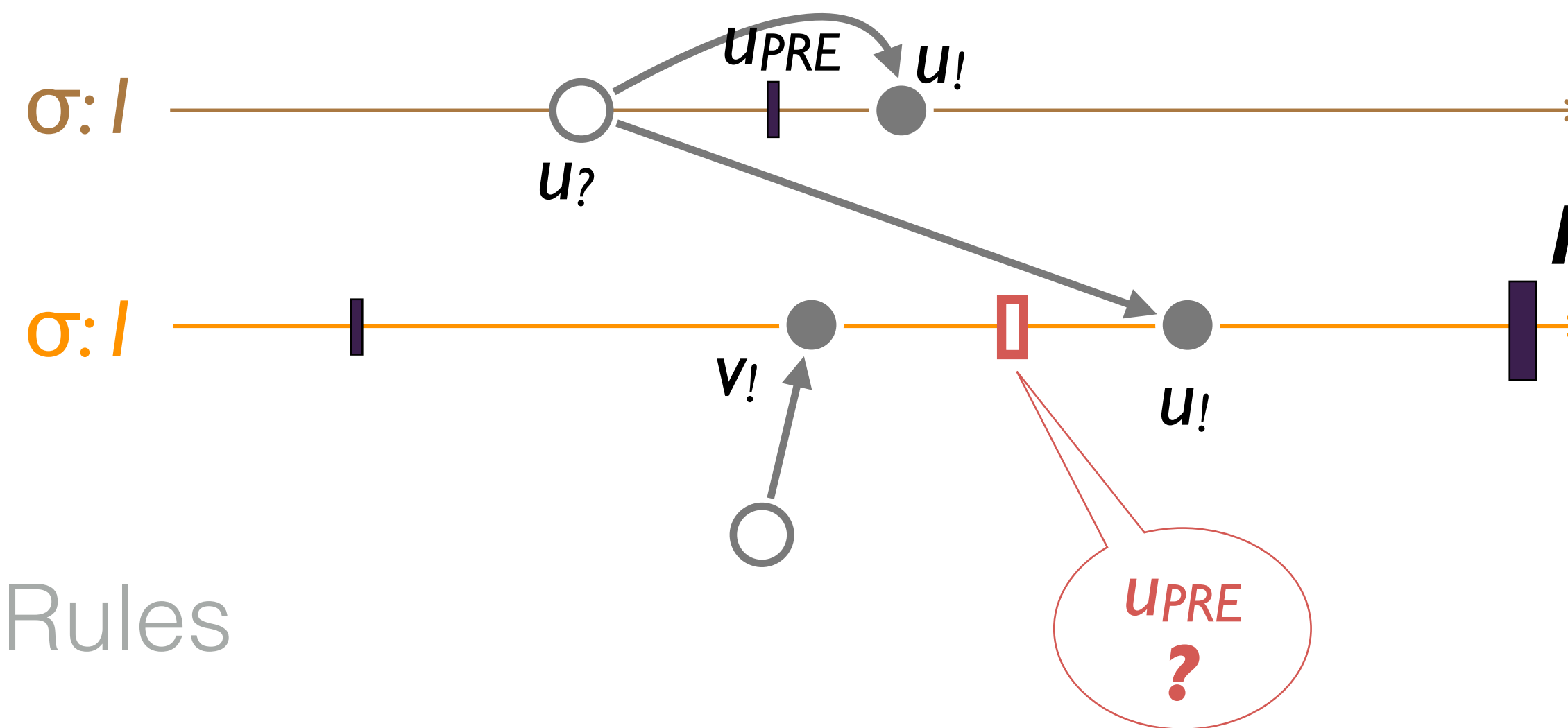‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
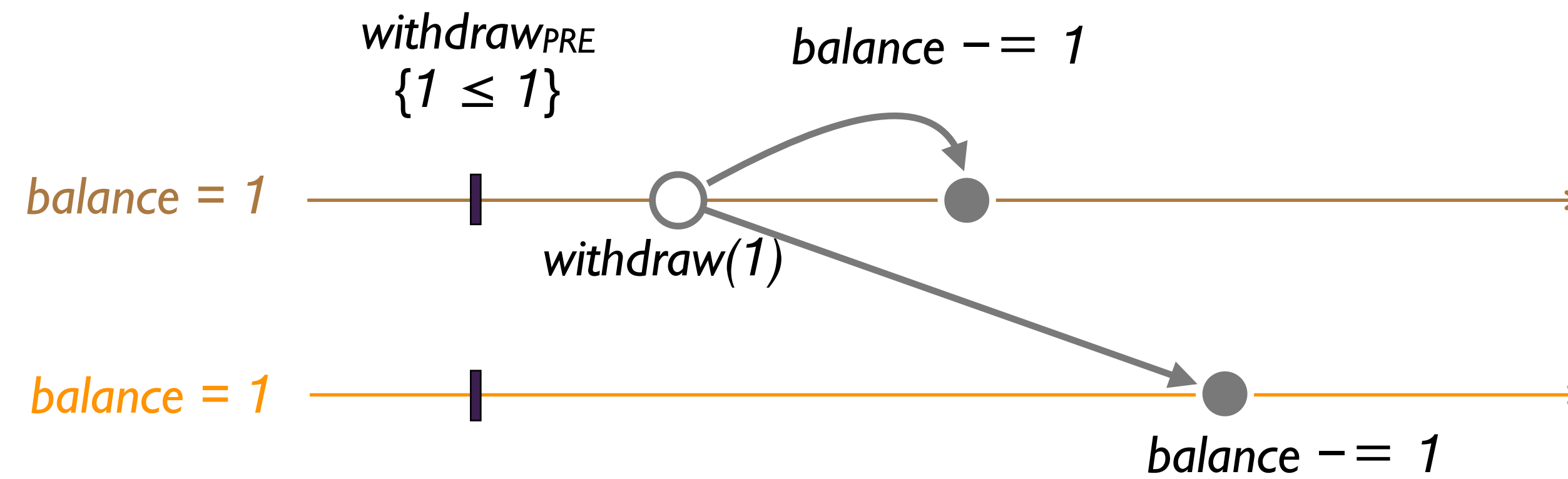‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
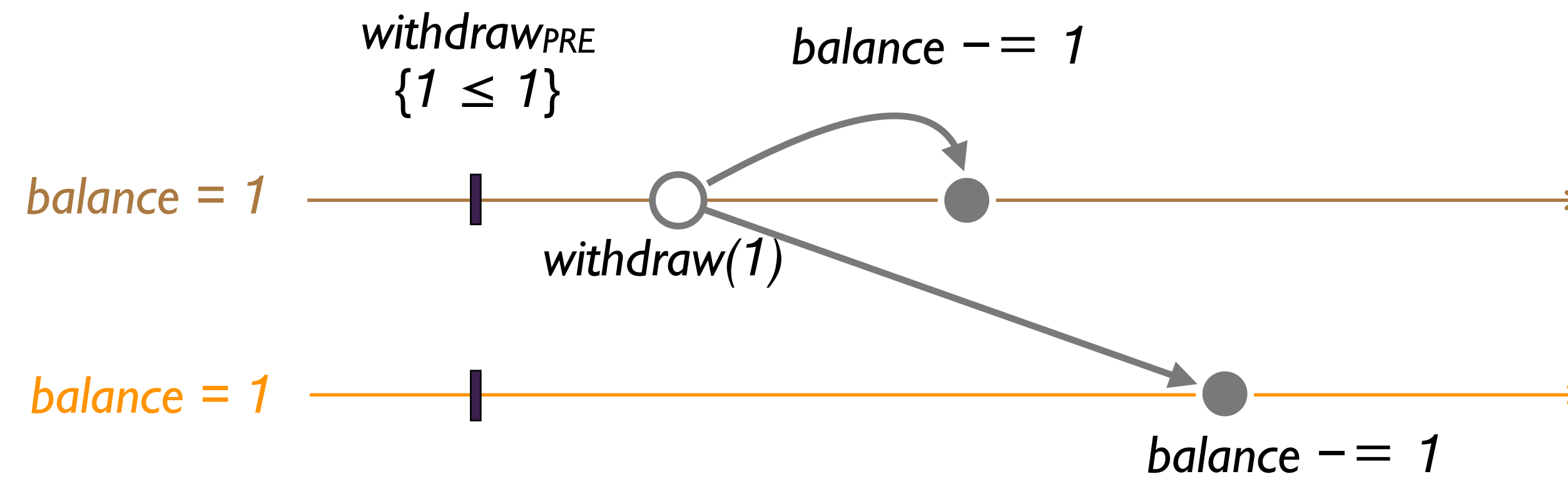‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

**3: Precondition Stability**
‣ **Every precondition is stable under every concurrent operation**

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
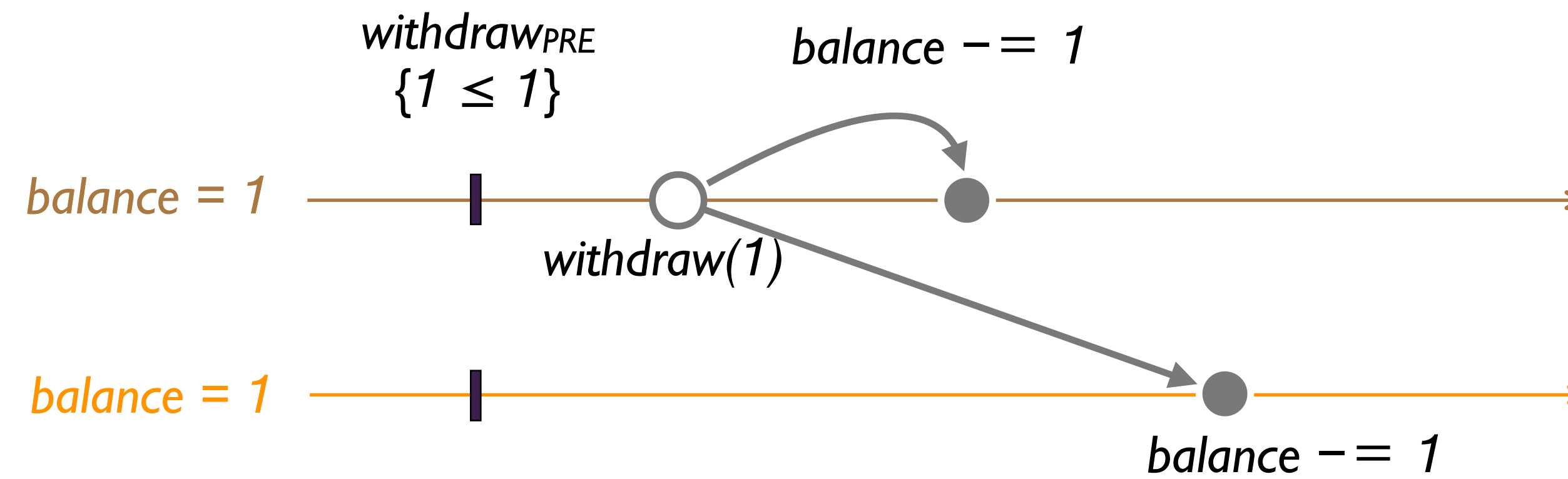‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

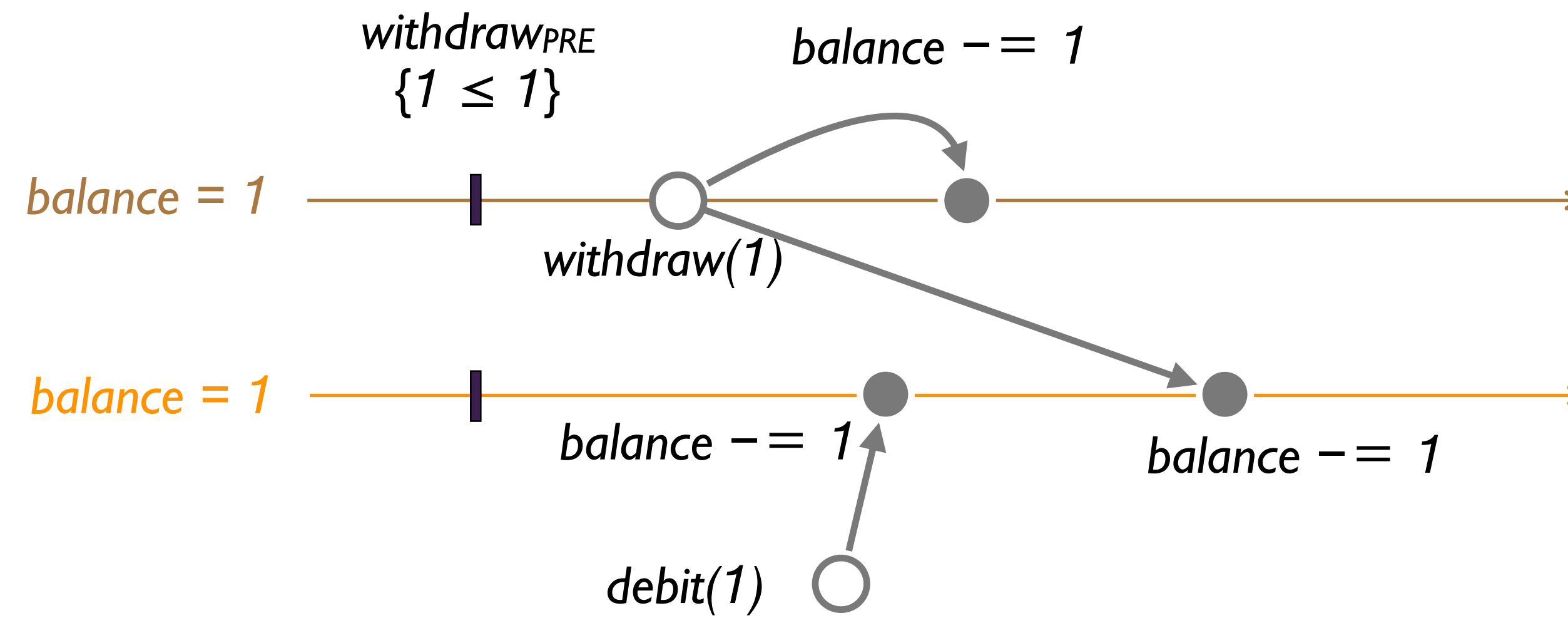2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation
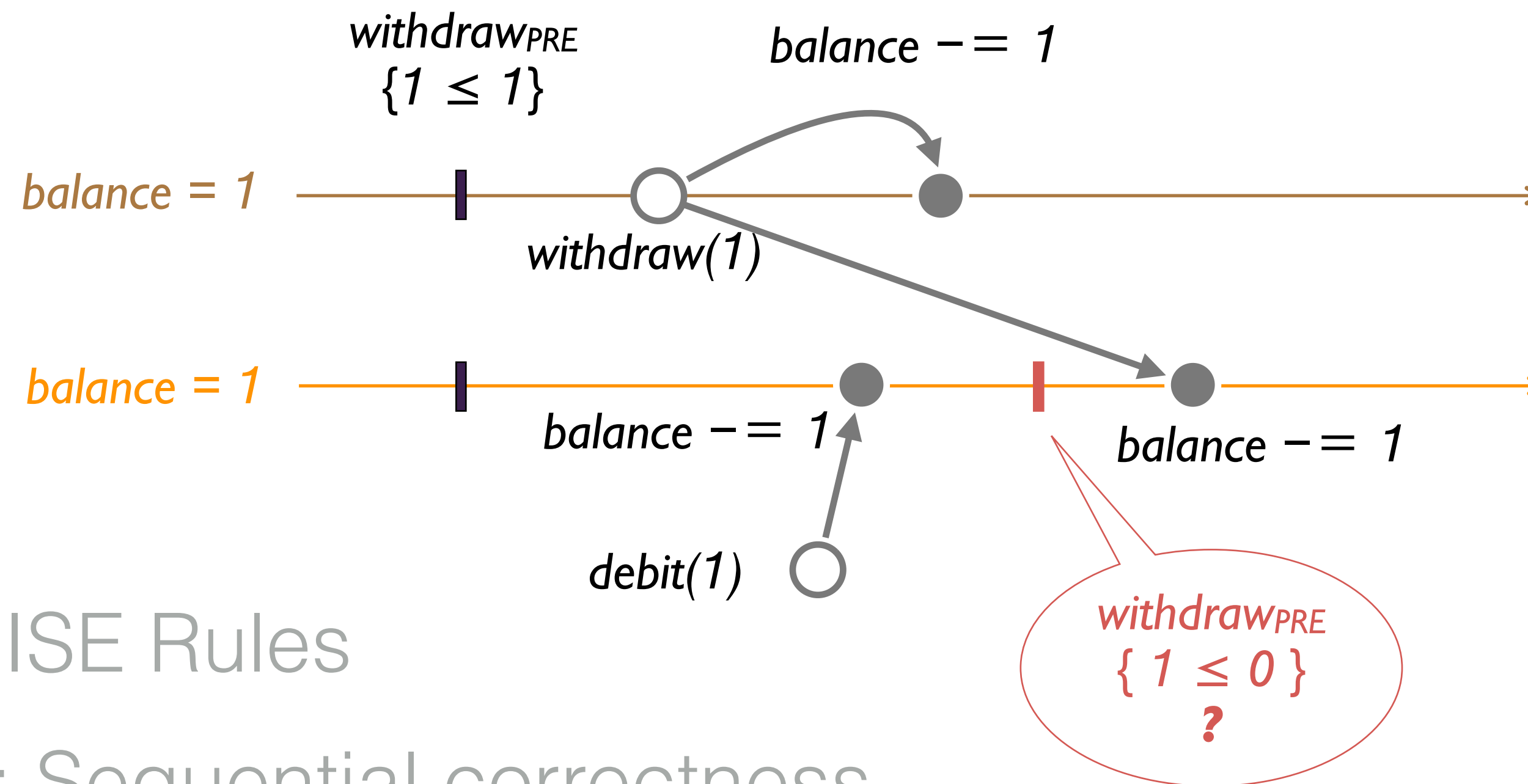
If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every
   concurrent operation
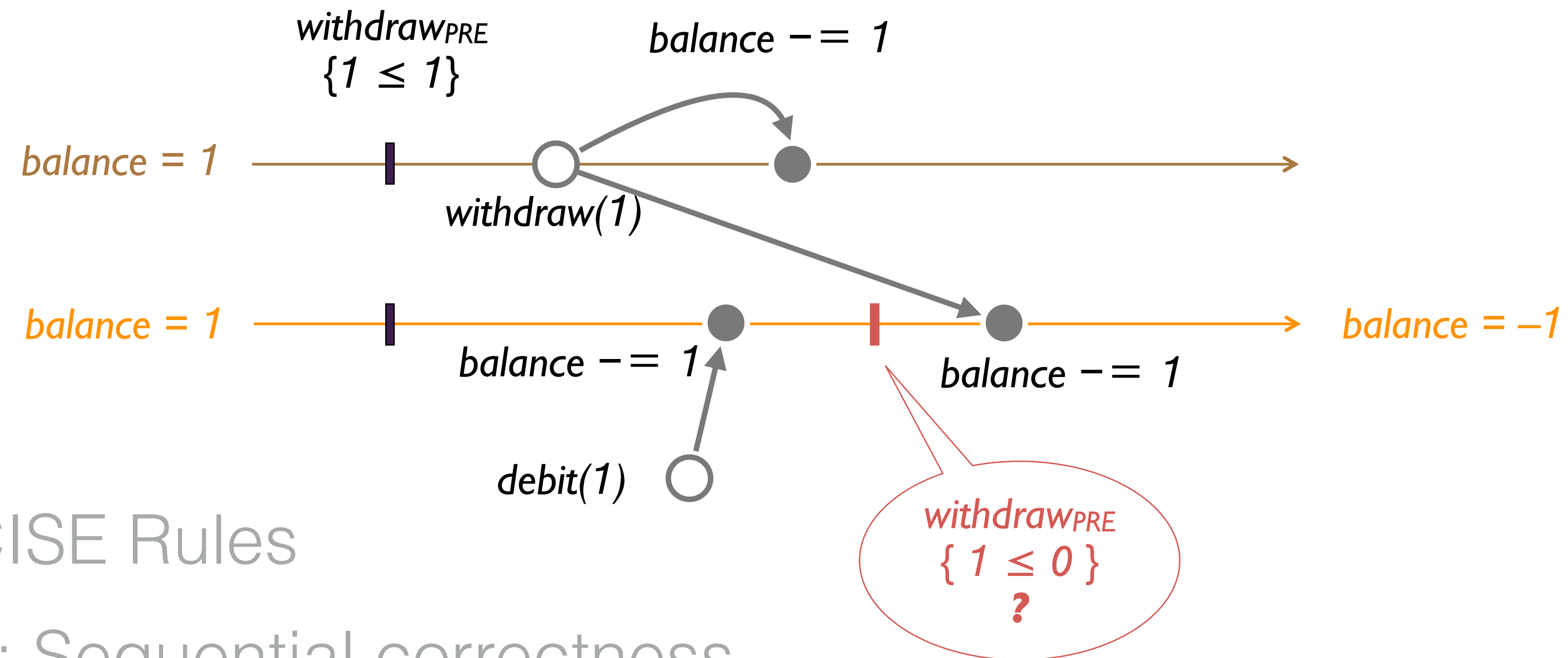
If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every
concurrent operation

If satisfied: invariant is guaranteed

$withdraw_{PRE}$
$\{1 \leq 1\}$

$balance\ -=\ 1$

balance = 1

withdraw(1)

balance = 1

$balance\ -=\ 1$

$balance\ -=\ 1$

debit(1)

## CISE Rules

### 1: Sequential correctness
‣ Any single operation maintains the invariant

### 2: Convergence
‣ Concurrent effectors commute

### 3: Precondition Stability
‣ Every precondition is stable under every concurrent operation
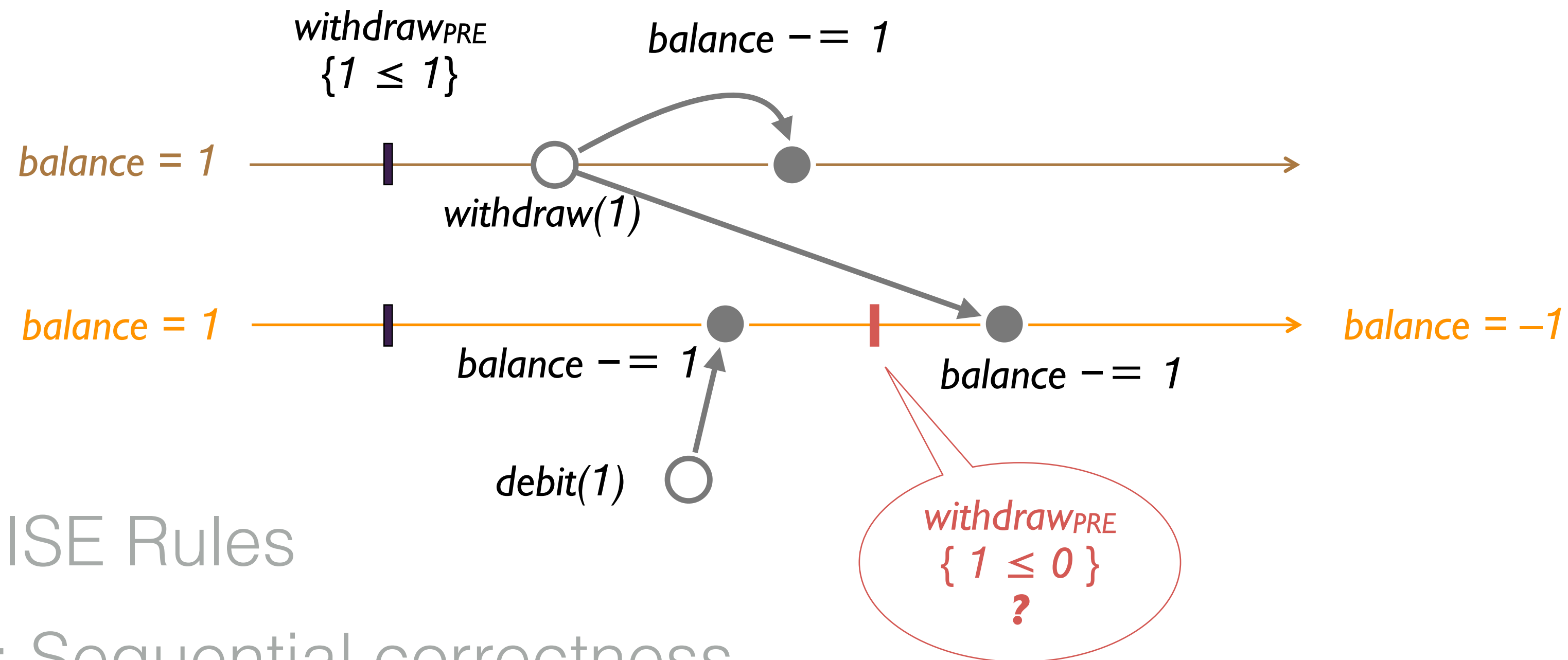
If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

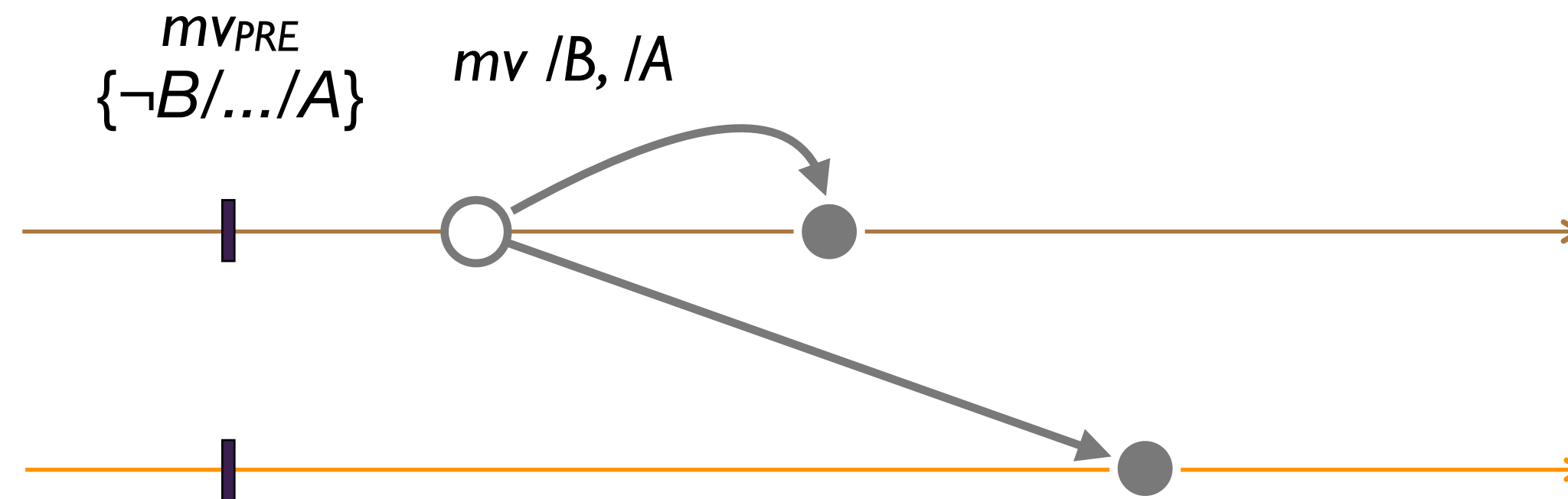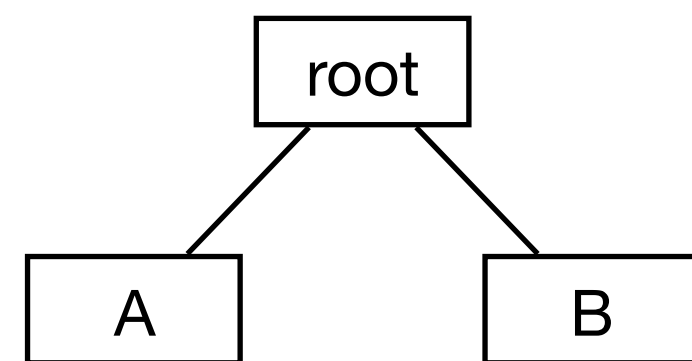If satisfied: invariant is guaranteed

withdraw$_{PRE}$ {1 ≤ 1}

balance −= 1

balance = 1

withdraw(1)

balance = 1    balance = −1

balance −= 1    balance −= 1

debit(1)

withdraw$_{PRE}$ { 1 ≤ 0 } ?

# CISE Rules

## 1: Sequential correctness
‣ Any single operation maintains the invariant

## 2: Convergence
‣ Concurrent effectors commute

## 3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

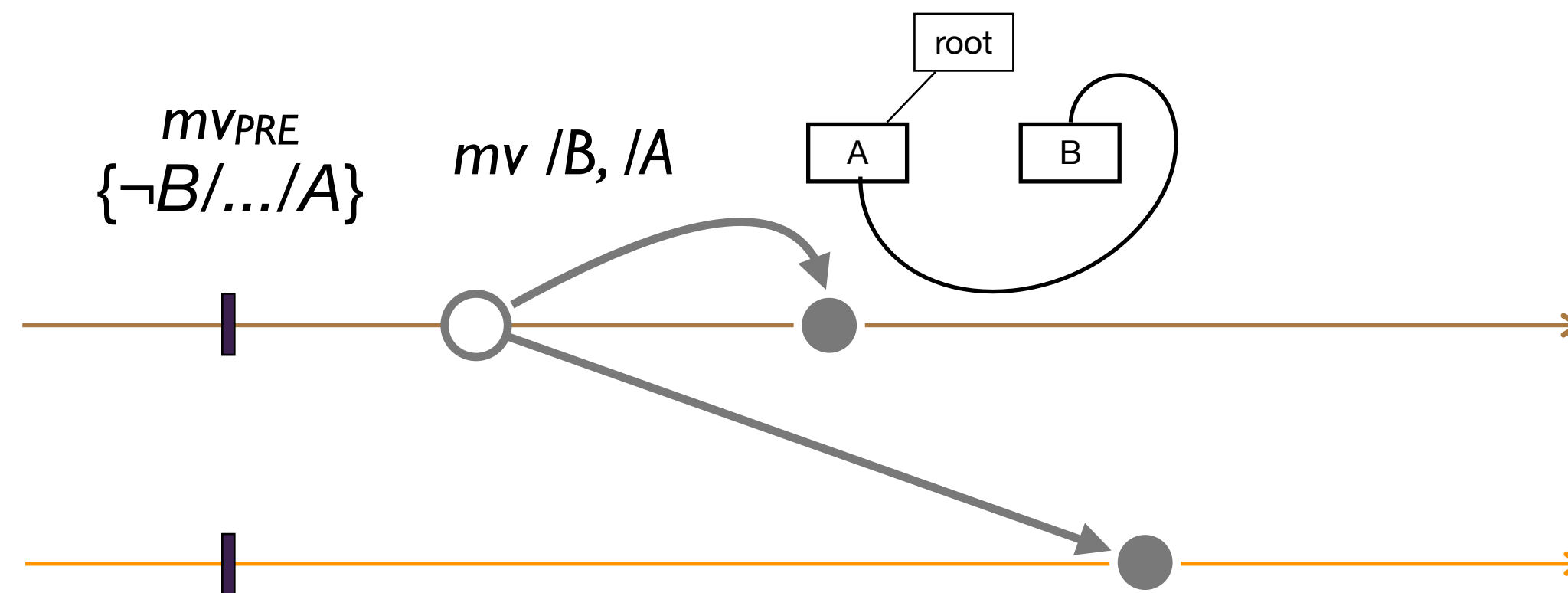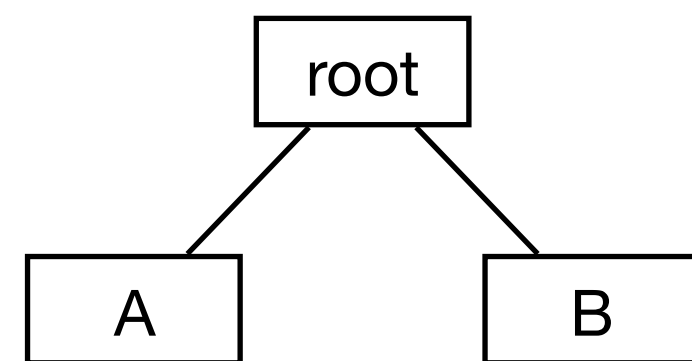1: Sequential correctness
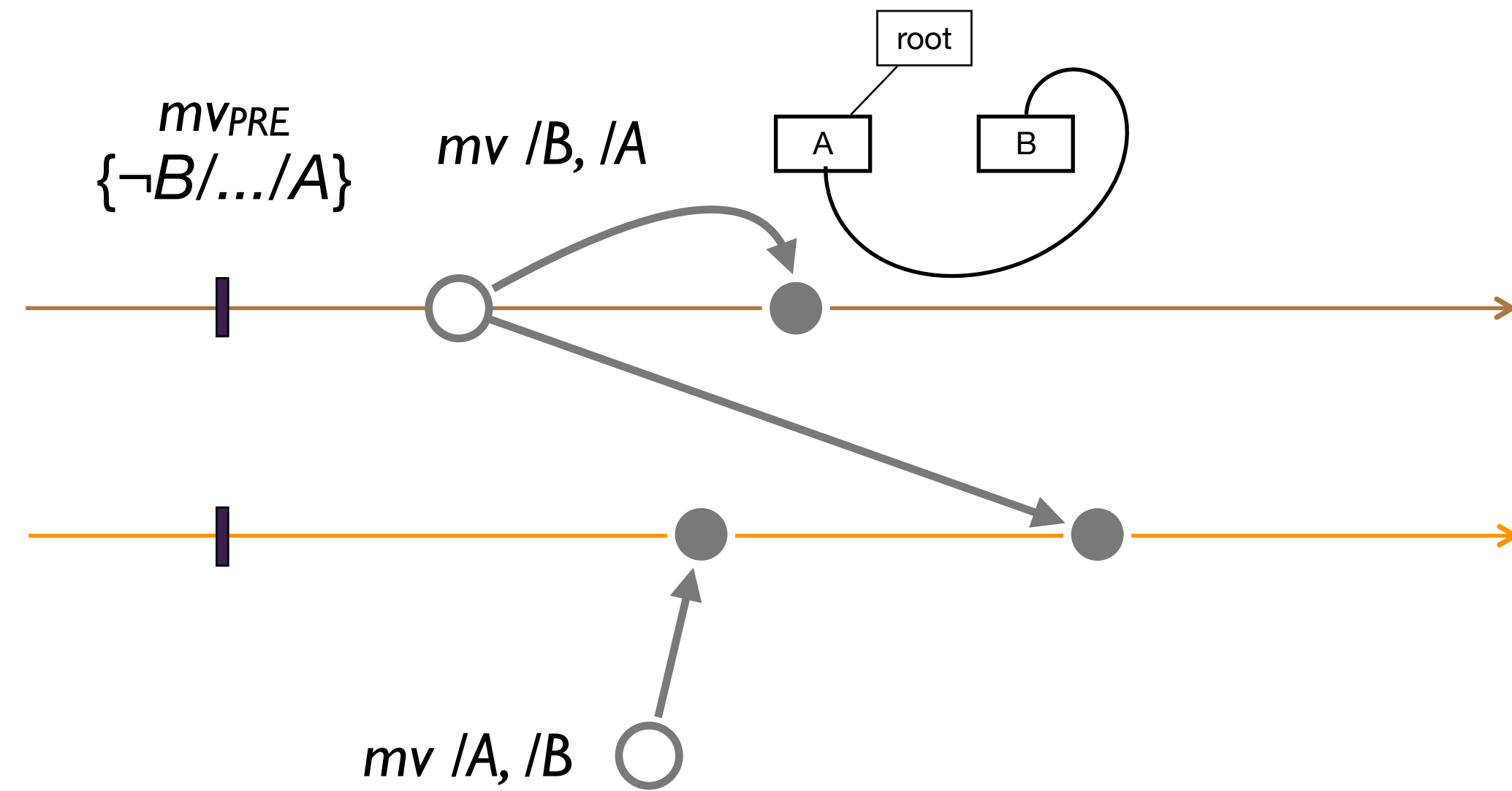‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

# Advanced example: file system

- Operations: *mkdir, rmdir, mv, write*, etc.
- Invariant: Tree
    - Rule 1 $\longrightarrow$ precondition on *mv*
    *"May not **move** node under self"*
    - Rule 2 $\longrightarrow$ Use CRDTs for *write || write*
    - Rule 3 $\longrightarrow$ *mv || mv* precondition unstable

# Advanced example: file system

- Operations: *mkdir, rmdir,* **mv**, *update*, etc.
- Invariant: Tree
  - Rule 1 ⟶ precondition on *mv*
    *"May not **move** node under self"*
  - Rule 2 ⟶ Use CRDTs for *update ‖ update*
- Rule 3 ⟶ *mv ‖ mv* precondition unstable

CISE Rules

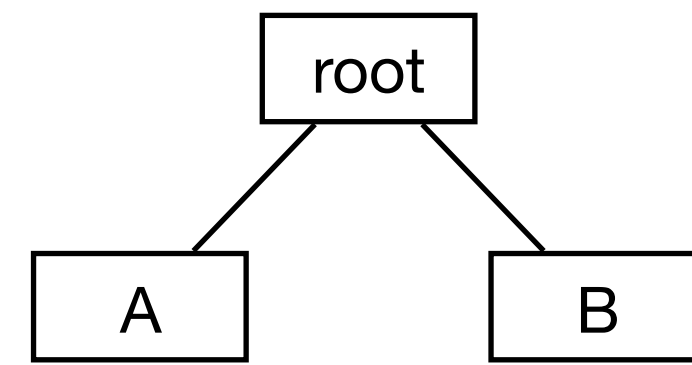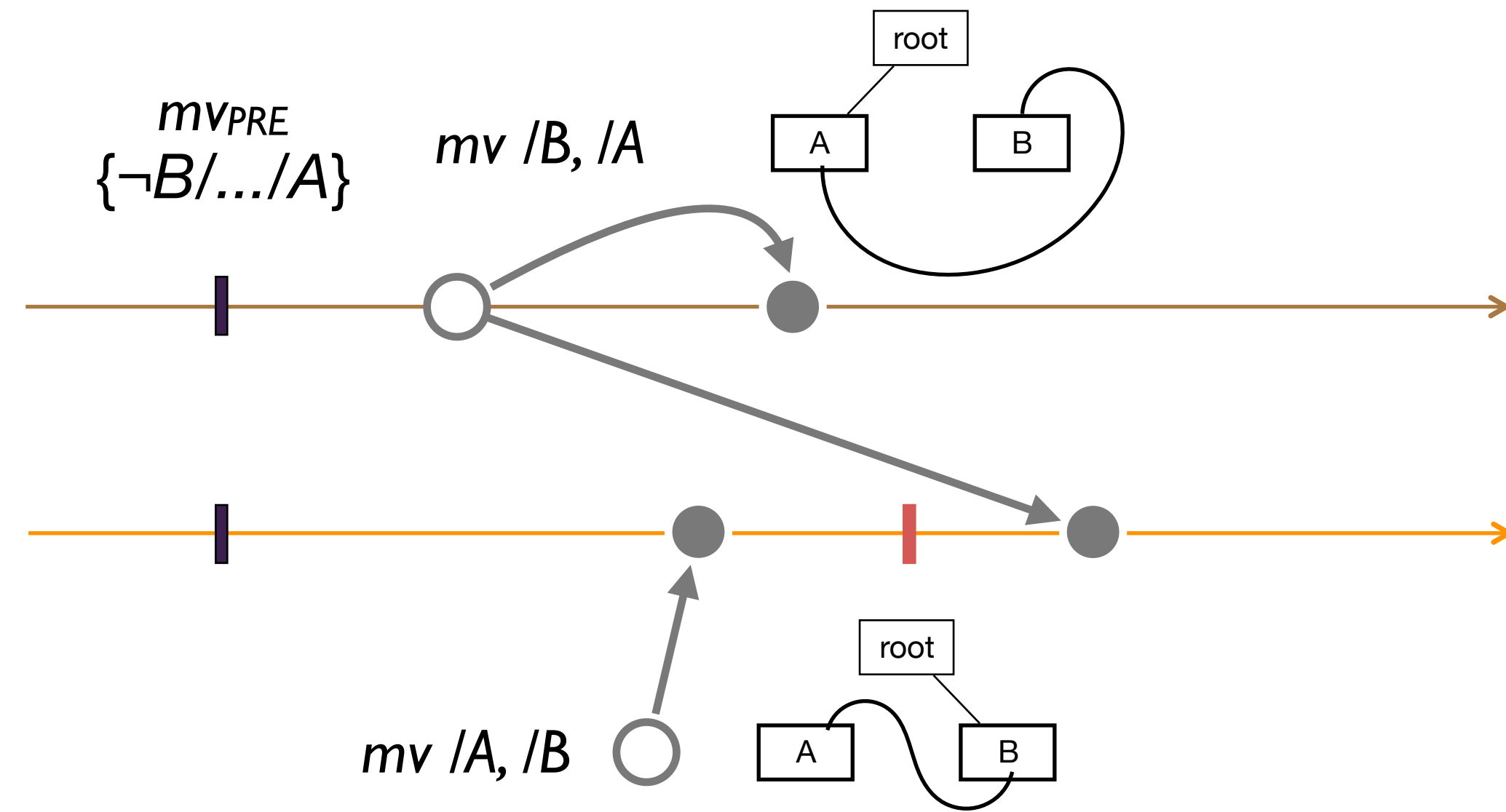1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

# CISE Rules

**1: Sequential correctness**
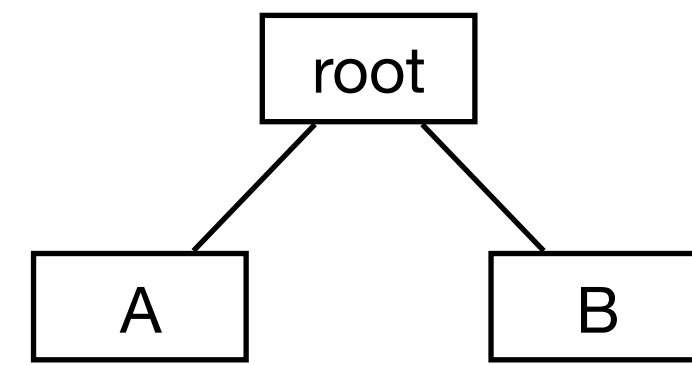‣ Any single operation maintains the invariant

**2: Convergence**
‣ Concurrent effectors commute

**3: Precondition Stability**
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

# CISE Rules

## 1: Sequential correctness
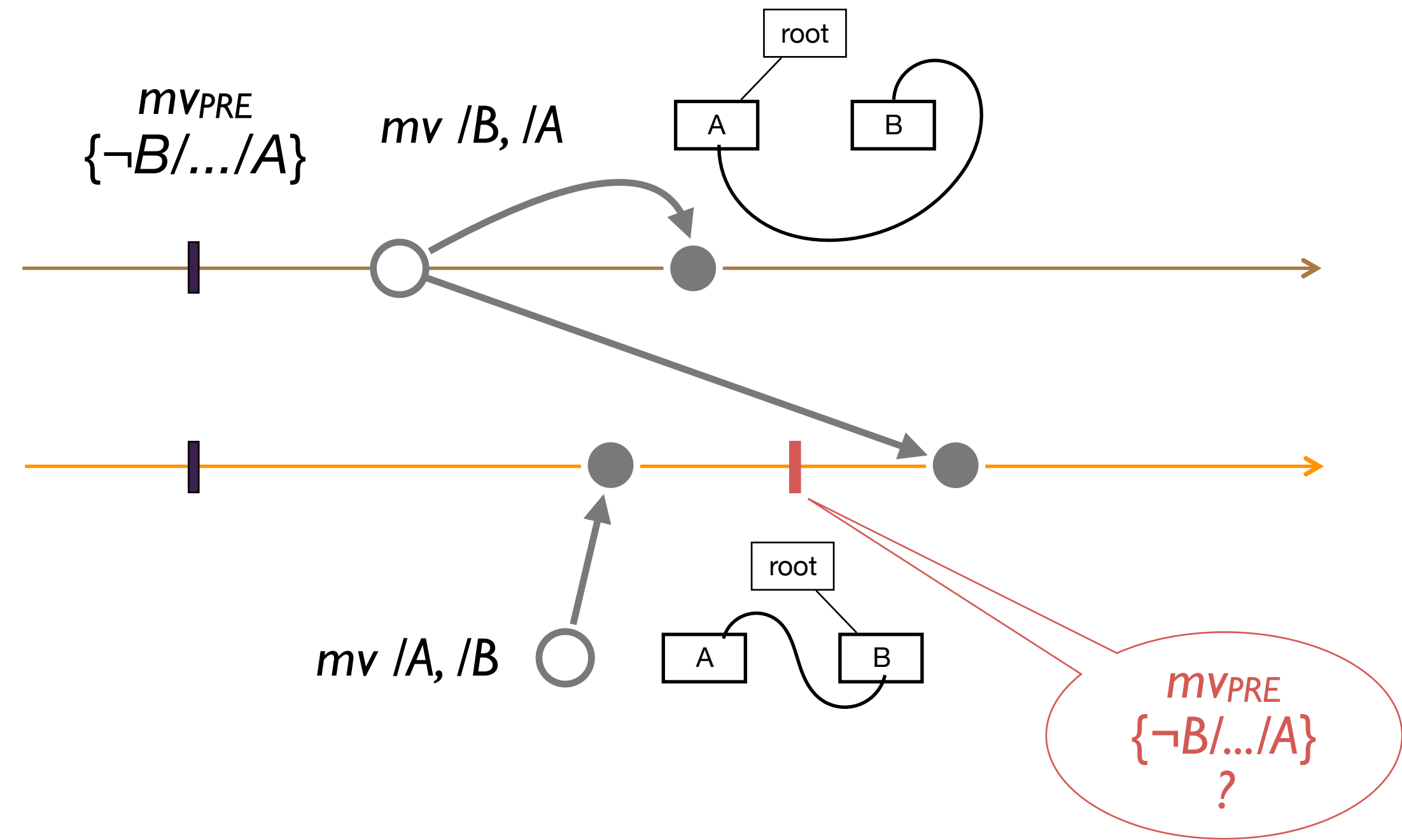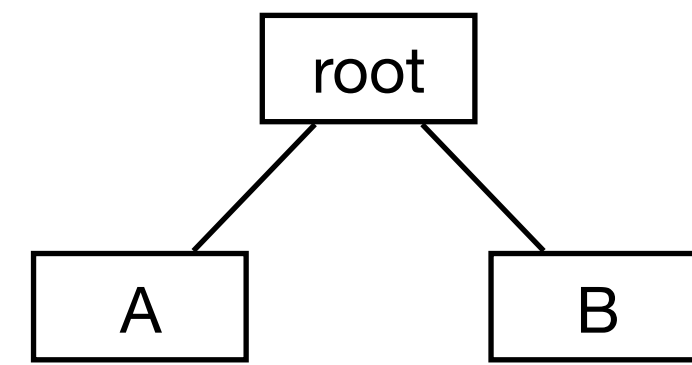‣ Any single operation maintains the invariant

## 2: Convergence
‣ Concurrent effectors commute

## 3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

$mv_{PRE}$
$\{\neg B/.../A\}$

$mv$ /B, /A

$mv$ /A, /B

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

$mv_{PRE}$
$\{\neg B/.../A\}$

$mv\ /B,\ /A$

$mv\ /A,\ /B$

$mv_{PRE}$
$\{\neg B/.../A\}$
?

CISE Rules

1: Sequential correctness
‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
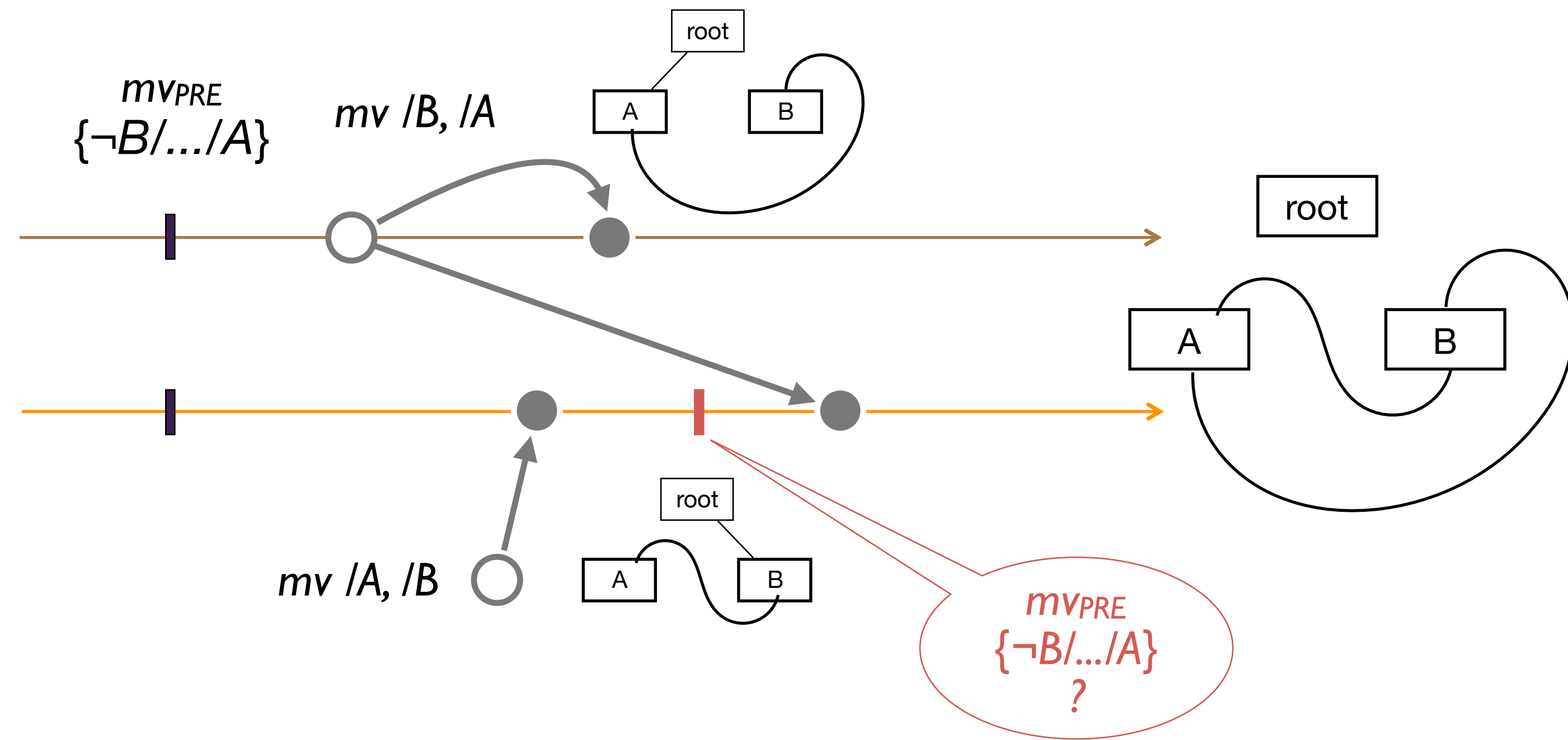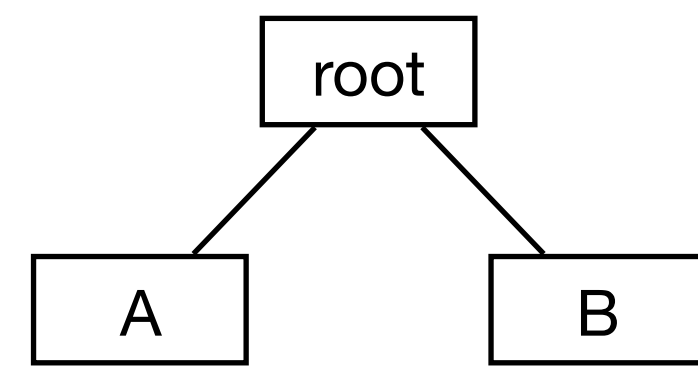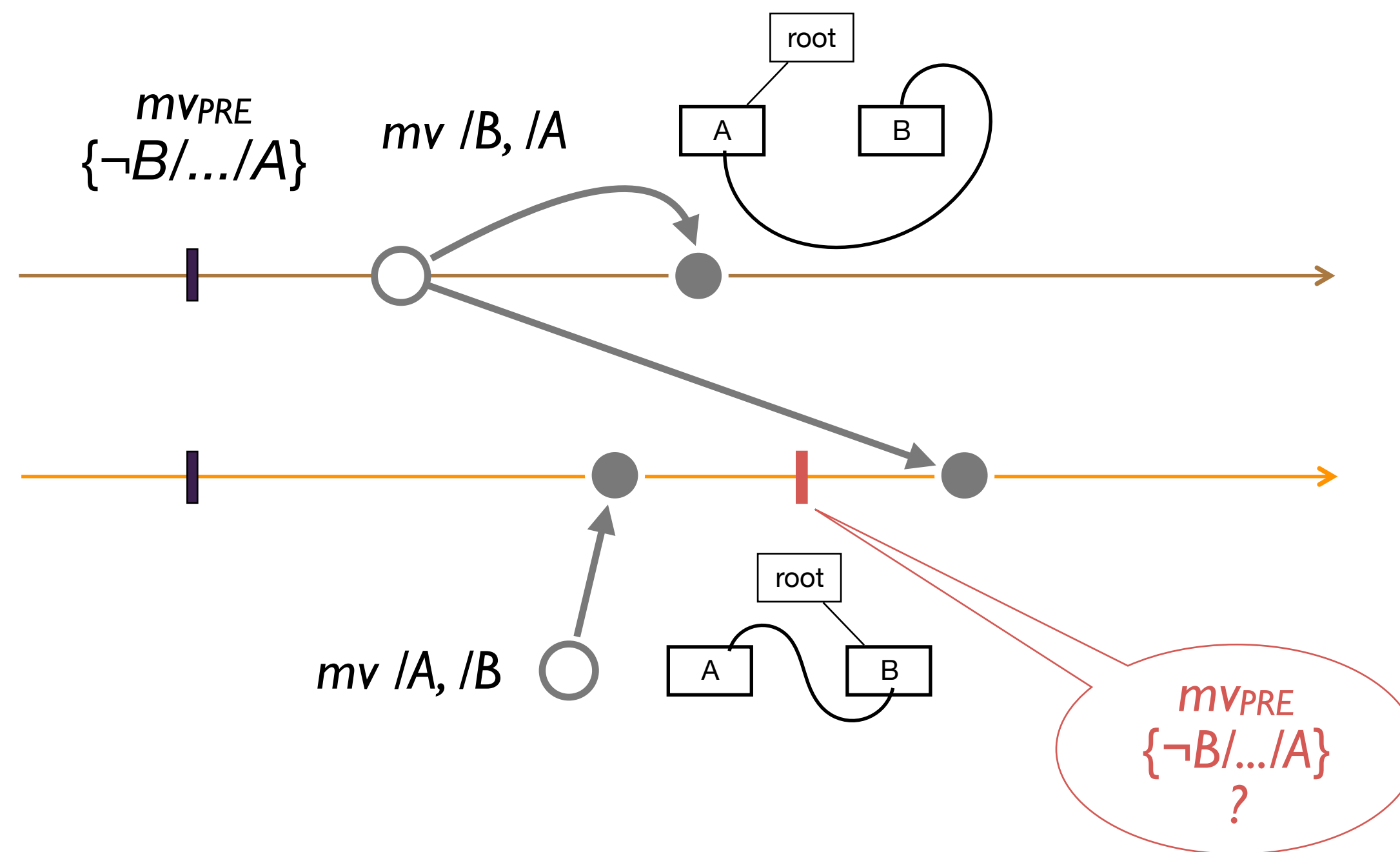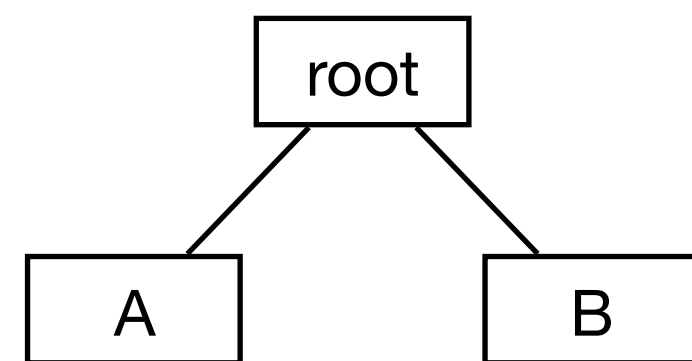‣ Any single operation maintains the invariant

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

CISE Rules

1: Sequential correctness
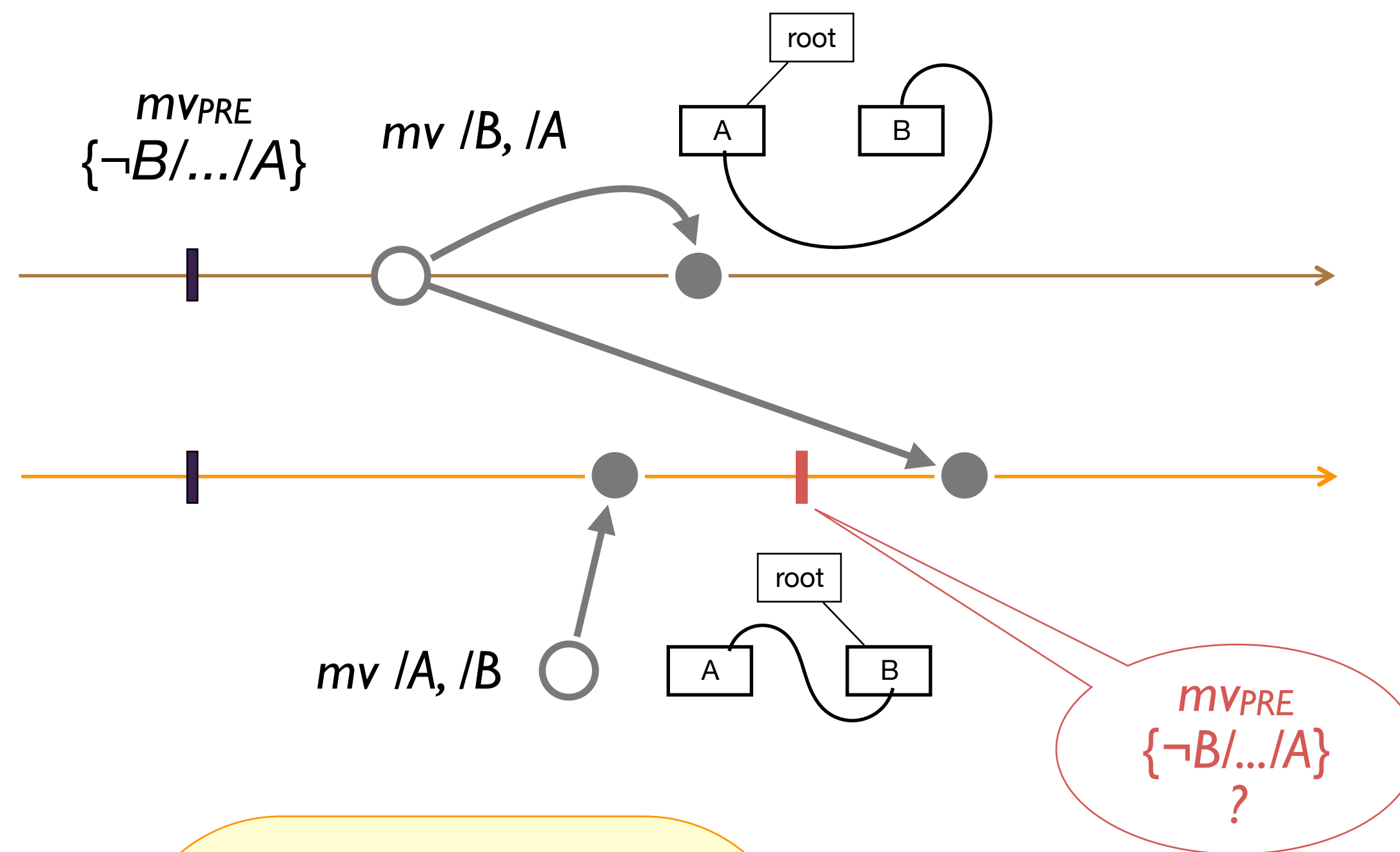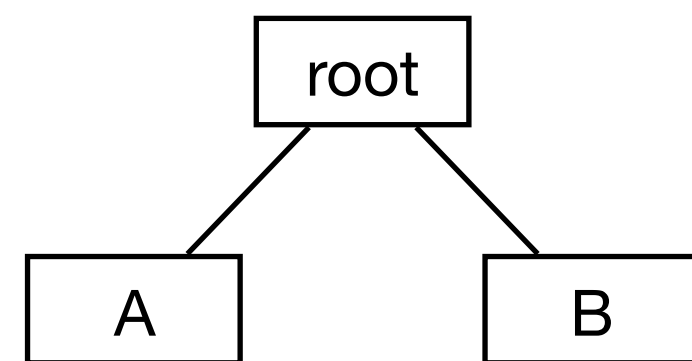‣ Any single operation maintains the invari...

2: Convergence
‣ Concurrent effectors commute

3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

# CISE Rules

## 1: Sequential correctness
‣ Any single operation maintains the i

## 2: Convergence
‣ Concurrent effectors commute

## 3: Precondition Stability
‣ Every precondition is stable under every concurrent operation

If satisfied: invariant is guaranteed

You can have your cake and eat it too

# CISE: The tool

Version of the tool (CEC) by Sreeja Nair

# Related Problems

▸ Going beyond single invariants

    ▸ Verify Pre/Post conditions of client programs

▸ State-Based implementations of CRDTs

▸ Composition of CRDTs

▸ … and much more :-)

# The END