

# A rigorous approach to consistency in cloud databases

Alexey Gotsman

IMDEA Software Institute, Madrid, Spain

Amazon.co.uk: Low Prices in Electronics, Books, Sports Equipment & more

Amazon.co.uk Your Amazon.co.uk Today's Deals Gift Cards Help

Shop by Department Search All Go Hello, Sign in Your Account Basket Wish List

¿Compras desde España? Shopping from Spain? Visita amazon.es Descubre

Amazon MP3 Cloud Player Kindle LOVEFILM Appstore for Android Audible

Meet the Kindle Fire

January Deals > Shop now

Two-Hour Flying Lesson Take to the skies! £99 (was £299) > See the deal amazonlocal

SHAMBALLA BRACELETS > Shop now

Google

https://www.google.com/?gws\_rd=ssl

Apple Yahoo! Google Maps YouTube Wikipedia News Popular

+You Gmail Images Sign in

# Google

Google Search I'm Feeling Lucky

Welcome to Facebook - Log In, Sign Up or Learn More

facebook

Email or Phone

Keep me logged in

## Sign Up

It's free and always stays free.

First name

Email

Re-enter email

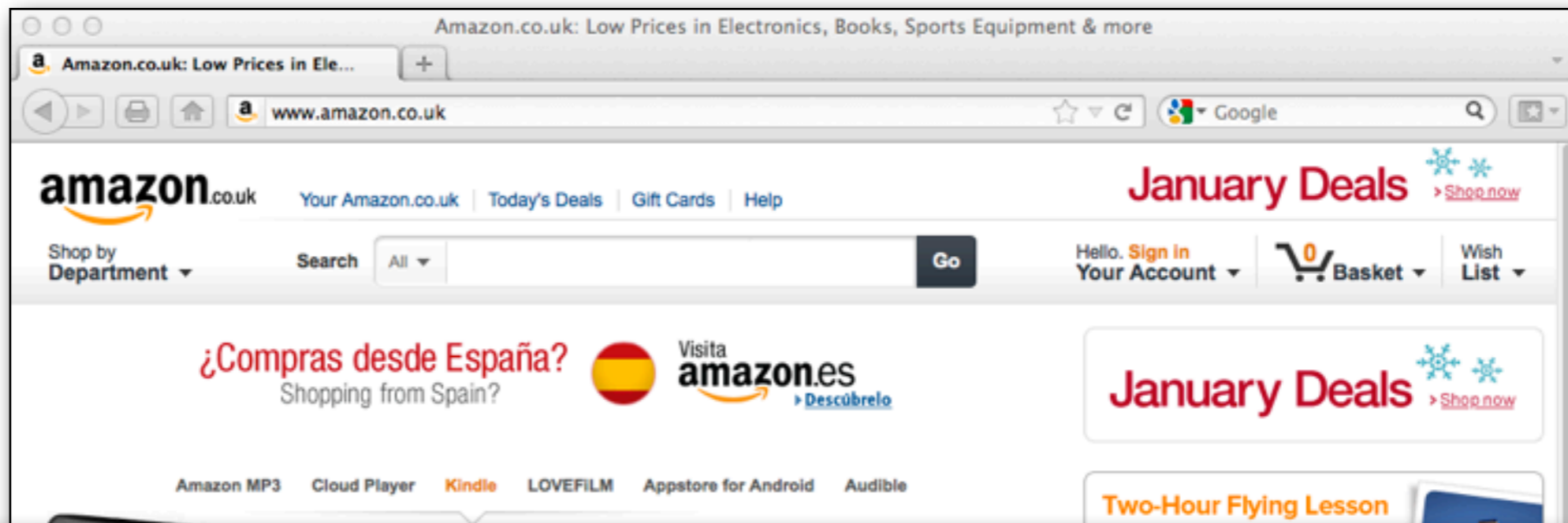
New password

Birthday

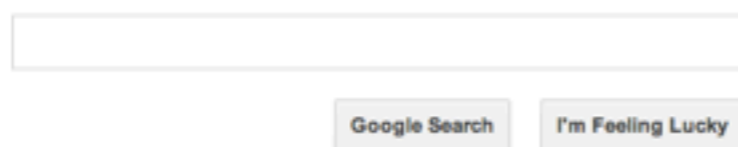
Connect with friends and the world around you on Facebook.

See photos and updates from friends in News Feed.

Share what's new in your life on your Timeline.



Data is replicated and partitioned across multiple nodes



Connect with friends and the world around you on Facebook.



See photos and updates from friends in News Feed.



Share what's new in your life on your Timeline.

Sign Up

It's free and always

First name

Email

Re-enter email

New password

Birthday

# Data centres across the world



Disaster-tolerance, minimising latency

# Data centres across the world



Disaster-tolerance, minimising latency

# Data centres across the world



Disaster-tolerance, minimising latency

# With thousands of machines inside



Load-balancing, fault-tolerance

# Replicas on mobile devices



Offline use

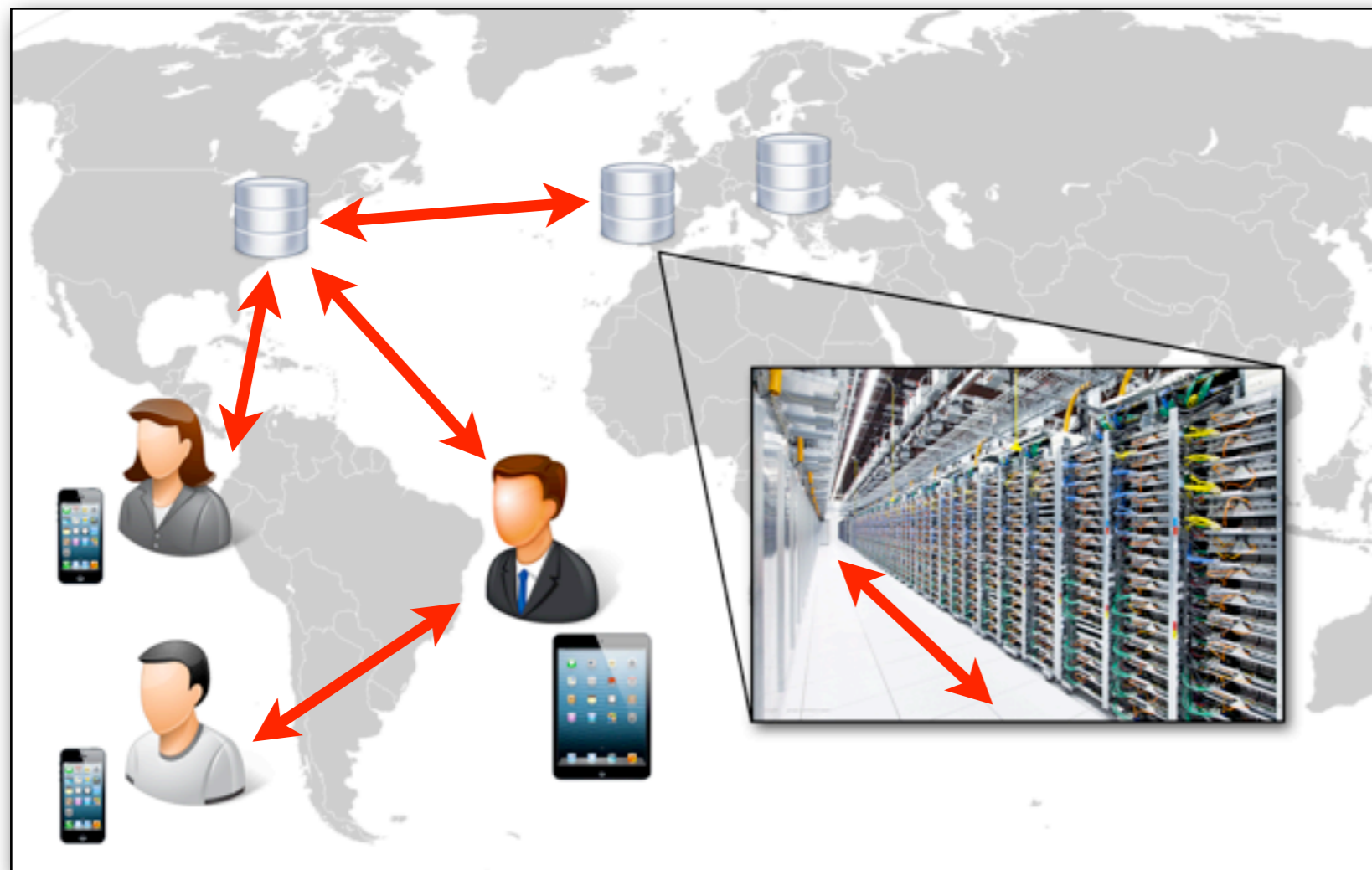




≈



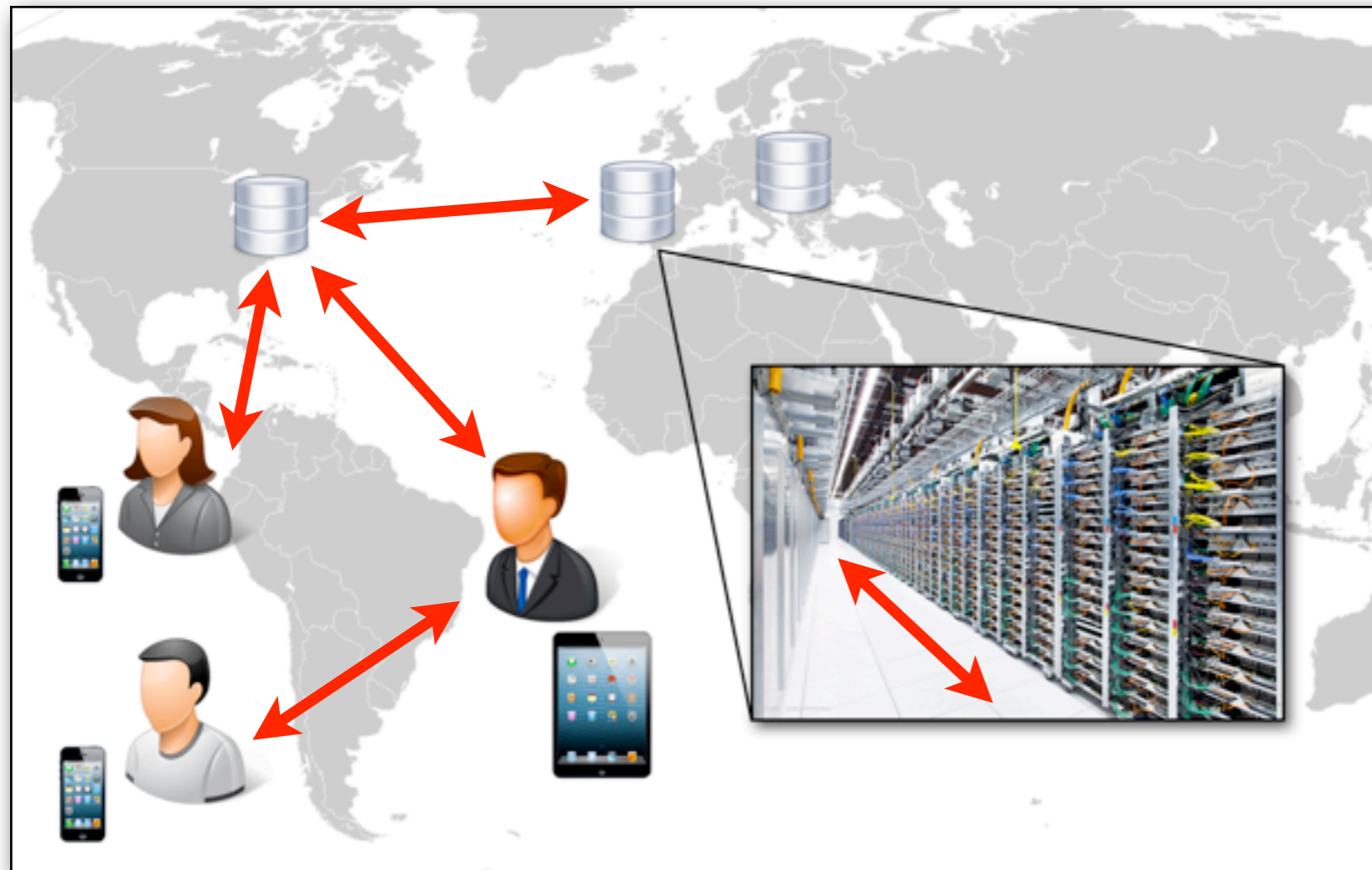
- **Strong consistency model:** the system behaves as if it processes requests serially on a centralised database - **linearizability, serializability**



≈



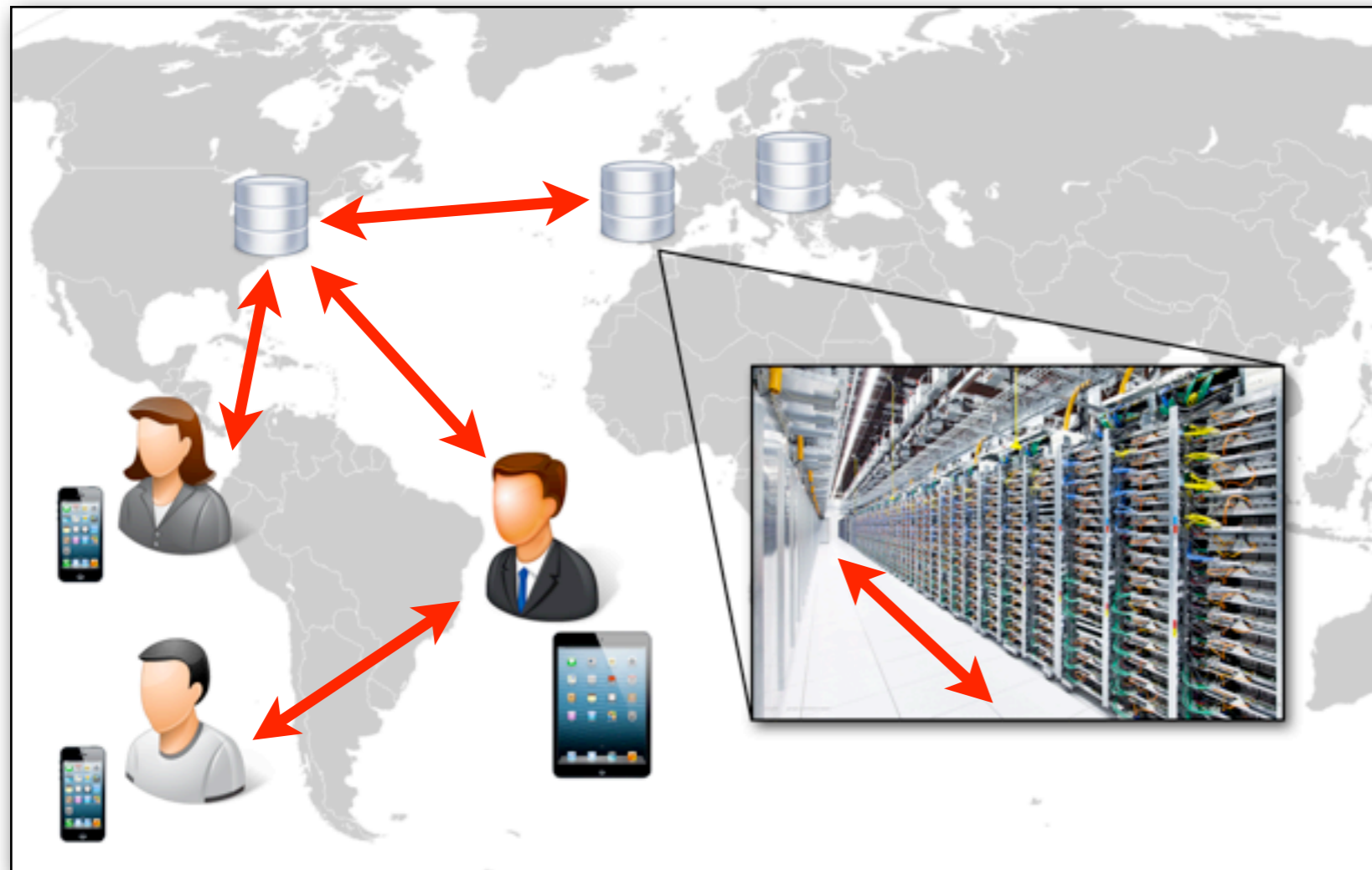
- **Strong consistency model:** the system behaves as if it processes requests serially on a centralised database - **linearizability, serializability**
- Requires **synchronisation:** contact other replicas when processing a request



≈



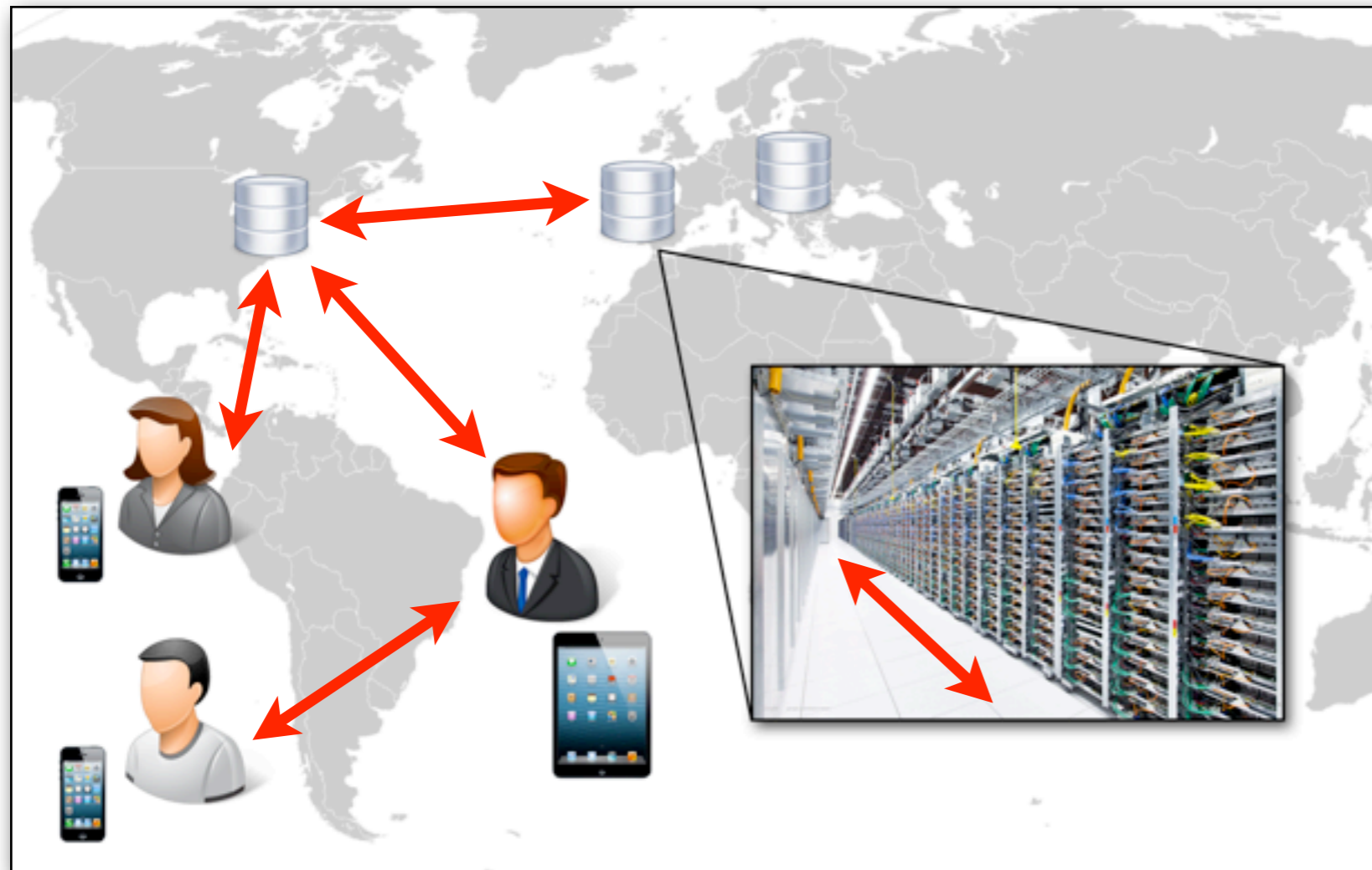
- Either strong **C**onsistency or **A**vailability in the presence of network **P**artitions [**CAP** theorem]



≈



- Either ~~strong Consistency~~ or Availability in the presence of network Partitions [CAP theorem]

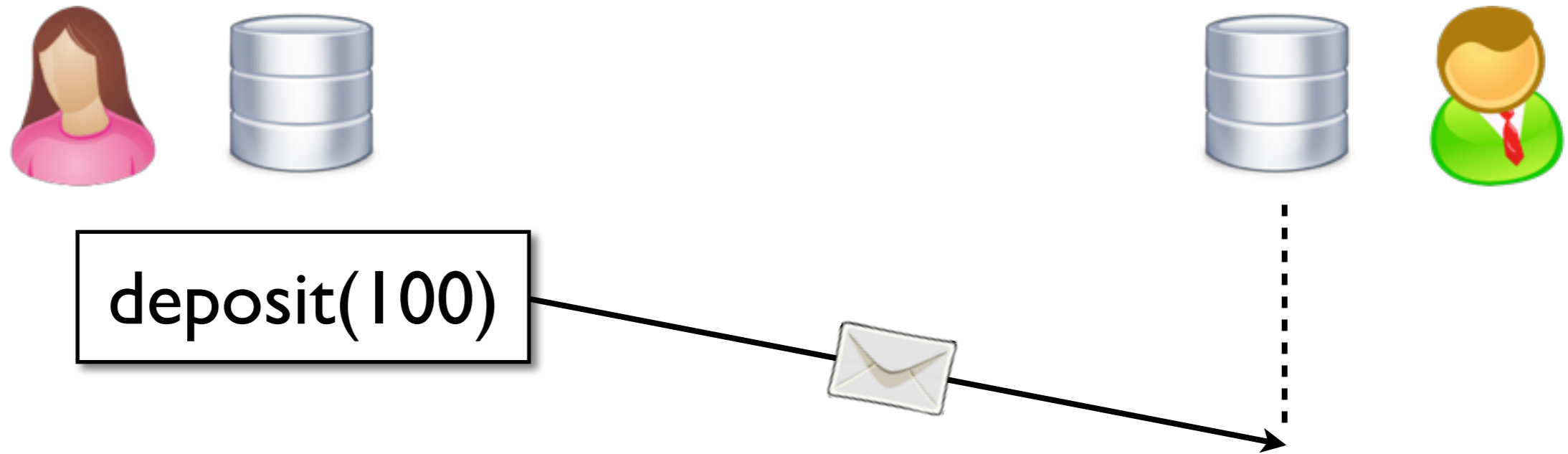


≈



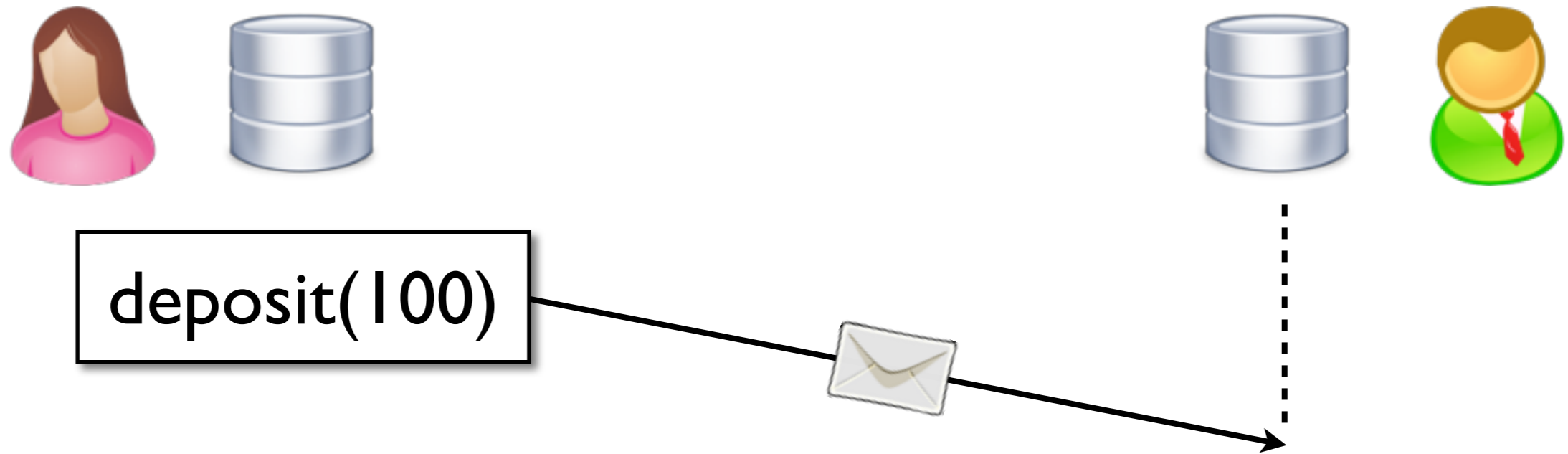
- Either ~~strong Consistency~~ or Availability in the presence of network Partitions [CAP theorem]
- Increased latency and resource consumption

# Relaxing synchronisation



Process an update locally, propagate effects to other replicas later

# Relaxing synchronisation



Process an update locally, propagate effects to other replicas later

- + Better scalability & availability
- **Weaken consistency**: deposit seen with a delay

# Anomalies ~ relaxed memory



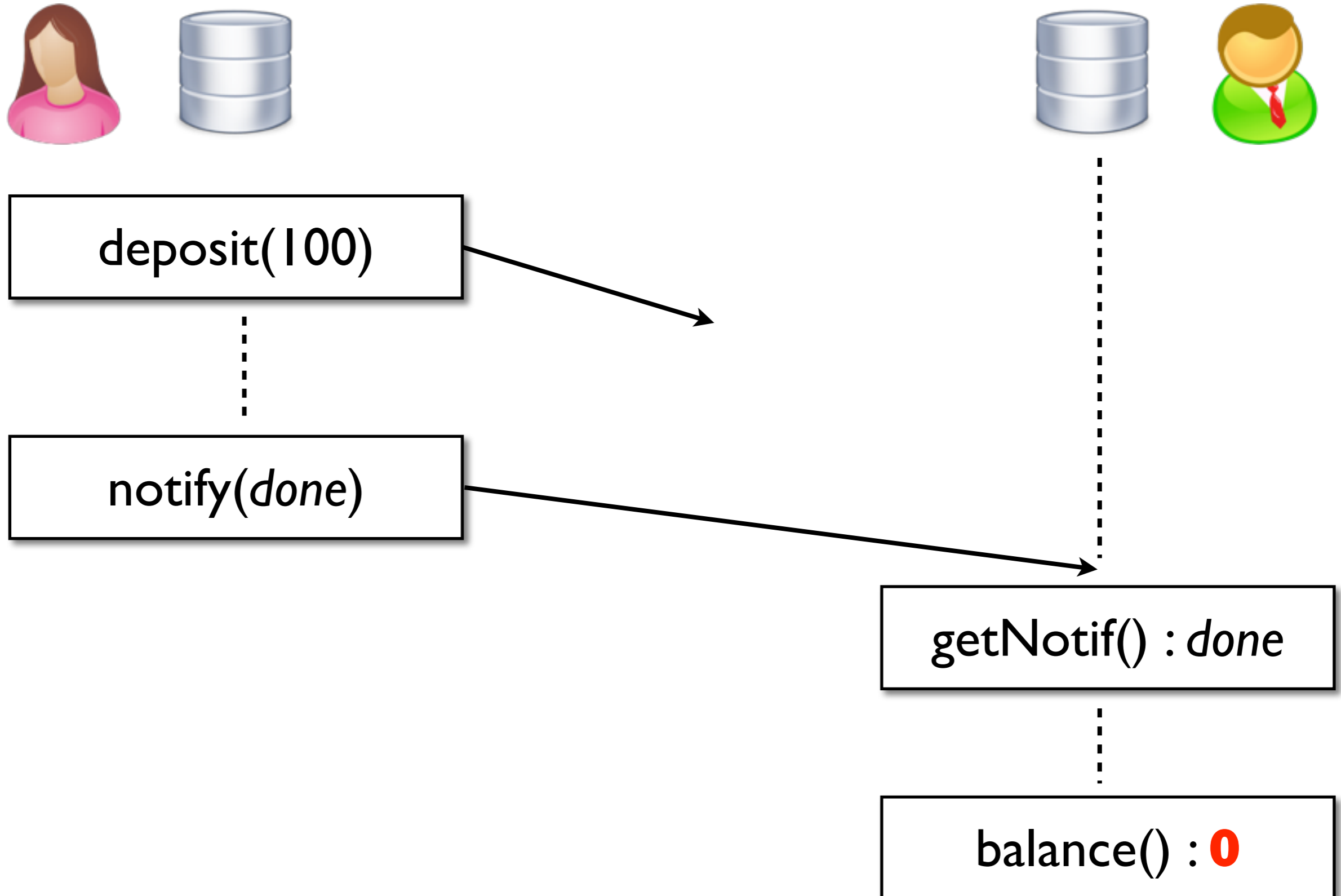
deposit(100)



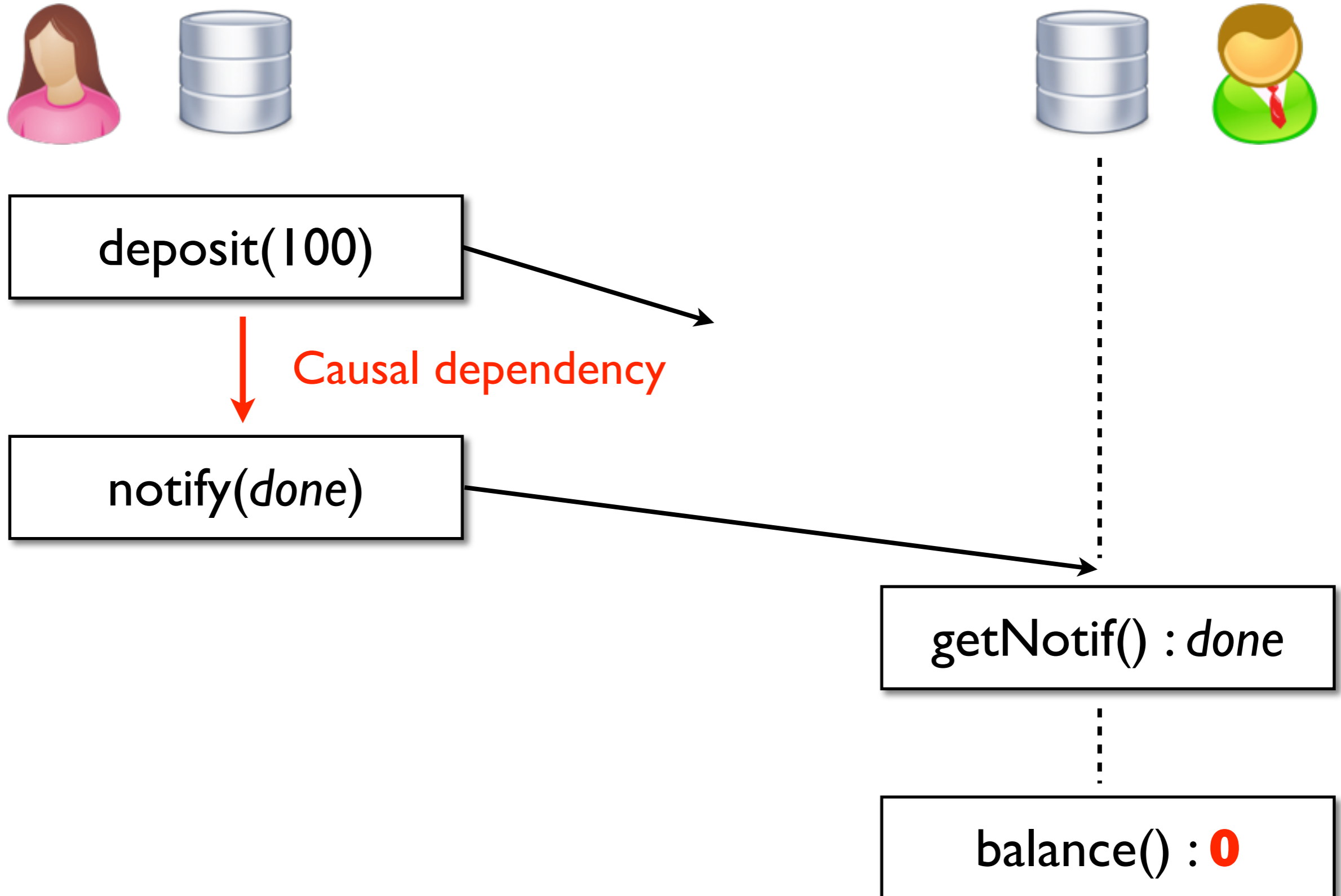
notify(*done*)



# Anomalies ~ relaxed memory



# Anomalies ~ relaxed memory





## Don't Settle for Eventual: Scalable Causal Consistency for Wide-Area Storage with COPS

Wyatt Lloyd\*, Michael J. Freedman\*, Michael Kaminsky†, and David G. Andersen‡  
\*Princeton University, †Intel Labs, ‡Carnegie Mellon University

## Highly Available Transactions: Virtues and Limitations

Peter Bailis, Aaron Davidson, Alan Fekete†, Ali Ghodsi, Joseph M. Hellerstein, Ion Stoica  
UC Berkeley and †University of Sydney

## Transactional storage for geo-replicated systems

Yair Sovran\* Russell Power\* Marcos K. Aguilera† Jinyang Li\*  
\*New York University †Microsoft Research Silicon Valley

## Eventually Consistent Transactions

Sebastian Burckhardt<sup>1</sup>, Daan Leijen<sup>1</sup>, Manuel Fähndrich<sup>1</sup>, and Mooly Sagiv<sup>2</sup>

<sup>1</sup> Microsoft Research

<sup>2</sup> Tel-Aviv University





## Don't Settle for Eventual: Scalable Causal Consistency for Wide-Area Storage with COPS

Wyatt Lloyd\*, Michael J. Freedman\*, Michael Kaminsky†, and David G. Andersen‡  
\*Princeton University, †Intel Labs, ‡Carnegie Mellon University

## Highly Available Transactions: Virtues and Limitations

Peter Bailis, Aaron Davidson, Alan Fekete†, Ali Ghodsi, Joseph M. Hellerstein, Ion Stoica  
UC Berkeley and †University of Sydney

## Transactional storage for geo-replicated systems

Yair Sovran\* Russell Power\* Marcos K. Aguilera† Jinyang Li\*  
\*New York University †Microsoft Research Silicon Valley

## Eventually Consistent Transactions

- A spectrum of consistency models:  
eventual consistency, causal consistency, ...
- Programming constructs for mitigating weakness



# Problem

- Is a given consistency model good for maintaining correctness in a given application?

*Does it weaken consistency too much, too little, just right?*

- How should programmers use the models and programming features correctly?

*No guidelines, patterns, static analysis tools, informal or inadequate specifications*

# Problem

- Is a given consistency model good for maintaining correctness in a given application?

*Does it weaken consistency too much to be just right?*

- How should programmers use

***“If no new updates are made to the database, then replicas will eventually reach a consistent state”***

practice

DOI:10.1145/1435417.1435433

**Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.**

BY WERNER VOGELS

## Eventually Consistent

AT THE FOUNDATION of Amazon's cloud computing are infrastructure services such as Amazon's S3 (Simple Storage Service), SimpleDB, and EC2 (Elastic Compute Cloud) that provide the resources for constructing Internet-scale computing platforms and a great variety of applications. The requirements placed on these infrastructure services are very strict; they need to score high marks in the areas of security, scalability, availability, performance, and cost-effectiveness, and they need to meet these requirements while serving millions of customers around the globe, continuously.

Under the covers these services are massive distributed systems that operate on a worldwide scale. This scale creates additional challenges, because when a system processes trillions and trillions of requests, events that normally have a low probability of occurrence are now guaranteed to happen and must be accounted for upfront in the design and architecture of the system. Given the worldwide scope of these systems, we use replication techniques ubiquitously to guarantee consistent performance and high availability. Although replication brings us closer to our goals, it cannot achieve them in a perfectly

transparent manner; under a number of conditions the customers of these services will be confronted with the consequences of using replication techniques inside the services.

One of the ways in which this manifests itself is in the type of data consistency that is provided, particularly when many widespread distributed systems provide an *eventual consistency* model in the context of data replication. When designing these large-scale systems at Amazon, we use a set of guiding principles and abstractions related to large-scale data replication and focus on the trade-offs between high availability and data consistency. Here, I present some of the relevant background that has informed our approach to delivering reliable distributed systems that must operate on a global scale. (An earlier version of this article appeared as a posting on the "All Things Distributed" Weblog and was greatly improved with the help of its readers.)

**Historical Perspective**

In an ideal world there would be only one consistency model: when an update is made all observers would see that update. The first time this surfaced as difficult to achieve was in the database systems of the late 1970s. The best "period piece" on this topic is "Notes on Distributed Databases" by Bruce Lindsay et al. It lays out the fundamental principles for database replication and discusses a number of techniques that deal with achieving consistency. Many of these techniques try to achieve *distribution transparency*—that is, to the user of the system it appears as if there is only one system instead of a number of collaborating systems. Many systems during this time took the approach that it was better to fall the complete system than to break this transparency.<sup>2</sup>

In the mid-1990s, with the rise of larger Internet systems, these practices were revisited. At that time people began to consider the idea that availability was perhaps the most impor-

# Problem

- Is a given consistency model good for maintaining correctness in a given application?

*Does it weaken consistency too much to be just right?*

- How should programmers use

***“If no new updates are made to the database, then replicas will eventually reach a consistent state”***

practice

DOI:10.1145/1435417.1435433

**Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.**

BY WERNER VOGELS

## Eventually Consistent

AT THE FOUNDATION of Amazon's cloud computing are infrastructure services such as Amazon's S3 (Simple Storage Service), SimpleDB, and EC2 (Elastic Compute Cloud) that provide the resources for constructing Internet-scale computing platforms and a great variety of applications. The requirements placed on these infrastructure services are very strict; they need to score high marks in the areas of security, scalability, availability, performance, and cost-effectiveness, and they need to meet these requirements while serving millions of customers around the globe, continuously.

Under the covers these services are massive distributed systems that operate on a worldwide scale. This scale creates additional challenges, because when a system processes trillions and trillions of requests, events that normally have a low probability of occurrence are now guaranteed to happen and must be accounted for upfront in the design and architecture of the system. Given the worldwide scope of these systems, we use replication techniques ubiquitously to guarantee consistent performance and high availability. Although replication brings us closer to our goals, it cannot achieve them in a perfectly

transparent manner; under a number of conditions the customers of these services will be confronted with the consequences of using replication techniques inside the services.

One of the ways in which this manifests itself is in the type of data consistency that is provided, particularly when many widespread distributed systems provide an *eventual consistency* model in the context of data replication. When designing these large-scale systems at Amazon, we use a set of guiding principles and abstractions related to large-scale data replication and focus on the trade-offs between high availability and data consistency. Here, I present some of the relevant background that has informed our approach to delivering reliable distributed systems that must operate on a global scale. (An earlier version of this article appeared as a posting on the "All Things Distributed" Weblog and was greatly improved with the help of its readers.)

**Historical Perspective**

In an ideal world there would be only one consistency model: when an update is made all observers would see that update. The first time this surfaced as difficult to achieve was in the database systems of the late 1970s. The best "period piece" on this topic is "Notes on Distributed Databases" by Bruce Lindsay et al. It lays out the fundamental principles for database replication and discusses a number of techniques that deal with achieving consistency. Many of these techniques try to achieve *distribution transparency*—that is, to the user of the system it appears as if there is only one system instead of a number of collaborating systems. Many systems during this time took the approach that it was better to fall the complete system than to break this transparency.<sup>2</sup>

In the mid-1990s, with the rise of larger Internet systems, these practices were revisited. At that time people began to consider the idea that availability was perhaps the most impor-

# Problem

- Is a given consistency model good for maintaining correctness in a given application?

*Does it weaken consistency too much to be just right?*

- How should programmers use

*“If no new updates are made to the database, then replicas will eventually reach **a consistent state**”*

practice

DOI:10.1145/1435417.1435433

**Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.**

BY WERNER VOGELS

## Eventually Consistent

AT THE FOUNDATION of Amazon's cloud computing are infrastructure services such as Amazon's S3 (Simple Storage Service), SimpleDB, and EC2 (Elastic Compute Cloud) that provide the resources for constructing Internet-scale computing platforms and a great variety of applications. The requirements placed on these infrastructure services are very strict; they need to score high marks in the areas of security, scalability, availability, performance, and cost-effectiveness, and they need to meet these requirements while serving millions of customers around the globe, continuously.

Under the covers these services are massive distributed systems that operate on a worldwide scale. This scale creates additional challenges, because when a system processes trillions and trillions of requests, events that normally have a low probability of occurrence are now guaranteed to happen and must be accounted for upfront in the design and architecture of the system. Given the worldwide scope of these systems, we use replication techniques ubiquitously to guarantee consistent performance and high availability. Although replication brings us closer to our goals, it cannot achieve them in a perfectly

transparent manner; under a number of conditions the customers of these services will be confronted with the consequences of using replication techniques inside the services.

One of the ways in which this manifests itself is in the type of data consistency that is provided, particularly when many widespread distributed systems provide an *eventual consistency* model in the context of data replication. When designing these large-scale systems at Amazon, we use a set of guiding principles and abstractions related to large-scale data replication and focus on the trade-offs between high availability and data consistency. Here, I present some of the relevant background that has informed our approach to delivering reliable distributed systems that must operate on a global scale. (An earlier version of this article appeared as a posting on the "All Things Distributed" Weblog and was greatly improved with the help of its readers.)

**Historical Perspective**

In an ideal world there would be only one consistency model: when an update is made all observers would see that update. The first time this surfaced as difficult to achieve was in the database systems of the late 1970s. The best "period piece" on this topic is "Notes on Distributed Databases" by Bruce Lindsay et al. It lays out the fundamental principles for database replication and discusses a number of techniques that deal with achieving consistency. Many of these techniques try to achieve *distribution transparency*—that is, to the user of the system it appears as if there is only one system instead of a number of collaborating systems. Many systems during this time took the approach that it was better to fall the complete system than to break this transparency.<sup>2</sup>

In the mid-1990s, with the rise of larger Internet systems, these practices were revisited. At that time people began to consider the idea that availability was perhaps the most impor-



# TOWARDS A CLOUD COMPUTING RESEARCH AGENDA

Ken Birman, Gregory Chockler, Robbert van Renesse

This particular example is a good one because, as we'll see shortly, if there was a single overarching theme within the keynote talks, it turns out to be that **strong synchronization** of the sort provided by a locking service **must be avoided like the plague**. This doesn't diminish the need for a tool like Chubby; when locking actually can't be avoided, one wants a reliable, standard, provably correct

# TOWARDS A CLOUD COMPUTING RESEARCH AGENDA

Ken Birman, Gregory Chockler, Robbert van Renesse

This particular example is a good one because, as we'll see shortly, if there was a single overarching theme within the keynote talks, it turns out to be that **strong synchronization** of the sort provided by a locking service **must be avoided like the plague**. This doesn't diminish the need for a tool like Chubby; when locking actually can't be avoided, one wants a reliable, standard, provably correct

## F1: A Distributed SQL Database That Scales

Jeff Shute  
Chad Whipkey  
David Menestrina

Radek Vingralek  
Eric Rollins  
Stephan Ellner  
Traian Stancescu

Bart Samwel  
Mircea Oancea  
John Cieslewicz  
Himani Apte

Ben Handy  
Kyle Littlefield  
Ian Rae\*

Google, Inc.

\*University of Wisconsin-Madison

### ABSTRACT

F1 is a distributed relational database system built at Google to support the AdWords business. F1 is a hybrid database that combines high availability, the scalability of NoSQL systems like Bigtable, and the consistency and us-

consistent and correct data.

Designing applications to cope with concurrency anomalies in their data is very error-prone, time-consuming, and ultimately not worth the performance gains.

# TOWARDS A CLOUD COMPUTING RESEARCH AGENDA

Ken Birman, Gregory Chockler, Robbert van Renesse

This particular example is a good one because, as we'll see shortly, if there was a single overarching theme within the keynote talks, it turns out to be that **strong synchronization** of the sort provided by a locking service **must be avoided like the plague**. This doesn't diminish the need for a tool like Chubby: when locking actually can't be avoided, one wants a reliable, standard, provably correct,

Need a **rigorous approach**: programming models and static analysis tools that allow relaxing synchronisation without compromising correctness

David Menestrina

Stephan Ellner

John Cieslewicz

Ian Rae\*

Traian Stancescu

Himani Apte

Google, Inc.

\*University of Wisconsin-Madison

## ABSTRACT

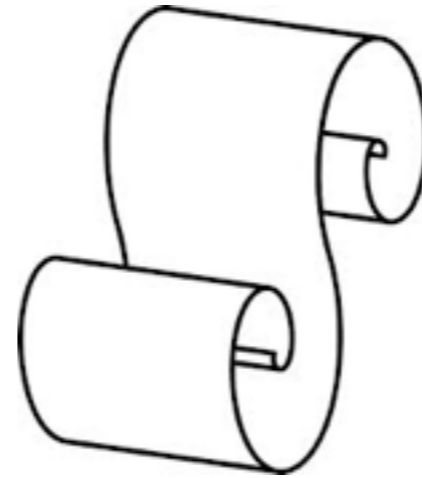
F1 is a distributed relational database system built at Google to support the AdWords business. F1 is a hybrid database that combines high availability, the scalability of NoSQL systems like Bigtable, and the consistency and us-

consistent and correct data.

Designing applications to cope with concurrency anomalies in their data is very error-prone, time-consuming, and ultimately not worth the performance gains.

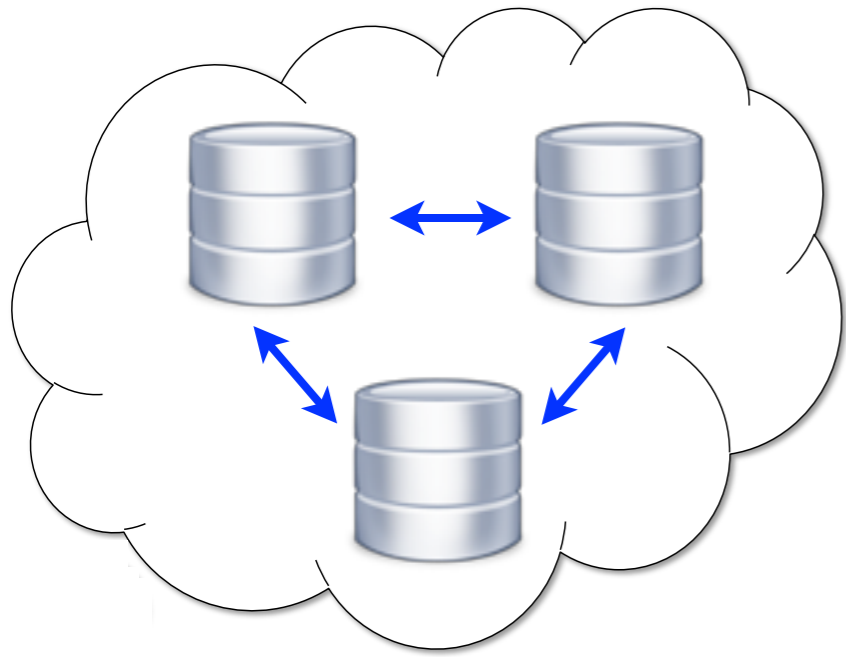
# Research agenda

# Research agenda

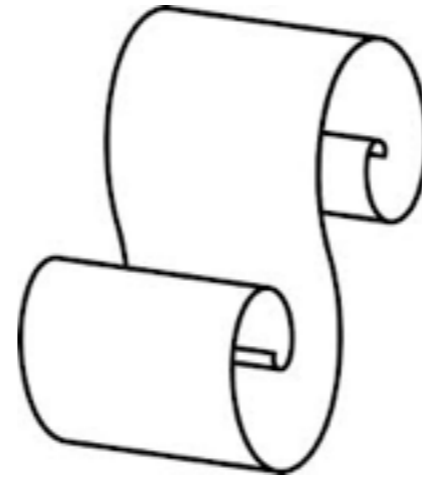


Formal semantics  
specification

# Research agenda

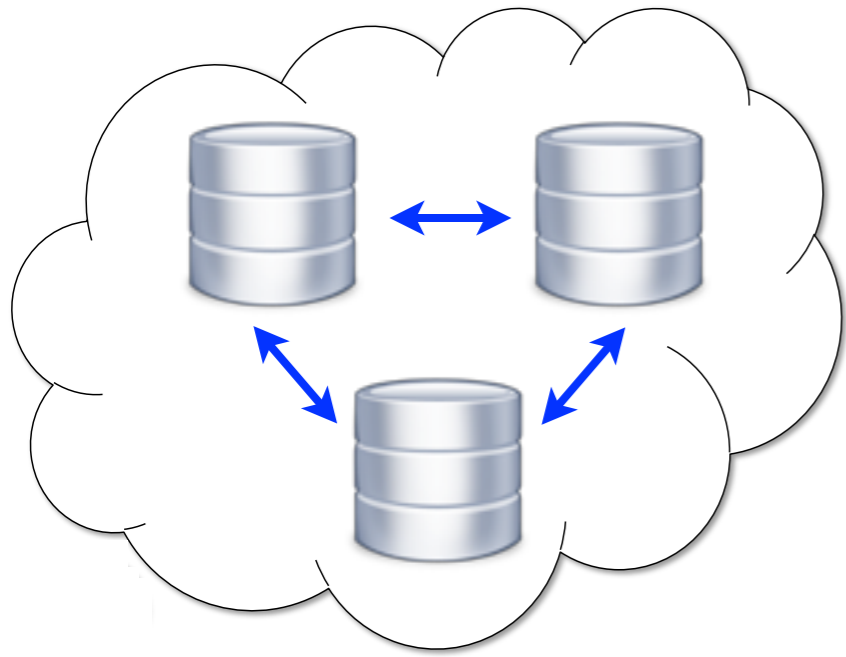


**Reasoning about DB  
implementations**

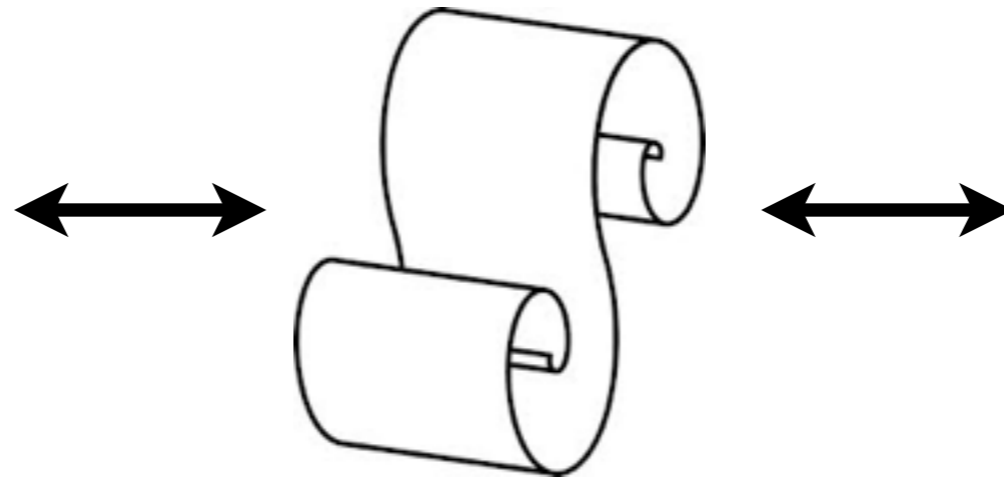


**Formal semantics  
specification**

# Research agenda



Reasoning about DB implementations

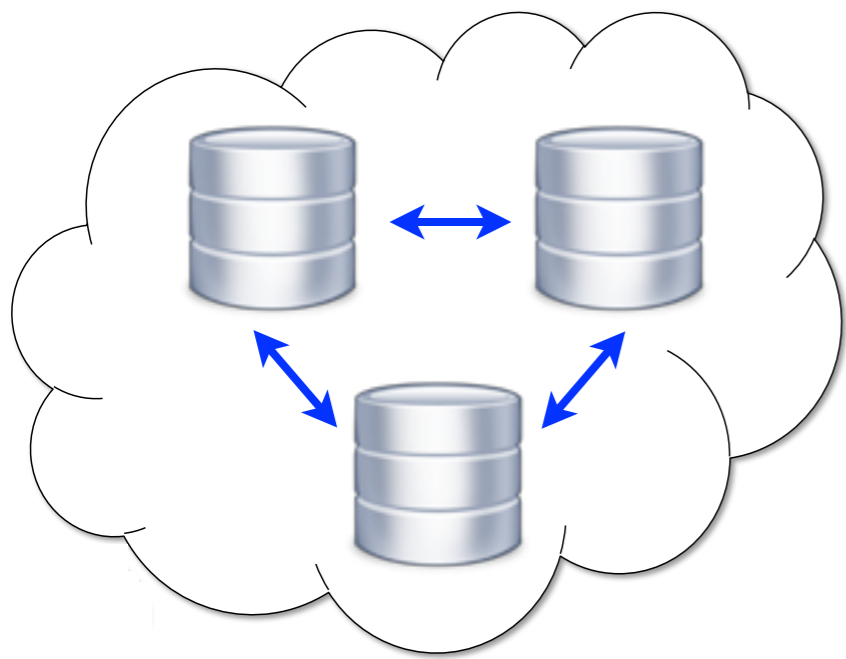


Formal semantics specification

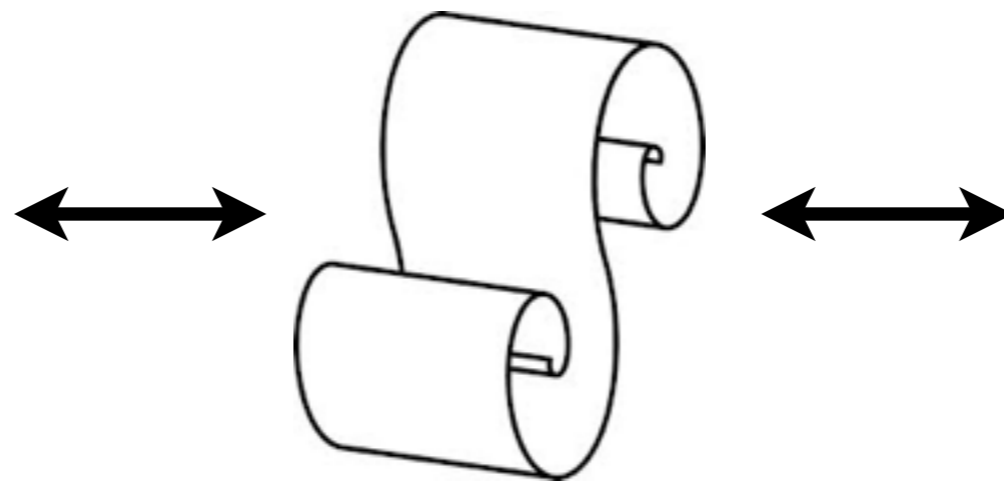


Reasoning about applications

# Research agenda



Reasoning about DB implementations



Formal semantics specification



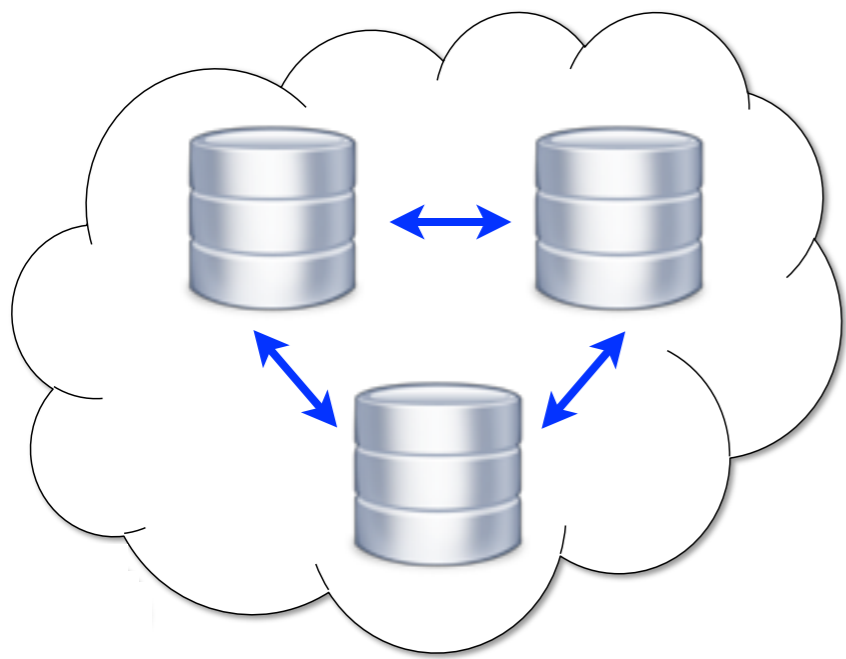
Reasoning about applications



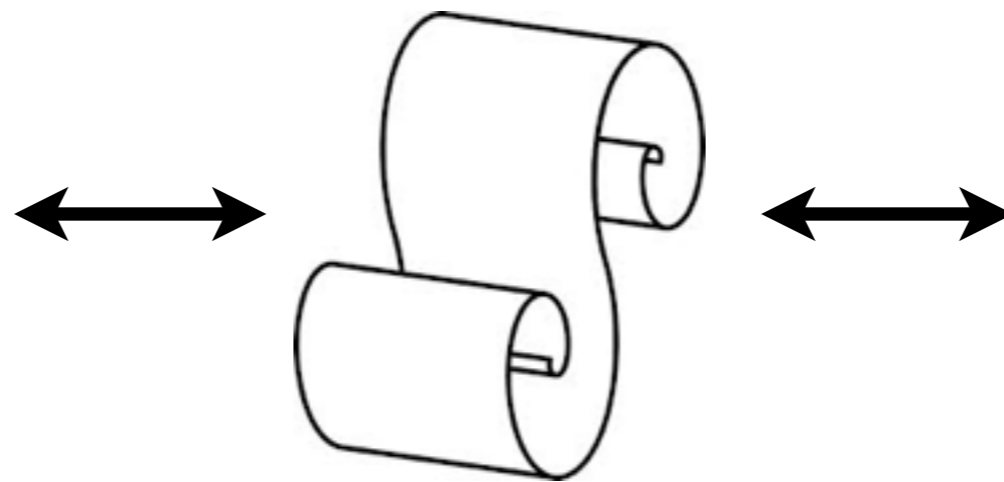
Improve DB programming models and implementations



# Research agenda



Reasoning about DB implementations



Formal semantics specification



Reasoning about applications

Joint work with Hongseok Yang (Oxford), Carla Ferreira (U Nova Lisboa), Mahsa Najafzadeh, Marc Shapiro (INRIA)

# Synchronisation can be necessary



balance = 100

balance  $\geq$  0



balance = 100

# Synchronisation can be necessary



balance = 100

withdraw(100) : ✓

balance = 0

balance  $\geq$  0



balance = 100

withdraw(100) : ✓

balance = 0

# Synchronisation can be necessary



balance = 100

balance  $\geq$  0



balance = 100

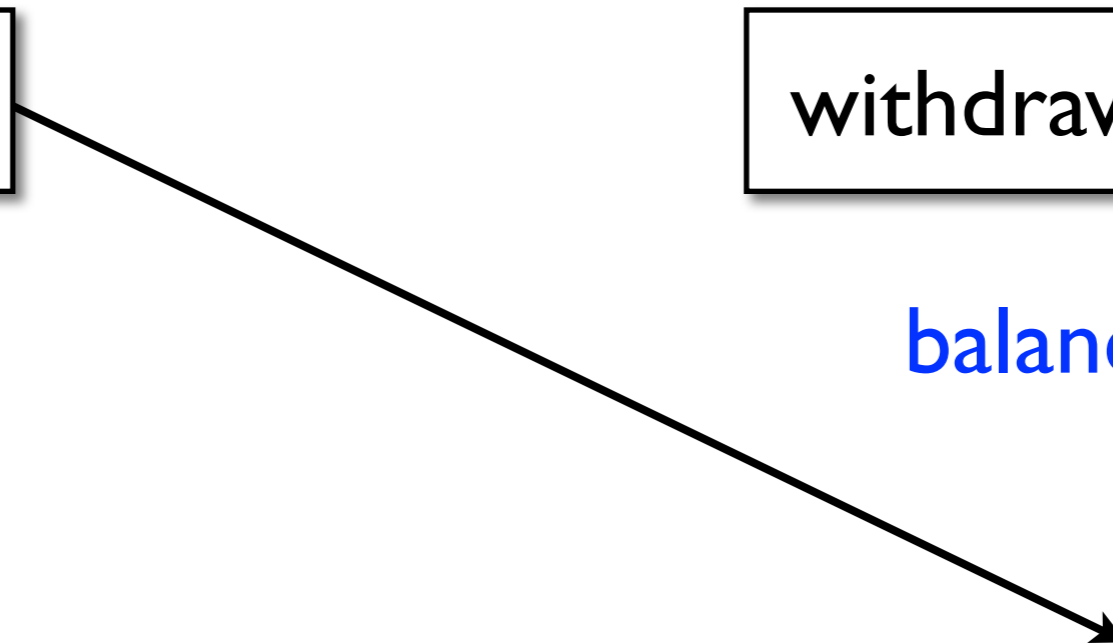
withdraw(100) : ✓

balance = 0

withdraw(100) : ✓

balance = 0

balance = -100



# Synchronisation can be necessary



balance = 100

balance  $\geq$  0



balance = 100

withdraw(100) : ✓

withdraw(100) : ✓

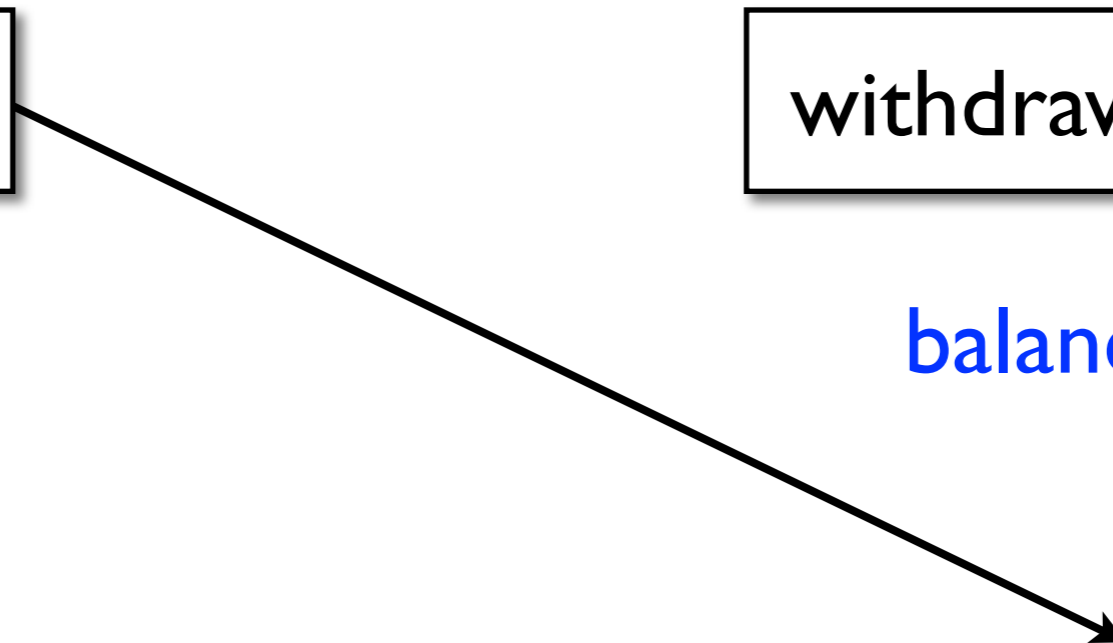
balance = 0

balance = 0



deposit(100)

balance = -100



# Consistency choices

- Choose consistency level for each operation:
  - ▶ Withdrawals strongly consistent
  - ▶ Deposits eventually consistent
- Pay for stronger semantics with latency, possible unavailability and money
- Hard to figure out the minimum consistency necessary to maintain correctness -  
proof rule and tool

# Consistency model

Generic model - not implemented, but can encode many existing models that are

- **Causal consistency** as a baseline:  
observe an update → observe the updates it depends on
- A construct for **strengthening consistency** on demand

# Operation semantics



$\sigma$



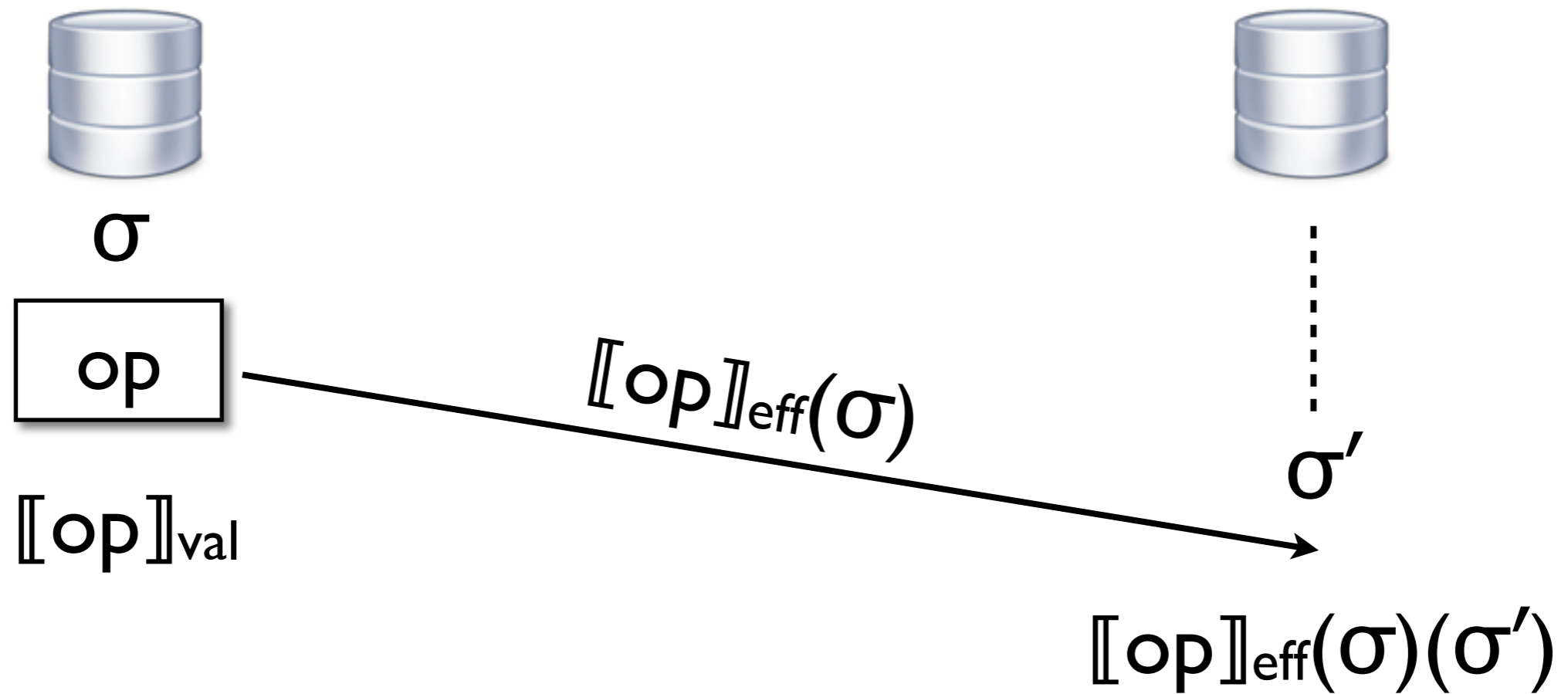
$\llbracket op \rrbracket_{val}$

Replica states:  $\sigma \in \text{State}$

Return value:  $\llbracket op \rrbracket_{val} \in \text{State} \rightarrow \text{Value}$



# Operation semantics

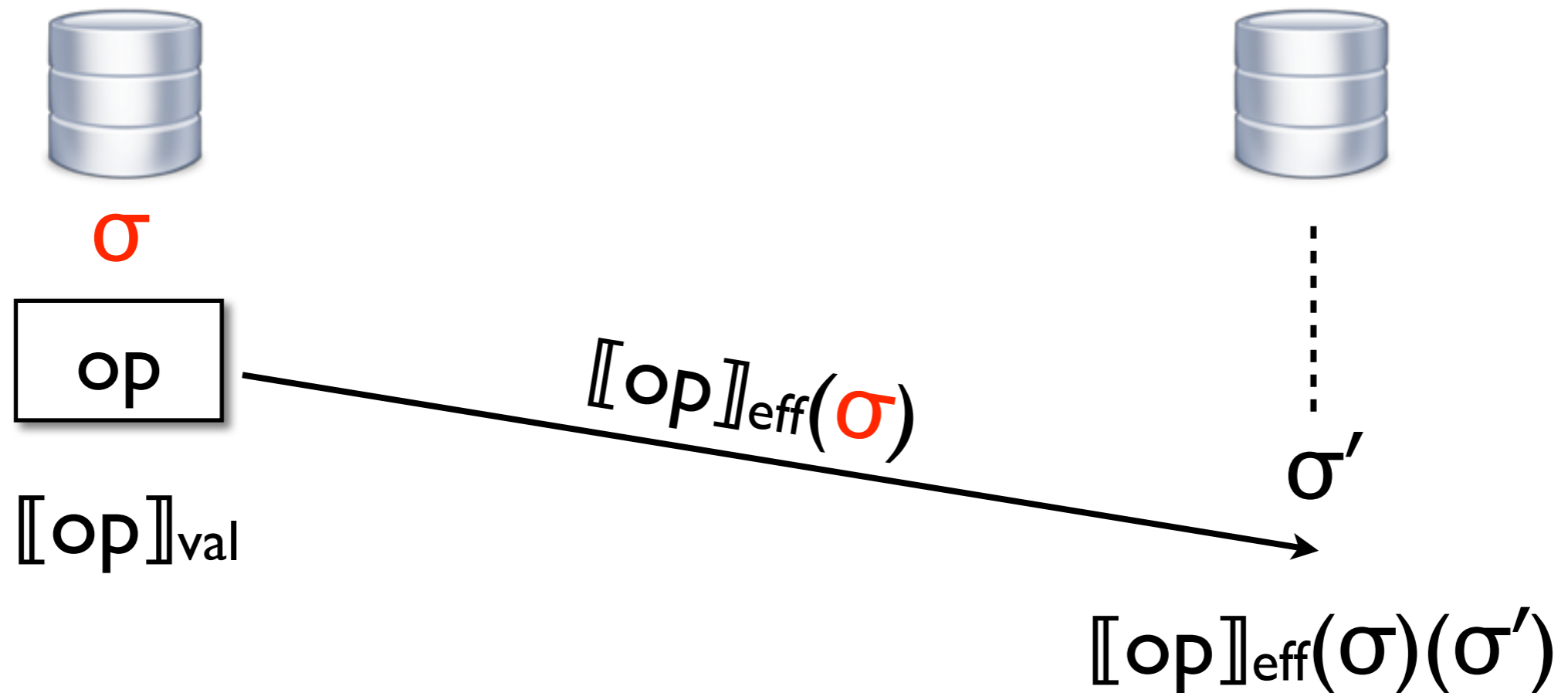


Replica states:  $\sigma \in \text{State}$

Return value:  $[[op]]_{val} \in \text{State} \rightarrow \text{Value}$

Effector:  $[[op]]_{eff} \in \text{State} \rightarrow (\text{State} \rightarrow \text{State})$

# Operation semantics

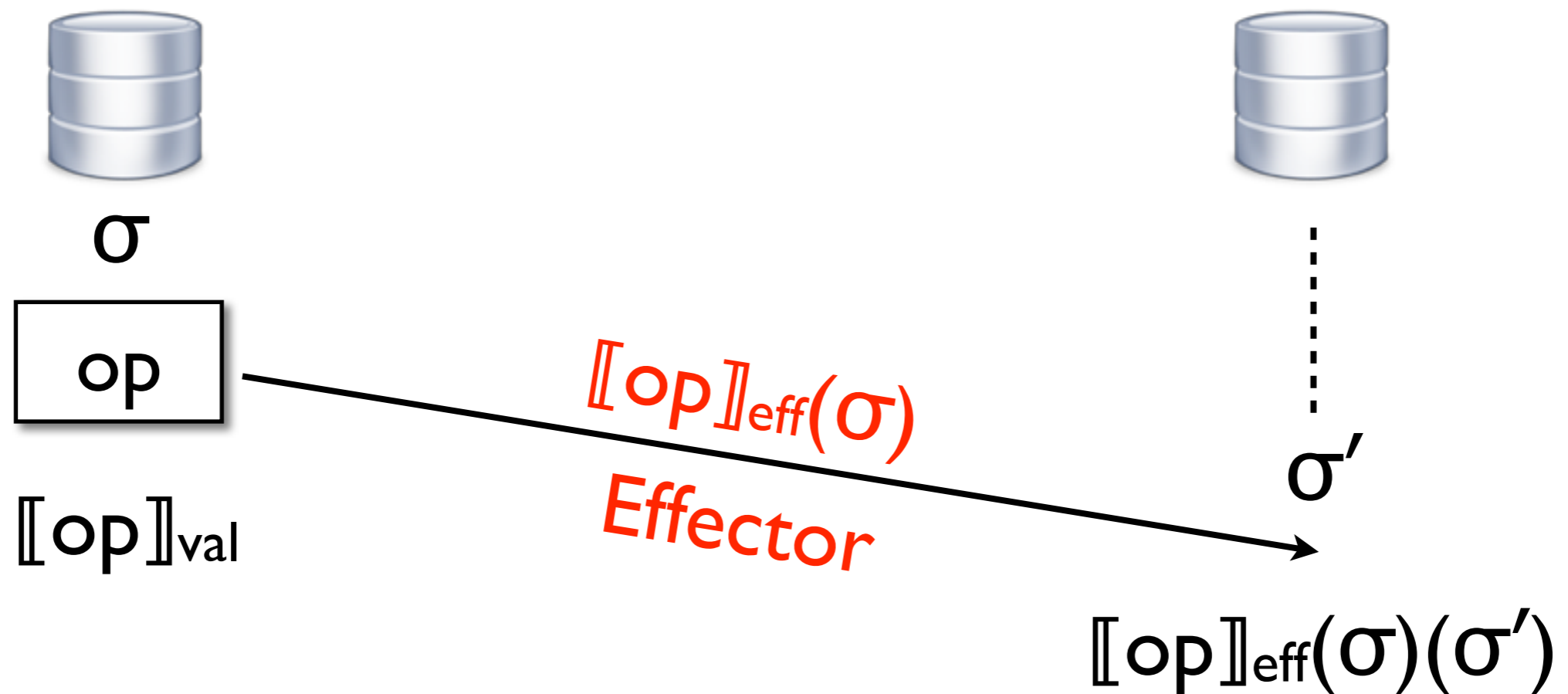


Replica states:  $\sigma \in \text{State}$

Return value:  $[[op]]_{val} \in \text{State} \rightarrow \text{Value}$

Effector:  $[[op]]_{eff} \in \text{State} \rightarrow (\text{State} \rightarrow \text{State})$

# Operation semantics

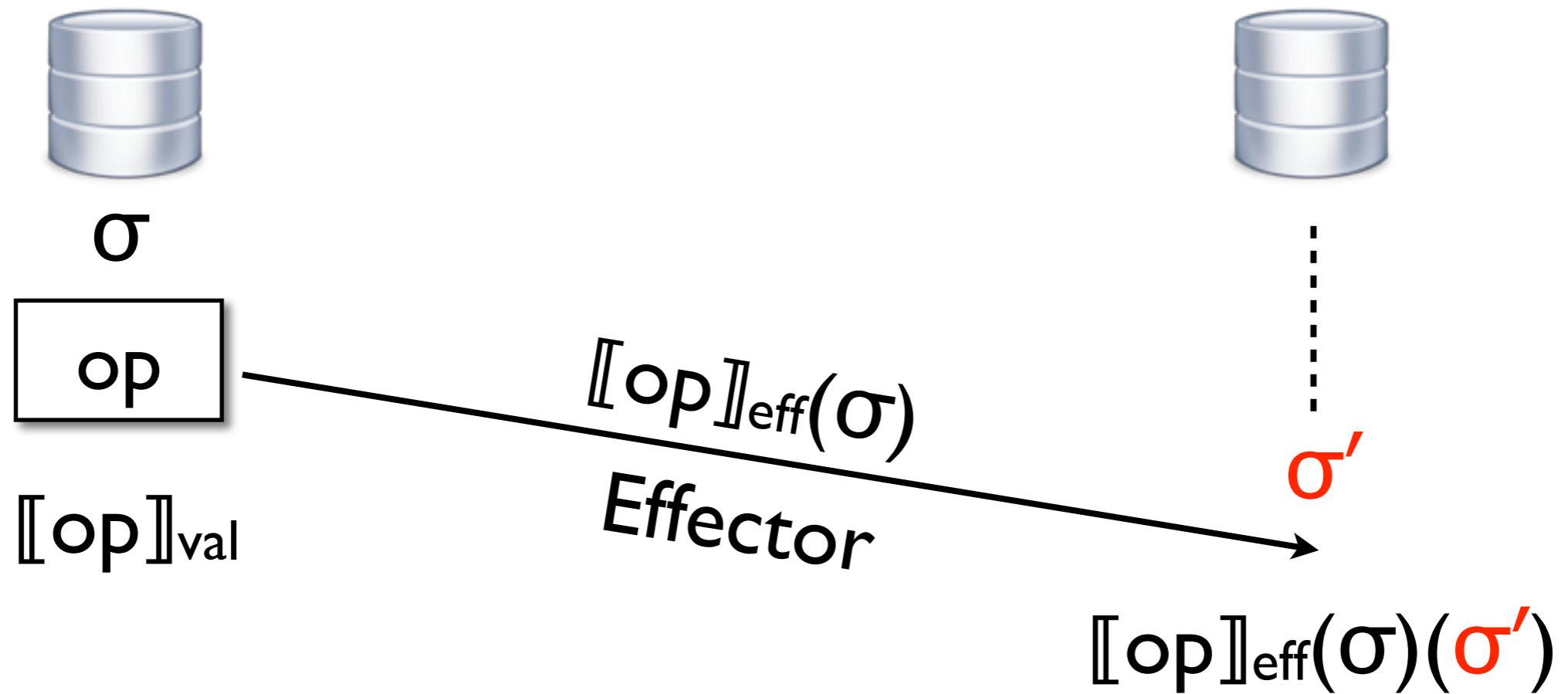


Replica states:  $\sigma \in \text{State}$

Return value:  $[[op]]_{val} \in \text{State} \rightarrow \text{Value}$

Effector:  $[[op]]_{eff} \in \text{State} \rightarrow (\text{State} \rightarrow \text{State})$

# Operation semantics

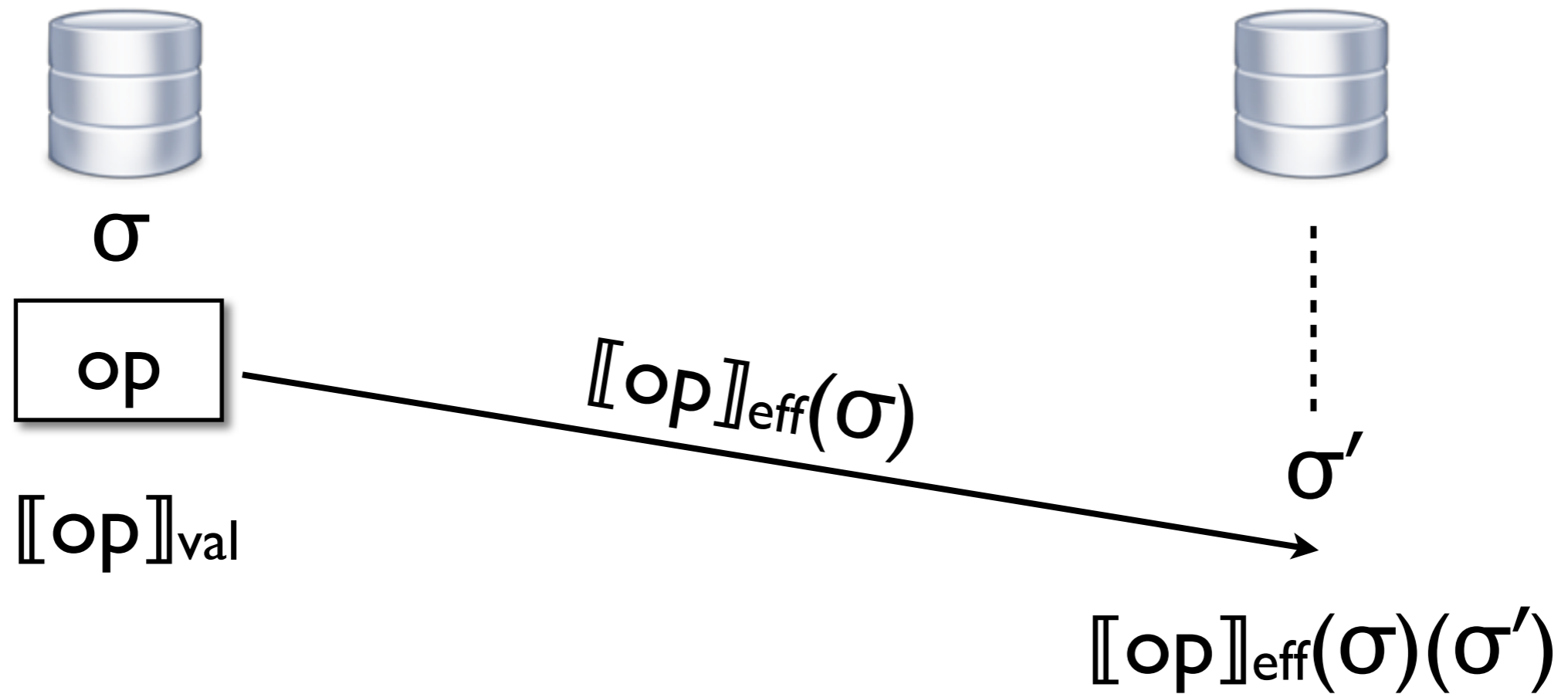


Replica states:  $\sigma \in \text{State}$

Return value:  $[[op]]_{\text{val}} \in \text{State} \rightarrow \text{Value}$

Effector:  $[[op]]_{\text{eff}} \in \text{State} \rightarrow (\text{State} \rightarrow \text{State})$

# Operation semantics

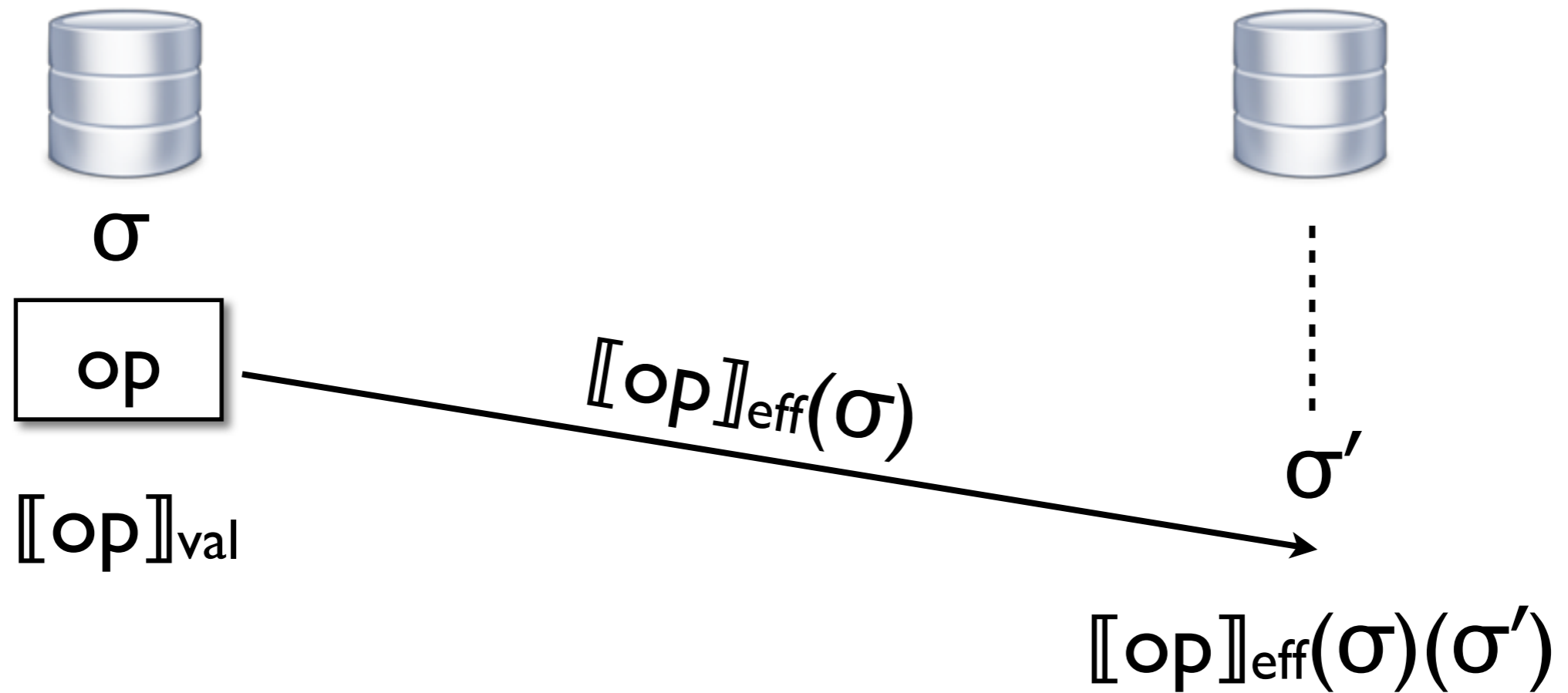


State =  $Z$

$$\llbracket \text{balance}() \rrbracket_{val}(\sigma) = \sigma$$

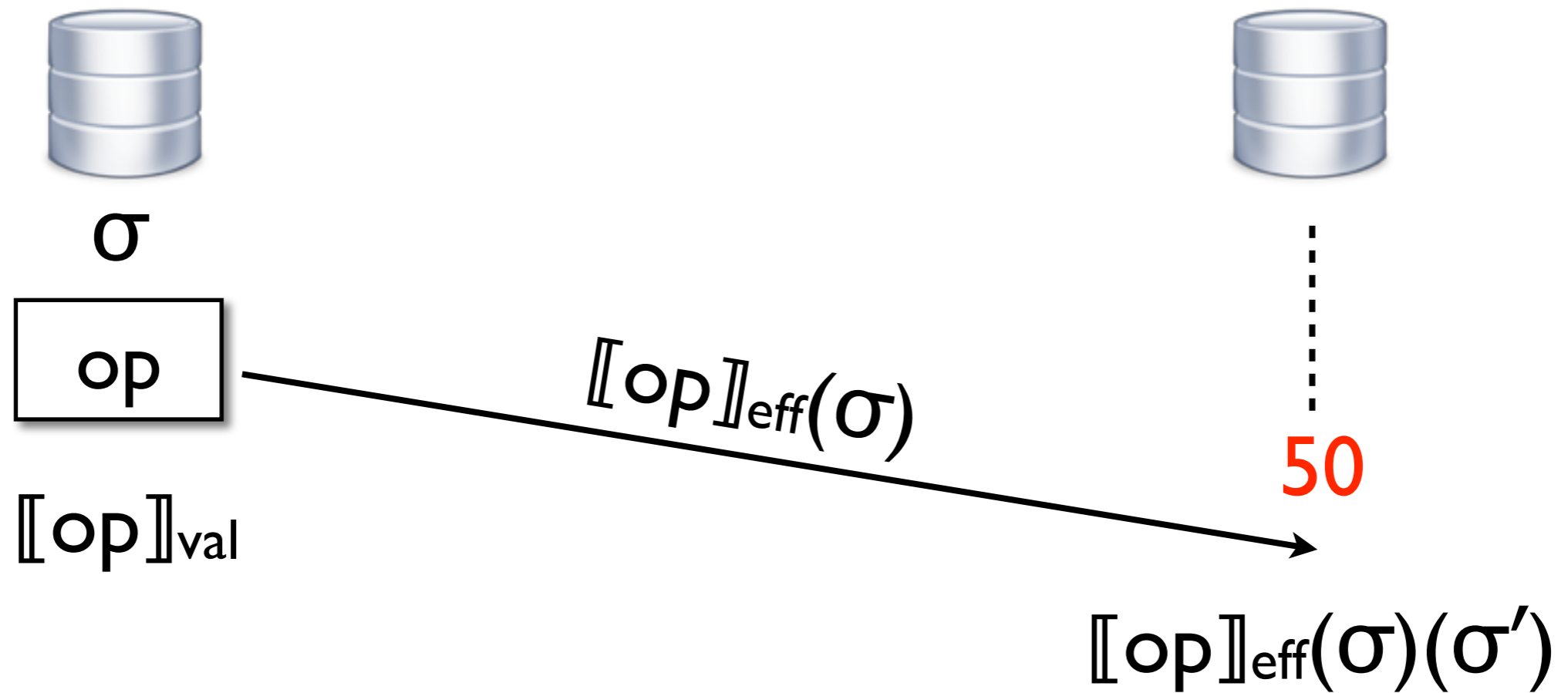
$$\llbracket \text{balance}() \rrbracket_{eff}(\sigma) = \lambda \sigma. \sigma$$

# Operation semantics



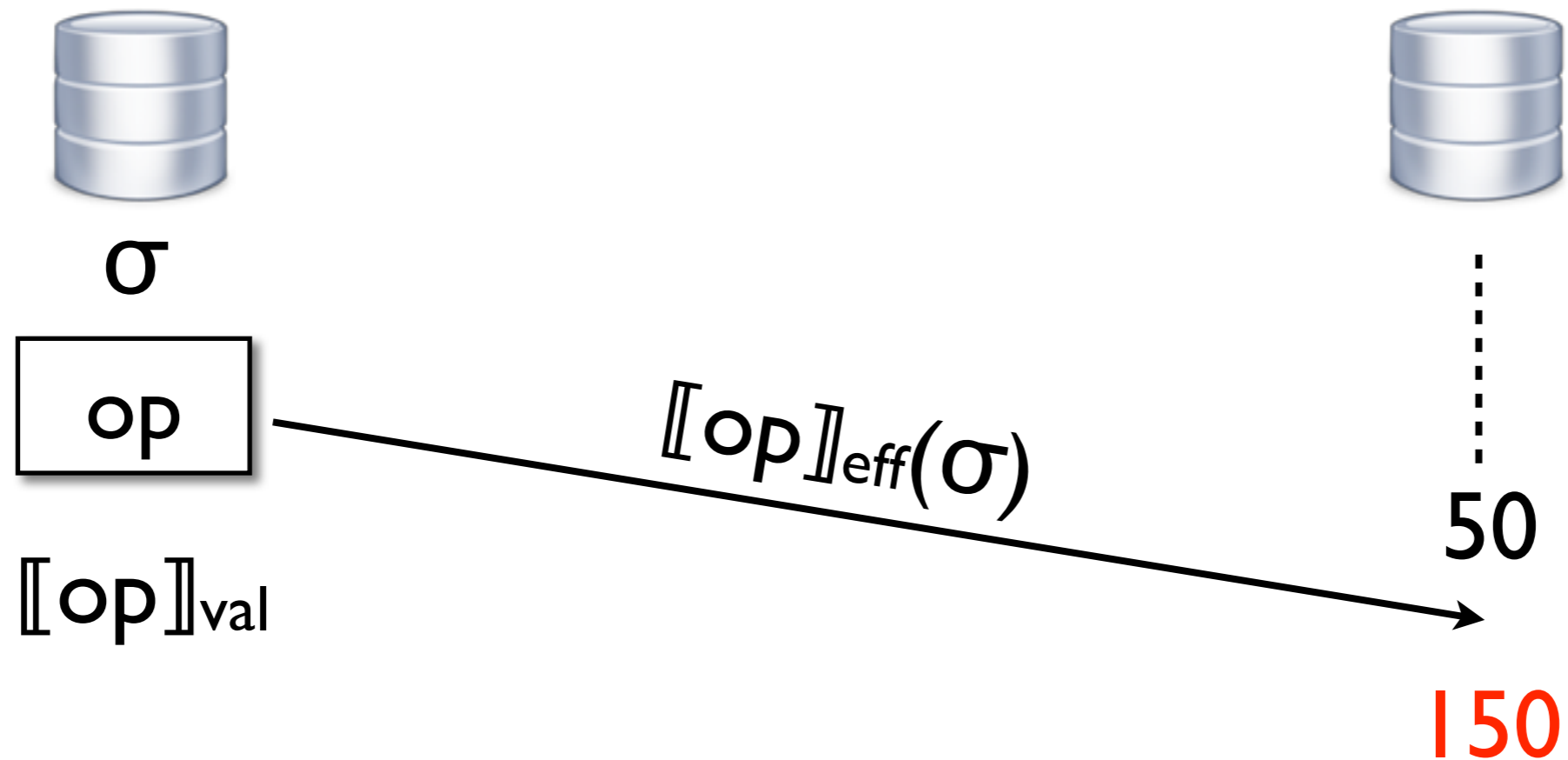
$$[[deposit(100)]]_{eff}(\sigma) = \lambda\sigma'. (\sigma' + 100)$$

# Operation semantics



$$\llbracket \text{deposit}(100) \rrbracket_{\text{eff}}(\sigma) = \lambda \sigma'. (\sigma' + 100)$$

# Operation semantics



$$[[\text{deposit}(100)]]_{eff}(\sigma) = \lambda\sigma'. (\sigma' + 100)$$



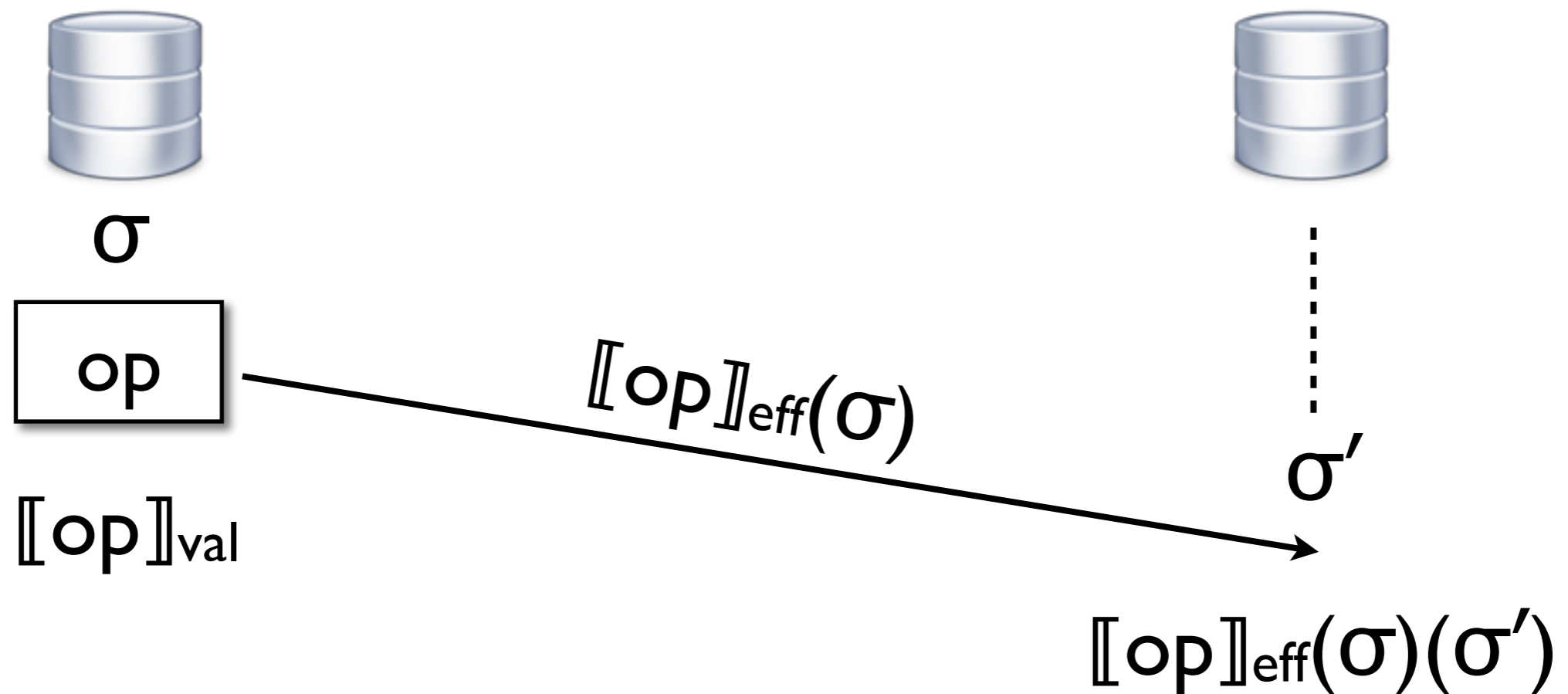
# Ensuring eventual consistency

- **Effectors have to commute:**

$$\forall op_1, op_2, \sigma_1, \sigma_2. \llbracket op_1 \rrbracket_{\text{eff}}(\sigma_1) ; \llbracket op_2 \rrbracket_{\text{eff}}(\sigma_2) = \llbracket op_2 \rrbracket_{\text{eff}}(\sigma_2) ; \llbracket op_1 \rrbracket_{\text{eff}}(\sigma_1)$$

- **Eventual consistency:** replicas receiving the same messages in different orders end up in the same state
- **Replicated data types** [Shapiro<sup>+</sup> 2011]: ready-made commutative implementations

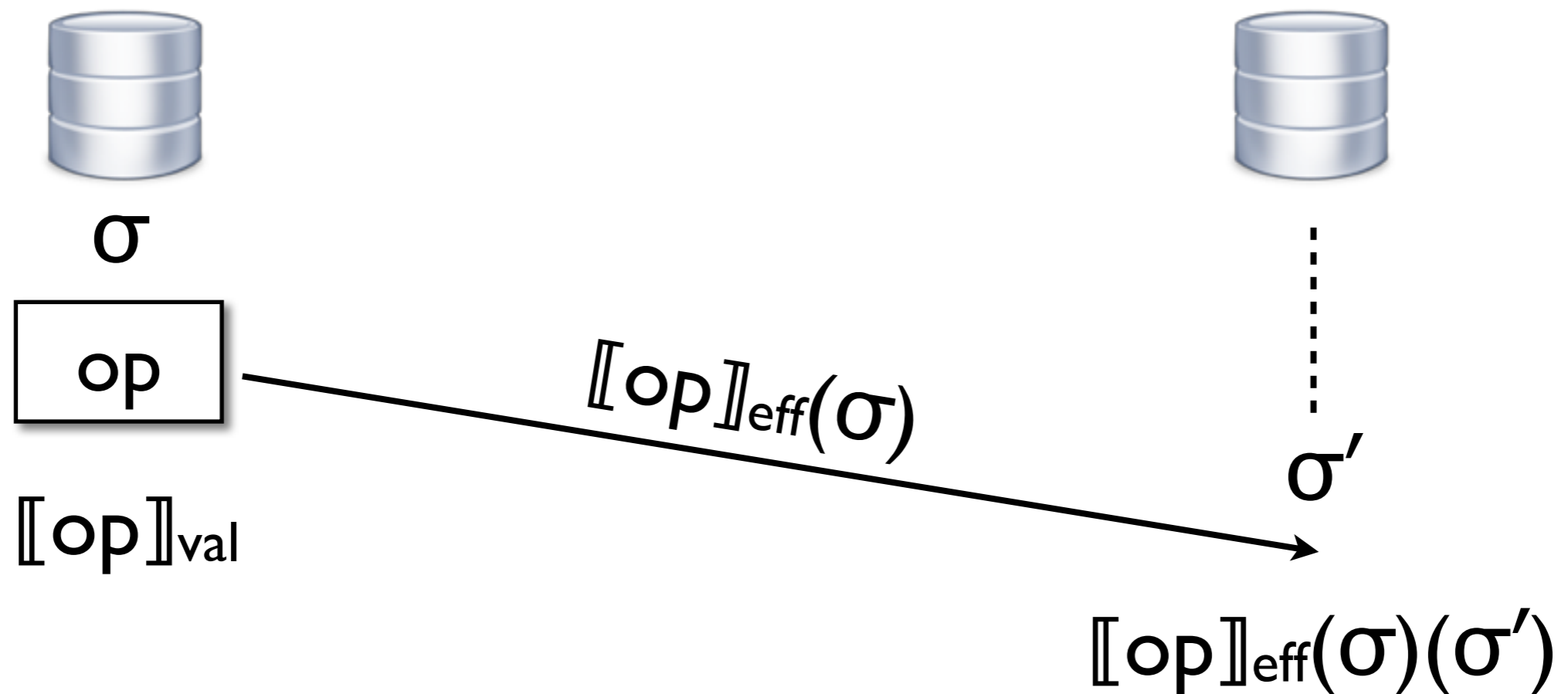
# Operation semantics



$[[withdraw(100)]]_{eff}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda\sigma'. \sigma' - 100)$  else  $(\lambda\sigma'. \sigma')$

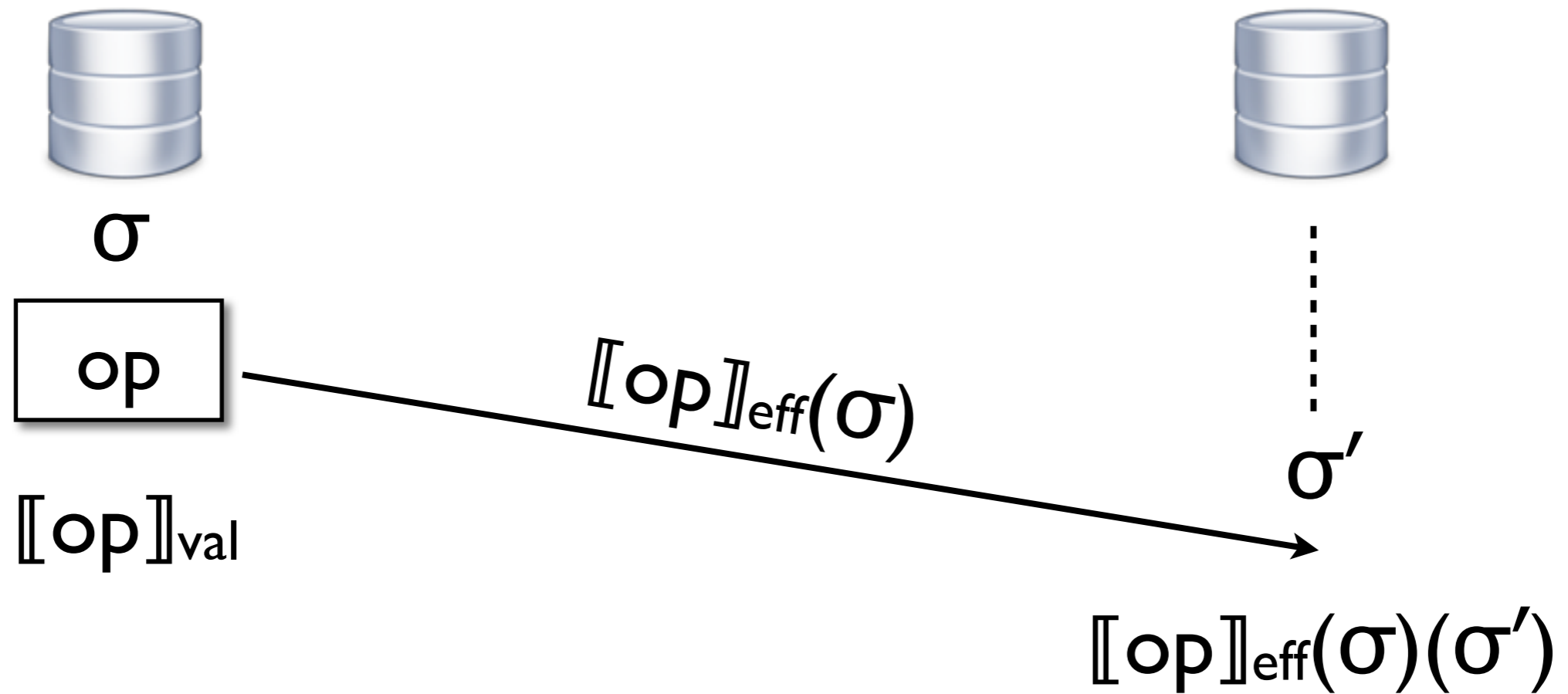
# Operation semantics



$\llbracket \text{withdraw}(100) \rrbracket_{eff}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda \sigma'. \sigma' - 100)$  else  $(\lambda \sigma'. \sigma')$

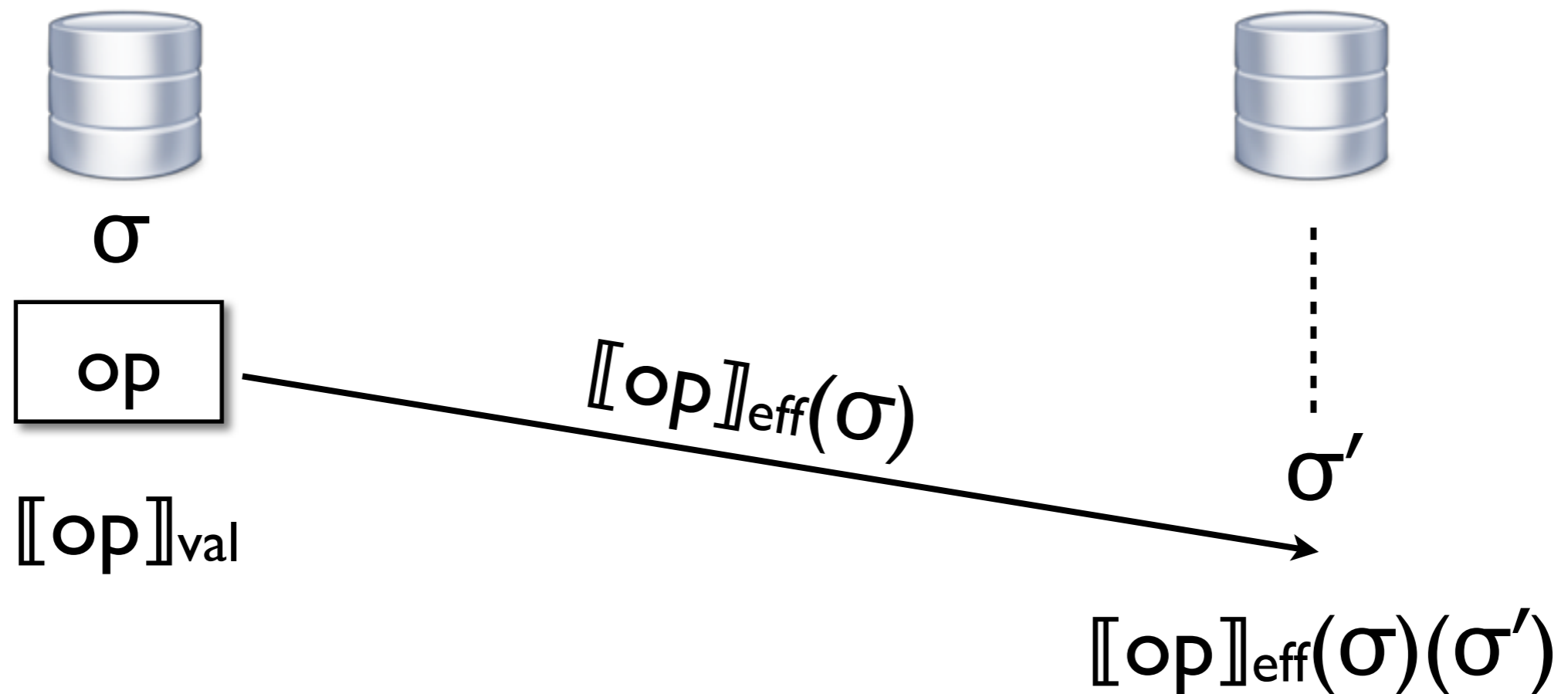
# Operation semantics



$[[withdraw(100)]]_{eff}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda\sigma'. \sigma' - 100)$  else  $(\lambda\sigma'. \sigma')$

# Operation semantics



$[[withdraw(100)]]_{eff}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda\sigma'. \sigma' - 100)$  else  $(\lambda\sigma'. \sigma')$



balance = 100

withdraw(100) : ✓

balance = 0

$\lambda\sigma'.\sigma' - 100$



balance = 100

withdraw(100) : ✓

balance = 0

$\llbracket \text{withdraw}(100) \rrbracket_{\text{eff}}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda\sigma'.\sigma' - 100)$  else  $(\lambda\sigma'.\sigma')$



balance = 100

withdraw(100) : ✓

balance = 0

$\lambda\sigma'.\sigma' - 100$



balance = 100

withdraw(100) : ✓

balance = 0

balance = -100

$\llbracket \text{withdraw}(100) \rrbracket_{\text{eff}}(\sigma) =$

if  $\sigma \geq 100$  then  $(\lambda\sigma'.\sigma' - 100)$  else  $(\lambda\sigma'.\sigma')$

# Strengthening consistency

Token system  $\approx$  locks on steroids:

- Token =  $\{\tau_1, \tau_2, \dots\}$
- Symmetric conflict relation  $\bowtie \subseteq \text{Token} \times \text{Token}$



# Strengthening consistency

Token system  $\approx$  locks on steroids:

- $\text{Token} = \{\tau_1, \tau_2, \dots\}$
- Symmetric conflict relation  $\bowtie \subseteq \text{Token} \times \text{Token}$

Example - mutual exclusion lock:

$\text{Token} = \{\tau\}; \tau \bowtie \tau$

# Strengthening consistency

Token system  $\approx$  locks on steroids:

- $\text{Token} = \{\tau_1, \tau_2, \dots\}$
- Symmetric conflict relation  $\bowtie \subseteq \text{Token} \times \text{Token}$

Example - mutual exclusion lock:

$\text{Token} = \{\tau\}; \tau \bowtie \tau$

Each operation associated with a set of tokens:

$\llbracket \text{op} \rrbracket_{\text{tok}} \in \text{State} \rightarrow \mathcal{P}(\text{Token})$

# Operations associated with conflicting tokens cannot be unaware of each other



balance = 100

$T \otimes T$



balance = 100

withdraw(100) : ✓

{T}

# Operations associated with conflicting tokens cannot be unaware of each other



balance = 100

withdraw(100) : ✓

{T}

T ✕ T



balance = 100



withdraw(100) : ?

{T}

Anything I don't know about?

# Operations associated with conflicting tokens cannot be unaware of each other



balance = 100

withdraw(100) : ✓

{T}

T ✗ T



balance = 100



balance = 0

withdraw(100) : ?

{T}

# Operations associated with conflicting tokens cannot be unaware of each other



balance = 100

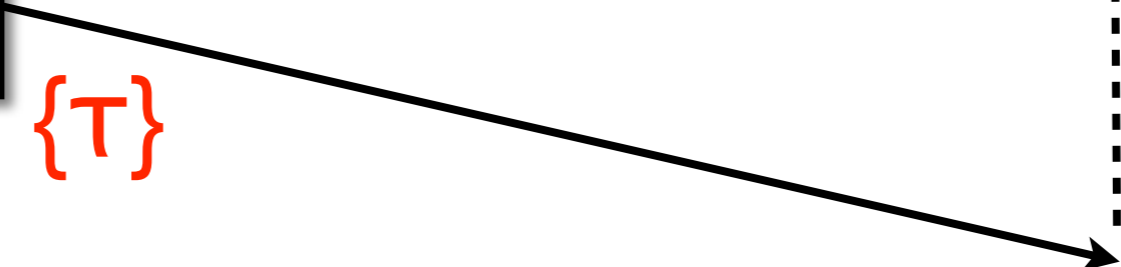
$T \times T$



balance = 100

withdraw(100) : ✓

{T}



balance = 0



withdraw(100) : ✗

{T}

# Operations associated with conflicting tokens cannot be unaware of each other



balance = 100

withdraw(100) : ✓

{T}



deposit(100)

∅

No synchronisation

T ✕ T



balance = 100



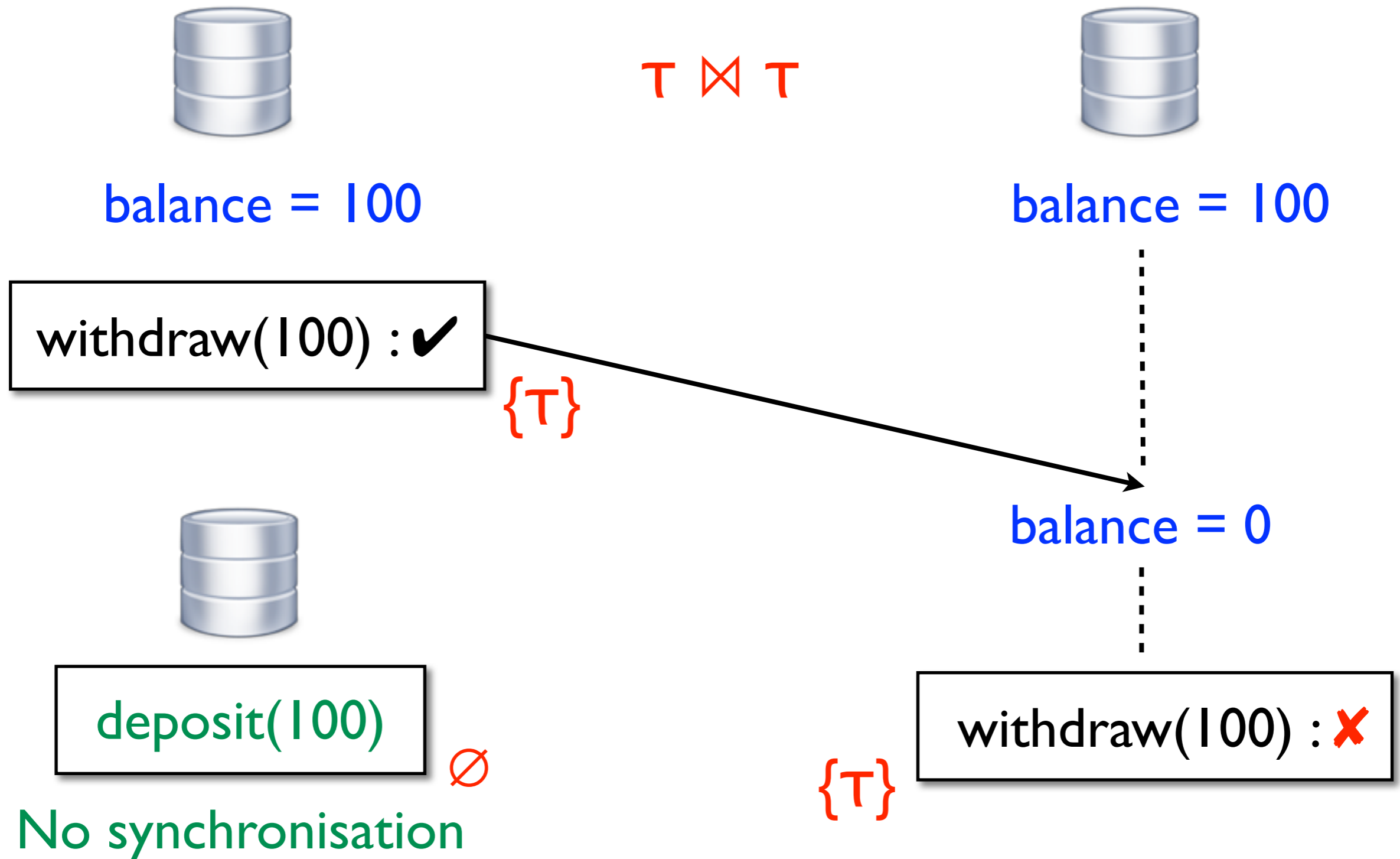
balance = 0

withdraw(100) : ✗

{T}

# Do we always have $I = (\text{balance} \geq 0)$ ?

## Rely-guarantee-based proof rule







$\sigma \in \mathcal{I}$

op



Assume invariant holds



Check it's preserved after  
executing op



$\sigma \in I$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in I?$

Effect applied in a different state!



$\sigma \in I$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in I?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{withdraw}(100) \rrbracket_{\text{eff}}(\sigma) =$

$\text{if } \sigma \geq 100 \text{ then } (\lambda\sigma'. \sigma' - 100) \text{ else } (\lambda\sigma'. \sigma')$



$\sigma \in \mathcal{I}$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in \mathcal{I}?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

- Effector safety:**  $f$  preserves  $\mathcal{I}$  when executed in any state satisfying  $P$ :  $\{\mathcal{I} \wedge P\} f \{\mathcal{I}\}$



$\sigma \in \mathcal{I}$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in \mathcal{I}?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

1. **Effector safety:**  $f$  preserves  $\mathcal{I}$  when executed in any state satisfying  $P$ :  $\{\mathcal{I} \wedge P\} f \{\mathcal{I}\}$

$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\sigma \in \mathcal{I}$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

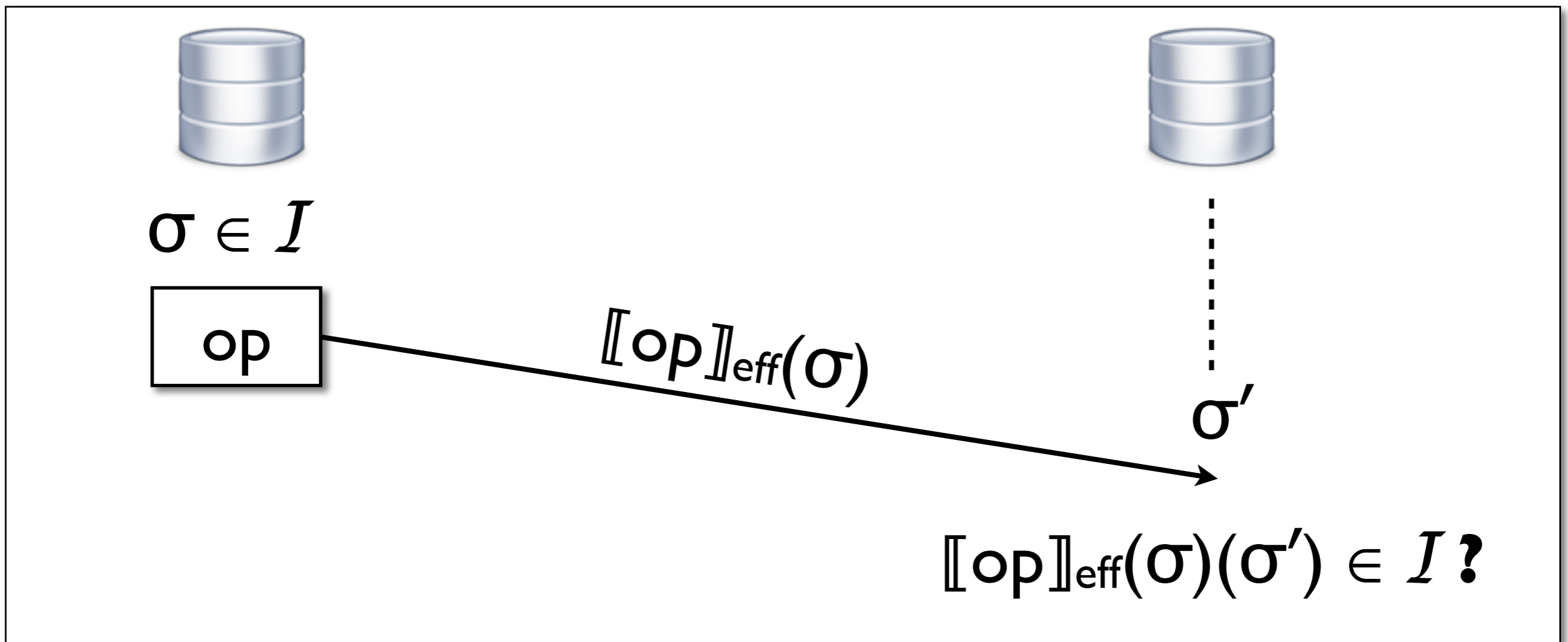
$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in \mathcal{I}?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

1. **Effector safety:**  $f$  preserves  $\mathcal{I}$  when executed in any state satisfying  $P$ :  $\{\mathcal{I} \wedge P\} f \{\mathcal{I}\}$

$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

I. **Effector safety:**  $f$  preserves  $\mathcal{J}$  when executed in any state satisfying  $P$ :  $\{\mathcal{J} \wedge P\} f \{\mathcal{J}\}$

$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\sigma \in \mathcal{I}$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in \mathcal{I}?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

1. **Effector safety:**  $f$  preserves  $\mathcal{I}$  when executed in any state satisfying  $P$ :  $\{\mathcal{I} \wedge P\} f \{\mathcal{I}\}$

$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$





$\sigma \in \mathcal{I}$

op

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)$



⋮

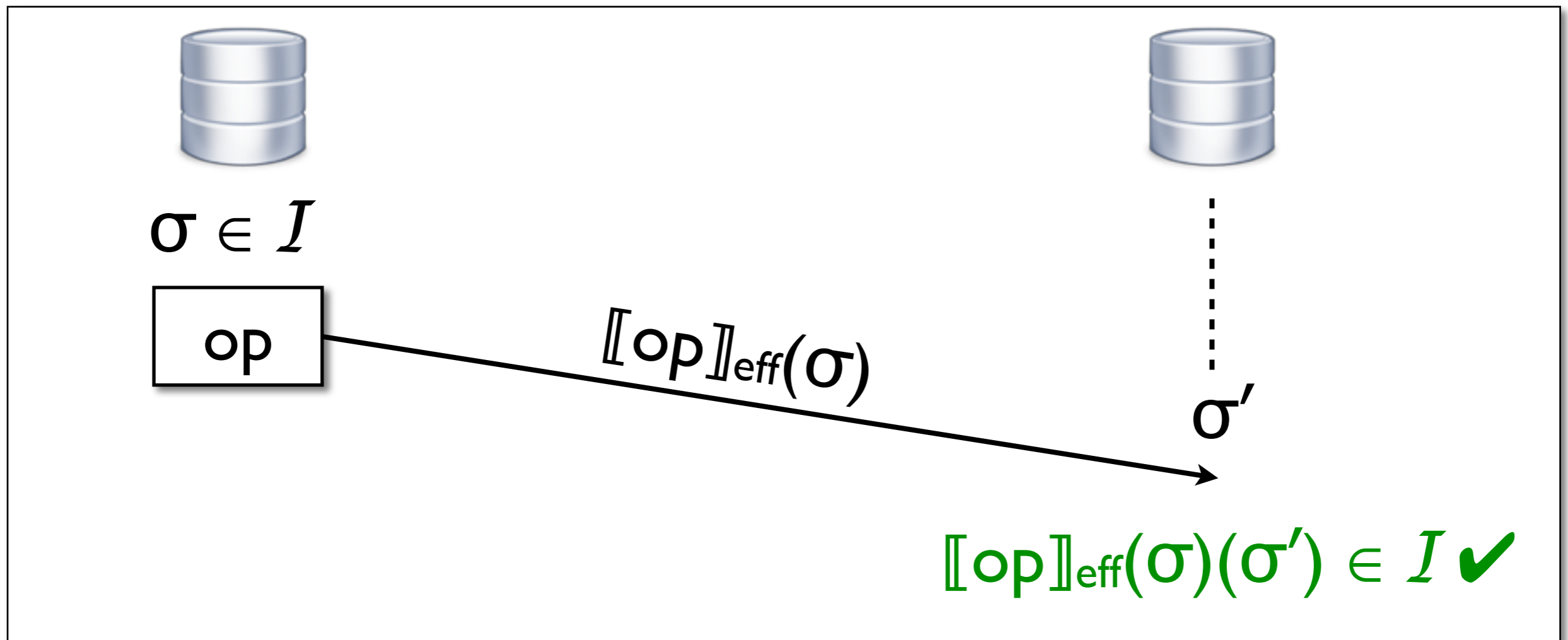
$\sigma'$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma)(\sigma') \in \mathcal{I}?$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

I. **Effector safety:**  $f$  preserves  $\mathcal{I}$  when executed in any state satisfying  $P$ :  $\{\mathcal{I} \wedge P\} f \{\mathcal{I}\}$

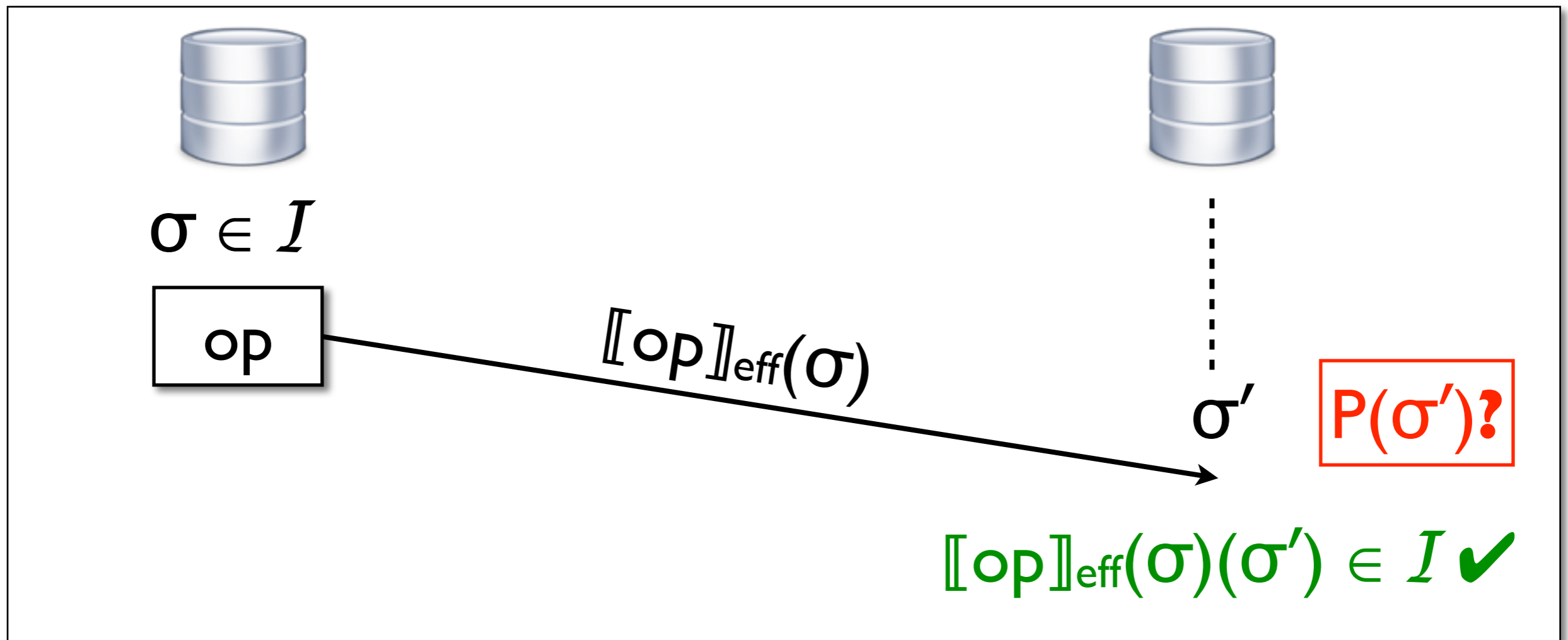
$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

I. **Effector safety:**  $f$  preserves  $\mathcal{J}$  when executed in any state satisfying  $P$ :  $\{\mathcal{J} \wedge P\} f \{\mathcal{J}\}$

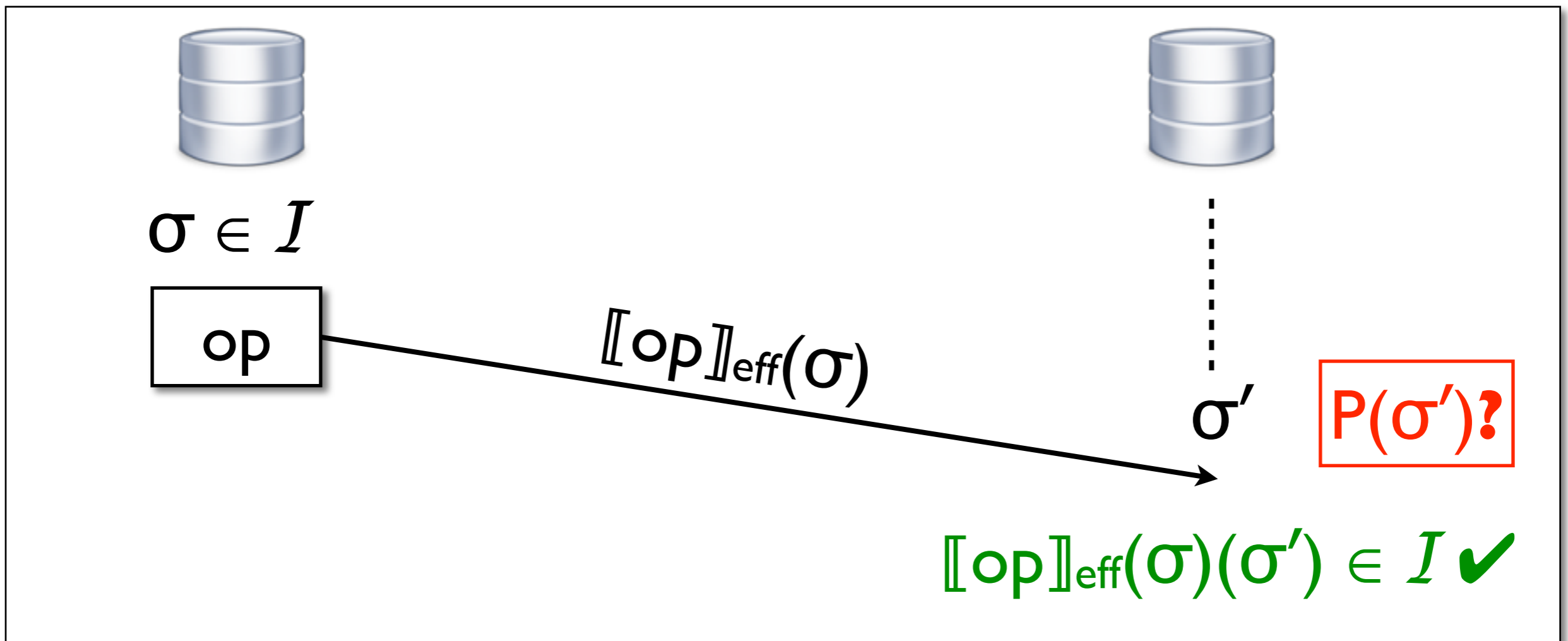
$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\llbracket op \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

- Effector safety:**  $f$  preserves  $\mathcal{J}$  when executed in any state satisfying  $P$ :  $\{\mathcal{J} \wedge P\} f \{\mathcal{J}\}$

$\{\text{bal} \geq 0 \wedge \text{bal} \geq 100\} \text{ bal} := \text{bal} - 100 \{\text{bal} \geq 0\}$



$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

1. **Effector safety:**  $f$  preserves  $\mathcal{J}$  when executed in any state satisfying  $P$ :  $\{\mathcal{J} \wedge P\} f \{\mathcal{J}\}$
2. **Precondition stability:**  $P$  will hold when  $f$  is applied at any replica

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{op} \rrbracket_{\text{tok}}(\sigma) = \text{if } P(\sigma) \text{ then } T \text{ else if...}$

$P$  is preserved by any effector  $f'$  of any operation that is not associated with a token conflicting with  $T$ :  $\{P\} f' \{P\}$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{op} \rrbracket_{\text{tok}}(\sigma) = \text{if } P(\sigma) \text{ then } T \text{ else if...}$

$P$  is preserved by any effector  $f'$  of any operation that is not associated with a token conflicting with  $T$ :  $\{P\} f' \{P\}$

$\llbracket \text{withdraw}(100) \rrbracket_{\text{tok}}(\sigma) = \{\tau\} \quad \tau \bowtie \tau$

$\llbracket \text{deposit}(100) \rrbracket_{\text{tok}}(\sigma) = \emptyset$

Check stability of withdraw's precondition against deposits:

$\{\text{bal} \geq 100\} \text{ bal} := \text{bal} + 100 \{\text{bal} \geq 100\}$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{op} \rrbracket_{\text{tok}}(\sigma) = \text{if } P(\sigma) \text{ then } T \text{ else if...}$

$P$  is preserved by any effector  $f'$  of any operation that is not associated with a token conflicting with  $T$ :  $\{P\} f' \{P\}$

$\llbracket \text{withdraw}(100) \rrbracket_{\text{tok}}(\sigma) = \{\tau\} \quad \tau \bowtie \tau$

$\llbracket \text{deposit}(100) \rrbracket_{\text{tok}}(\sigma) = \emptyset$

Check stability of withdraw's precondition against deposits:

$\{\text{bal} \geq 100\} \text{ bal} := \text{bal} + 100 \{\text{bal} \geq 100\}$

$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{op} \rrbracket_{\text{tok}}(\sigma) = \text{if } P(\sigma) \text{ then } T \text{ else if...}$

$P$  is preserved by any effector  $f'$  of any operation that is not associated with a token conflicting with  $T$ :  $\{P\} f' \{P\}$

$\llbracket \text{withdraw}(100) \rrbracket_{\text{tok}}(\sigma) = \{\tau\} \quad \tau \bowtie \tau$

$\llbracket \text{deposit}(100) \rrbracket_{\text{tok}}(\sigma) = \emptyset$

Check stability of withdraw's precondition against deposits:

$\{\text{bal} \geq 100\} \text{ bal} := \text{bal} + 100 \{\text{bal} \geq 100\}$



$\llbracket \text{op} \rrbracket_{\text{eff}}(\sigma) = \text{if } P(\sigma) \text{ then } f(\sigma) \text{ else if...}$

$\llbracket \text{op} \rrbracket_{\text{tok}}(\sigma) = \text{if } P(\sigma) \text{ then } T \text{ else if...}$

$P$  is preserved by any effector  $f'$  of any operation that is not associated with a token conflicting with  $T$ :  $\{P\} f' \{P\}$

$\llbracket \text{withdraw}(100) \rrbracket_{\text{tok}}(\sigma) = \{\tau\} \quad \tau \bowtie \tau$

$\llbracket \text{deposit}(100) \rrbracket_{\text{tok}}(\sigma) = \emptyset$

Check stability of withdraw's precondition against deposits:

$\{\text{bal} \geq 100\} \text{ bal} := \text{bal} + 100 \{\text{bal} \geq 100\}$

# Prototype tool

- Automates the proof rule
- Discharges verification conditions using SMT
- Case studies:
  - ▶ fragments of web applications
  - ▶ currently applying to a distributed file system

# Conclusion

- Weak consistency poses challenges for programmability
- But pay-off often worth it: availability, cost-effectiveness
- Verification methods enable weakening consistency without compromising correctness