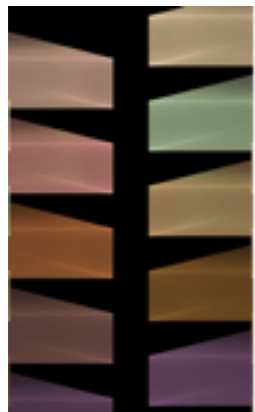


Transaction Chopping for Parallel Snapshot Isolation

Andrea Cerone,
Alexey Gotsman

Hongseok Yang



institute
imdea
software



DEPARTMENT OF
**COMPUTER
SCIENCE**

DISC - Tokyo - October 8th, 2015

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 Descubrelo

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

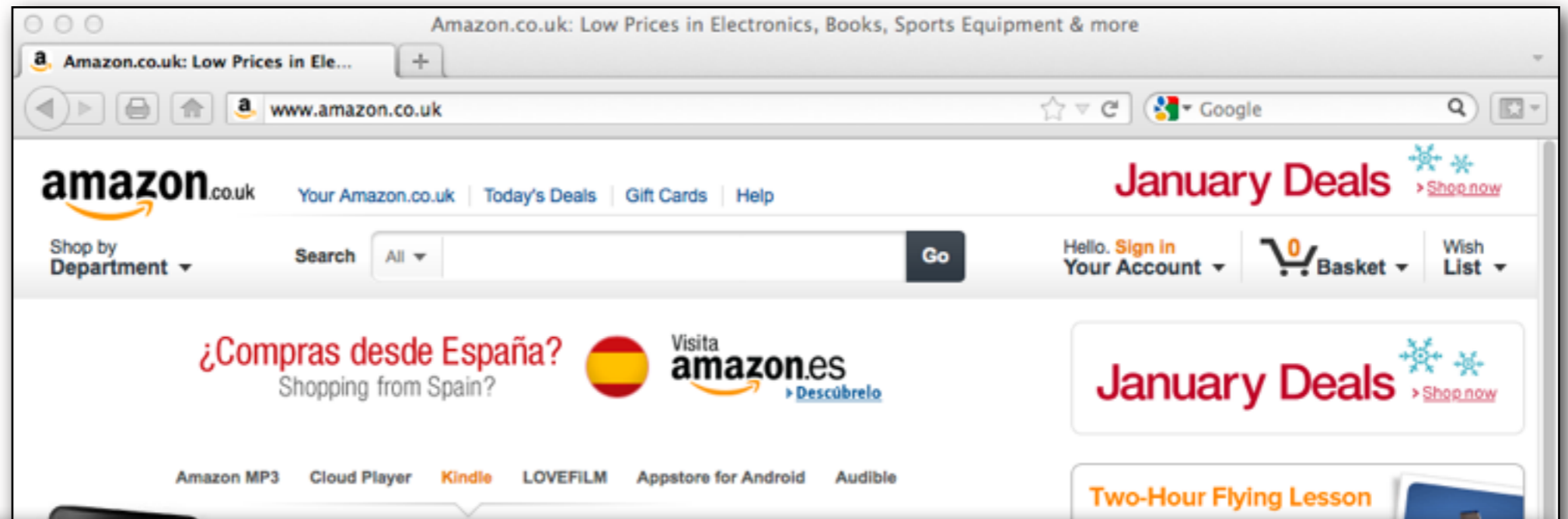
New password

Birthday

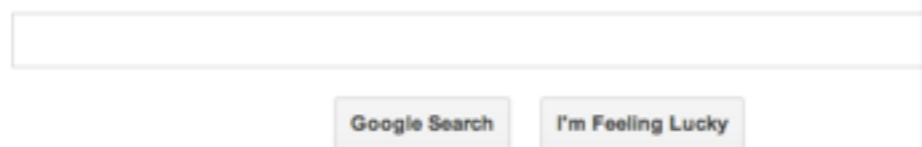
Connect with friends and the world around you on Facebook.

See photos and updates from friends in News Feed.

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



Data is replicated
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

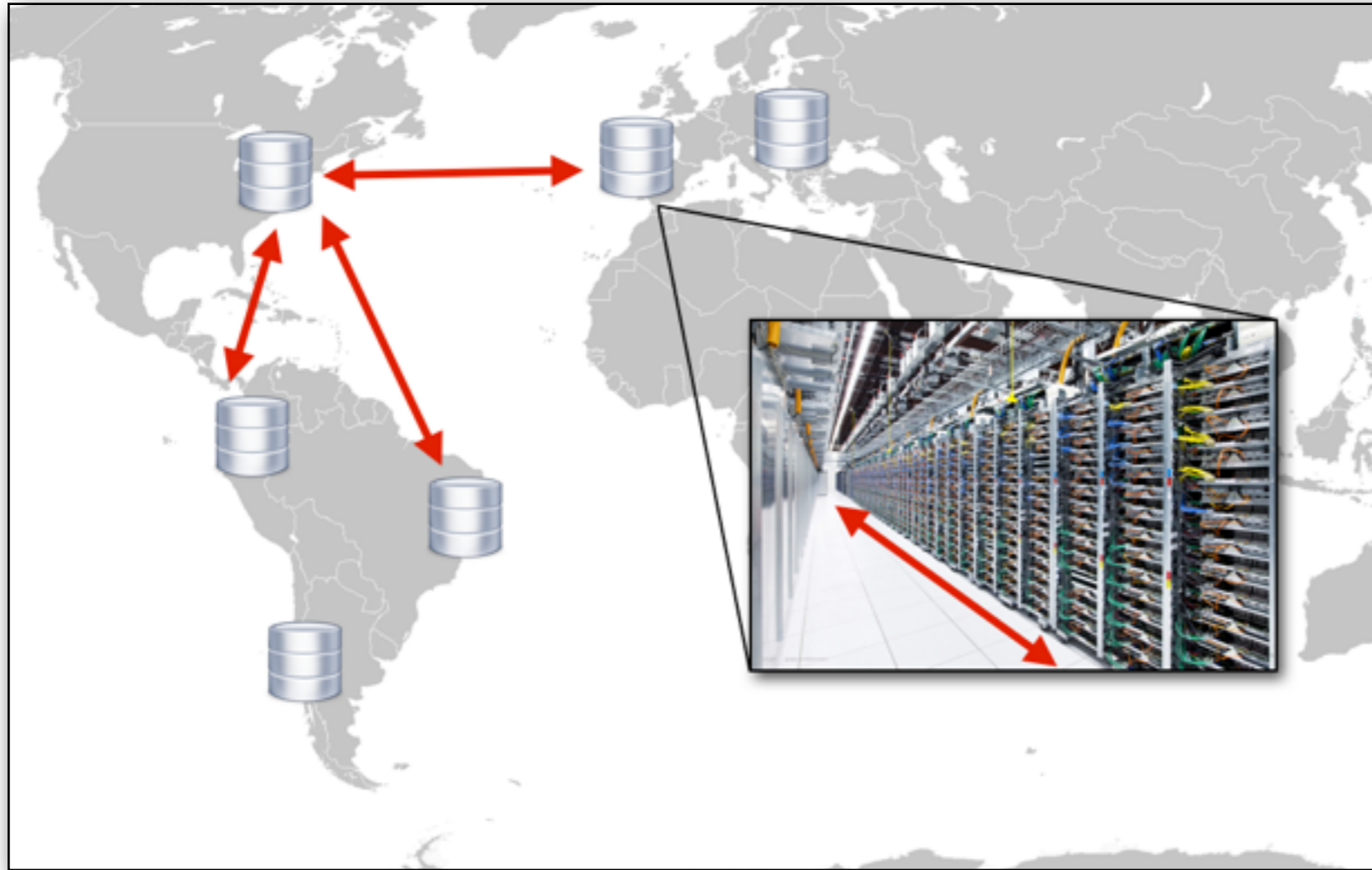


Disaster-tolerance, minimising latency

With thousands of machines inside



Fault-tolerance, load-balancing



≈



- **Serialisability:** the system behaves like a serial processor of transactions on a centralised database
- Requires **synchronisation:** expensive

Rethinking consistency in large-scale

Scalable Atomic Visibility with RAMP Transactions

Peter ... Alan Fekete[†], Ali Ghodsi, Joseph M. Hellerstein, Ion Stoica
[†] University of Sydney

Highly Available Transactions: Virtues and Limitations

Peter Bailis, Aaron Davidson, Alan Fekete[†], Ali Ghodsi, Joseph M. Hellerstein, Ion Stoica

Transactional storage for geo-replicated systems

Marcos K. Aguilera[†], Jinyang Li^{*}
[†] Silicon Valley

Yair Sovran^{*}
^{*} New

Eventually Consistent Transactions

Sebastian Burckhardt¹, Daan Leijen¹, Manuel Fähndrich¹, and Mooly Sagiv²

¹ Microsoft Research

² Tel-Aviv University

The database gives weaker guarantees to programmers

Weak Consistency Models



Performance boost

- require less synchronisation between replicas

Weak Consistency Models



Performance boost

- require less synchronisation between replicas



Anomalous behaviour

- executions which are not allowed by a serialisable database
- reasoning techniques for serialisable databases do not apply

Challenges

- Are applications OK with the proposed consistency models?

Do the non-serialisable behaviours exposed violate correctness?

- can we boost the performances of an application without violating its correctness?

Are we overpaying in performance penalties?

Challenges

- Are applications OK with the proposed consistency models?


Do the non-serialisable behaviours exposed violate correctness?

- can we boost the performances of an application without violating its correctness?

Are we overpaying in performance penalties?

**Parallel Snapshot Isolation:
Specification and Transaction Chopping**

Transaction chopping

- ⊕ **Static Analysis Technique**
Determines whether a transaction in a program can be chopped into a sequence of transactions
- ⊕ **Improves performance**
Smaller transactions lead to less conflicts
- ⊕ **Sound criterion for serialisable database**
 Shasha et al. 1995
- ⊖ **Soundness is consistency level dependant**
Soundness under PSI does not follow directly from the proof for serialisable DB

PSI

Operational Model



x, y



x, y



x, y

- Database consists of replicas storing **objects**
- Every object at every replica
- Clients issue **transactions** to be executed at replicas



```
start( $t_1$ )  
x.write(post)  
y.write(comment)  
commit( $t_1$ )
```



```
start( $t_2$ )  
x.write(other post)  
abort( $t_2$ )
```

- Write write conflict detection (concurrent transactions do not write to one same object)



```
start(t1)  
x.write(post)  
y.write(comment)  
commit(t1)
```

deliver

```
t1: x.write(post);  
y.write(comment)
```

Upon commit: send all tx updates to other replicas



```
start(t2)  
x.write(other post)  
abort(t2)
```

⋮

```
start(t3)  
y.read(empty)  
commit(t3)
```

⋮

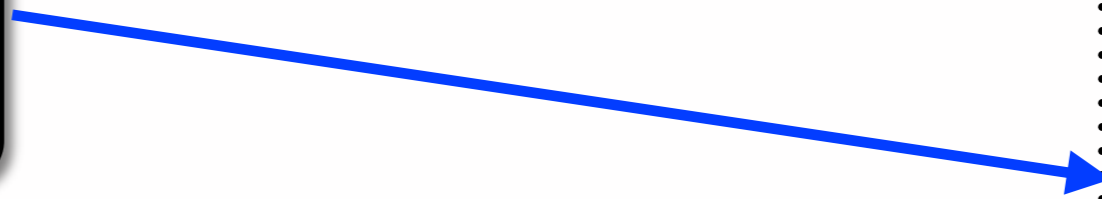
```
start(t4)  
y.read(comment)  
commit(t4)
```



start(t_1)
x.write(*post*)
commit(t_1)

⋮

start(t_2)
y.write(*comment*)
commit(t_2)



⋮

start(t_3)
y.read(*comment*)
x.read()
commit(t_3)

Message delivery:
causality is preserved



start(t_1)
x.write(*post*)
commit(t_1)

⋮

start(t_2)
y.write(*comment*)
commit(t_2)



⋮

start(t_3)
y.read(*comment*)
x.read(*post*)
commit(t_3)

Message delivery:
causality is preserved



`x.write(l)`

`y.write(l)`



Long fork

Disallowed by classical
snapshot isolation



`x.write(1)`

`y.write(1)`

`x.read(1)`

`y.read(0)`

Long fork

Disallowed by classical
snapshot isolation

*x written
before y*



`x.write(1)`

`y.write(1)`

`x.read(1)`
`y.read(0)`

`y.read(1)`
`x.read(0)`

Long fork

Disallowed by classical
snapshot isolation

*x written
before y*

*y written
before x*

Problems of the Operational Model

- Implementation Dependent
- Difficult to reason about

An alternative: Abstract Model

- Exploits the relationships between events

From operational model to abstract executions



start(t_1)
x.write(*post*)
y.write(*comment*)
commit(t_1)

deliver



start(t_2)
x.write(*other post*)
abort(t_2)
⋮

start(t_3)
y.read(*empty*)
commit(t_3)
⋮

start(t_4)
y.read(*comment*)
commit(t_4)

Abstraction from
DB events and
aborted transactions

From operational model to abstract executions

`x.write(post)`

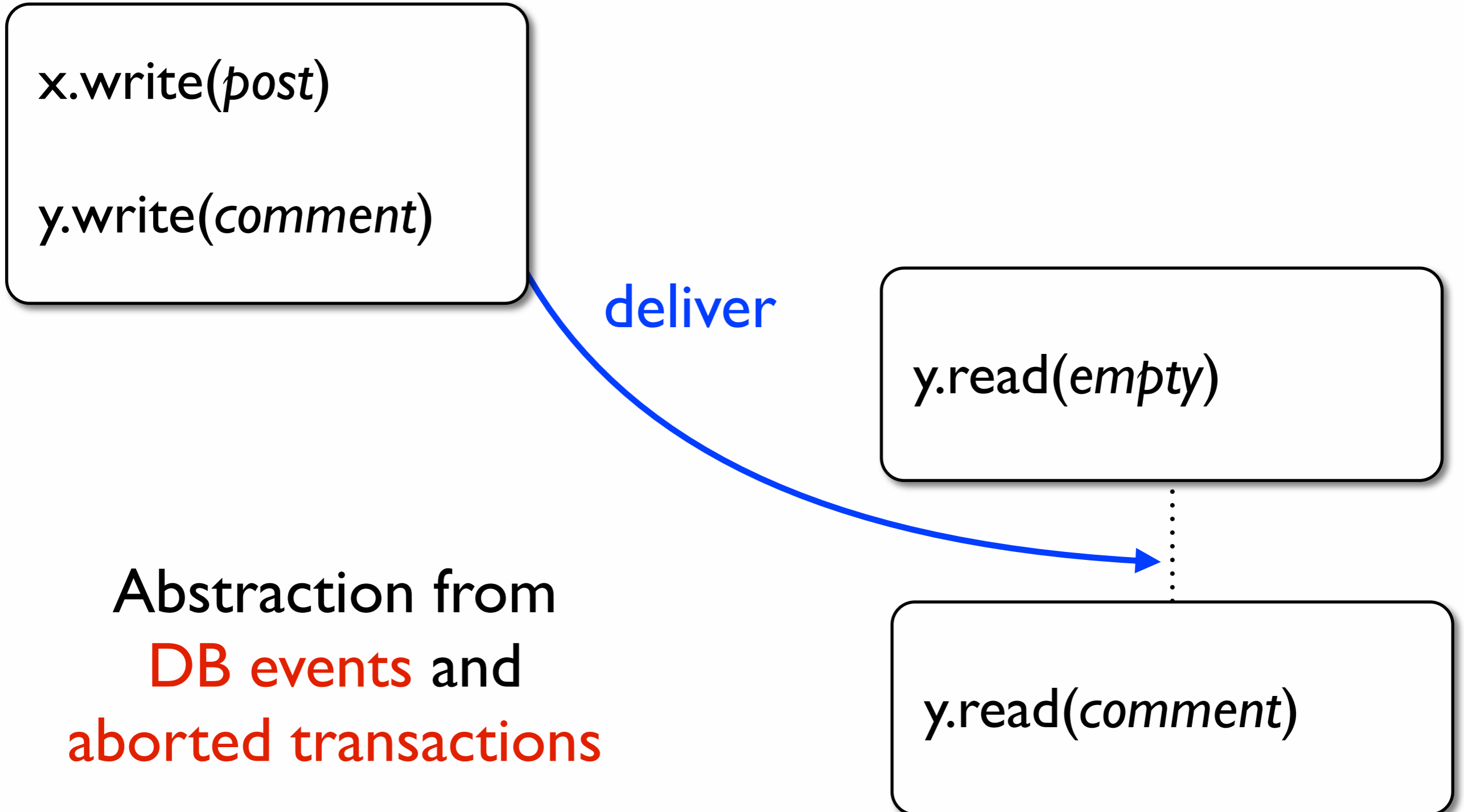
`y.write(comment)`

deliver

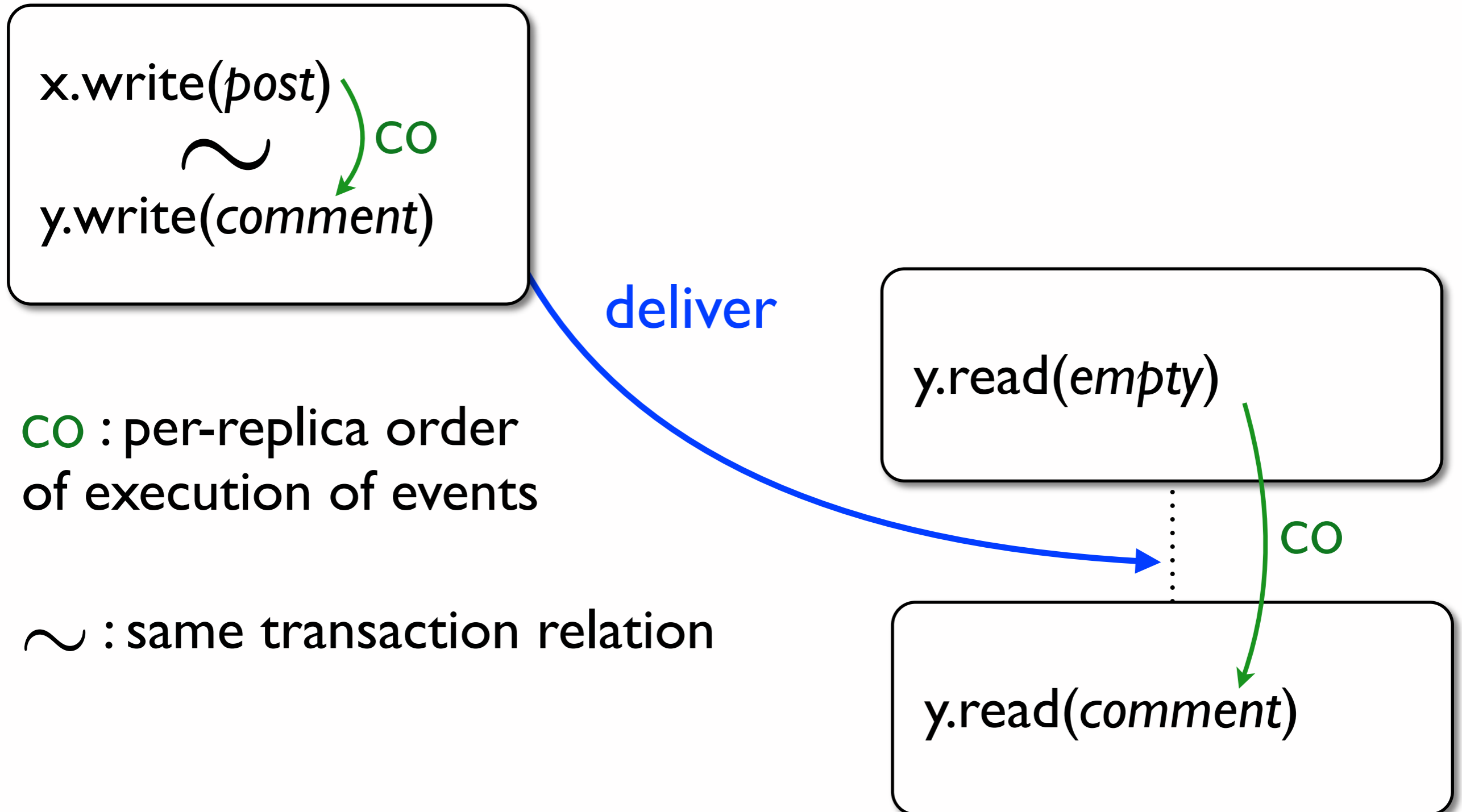
`y.read(empty)`

`y.read(comment)`

Abstraction from
DB events and
aborted transactions




From operational model to abstract executions



From operational model to abstract executions

$x.write(post)$
 \sim
 $y.write(comment)$




hb



co : per-replica order
of execution of events

$e \xrightarrow{hb} f$: either $e \xrightarrow{co} f$
or when f executes its replica
has received the effects of e

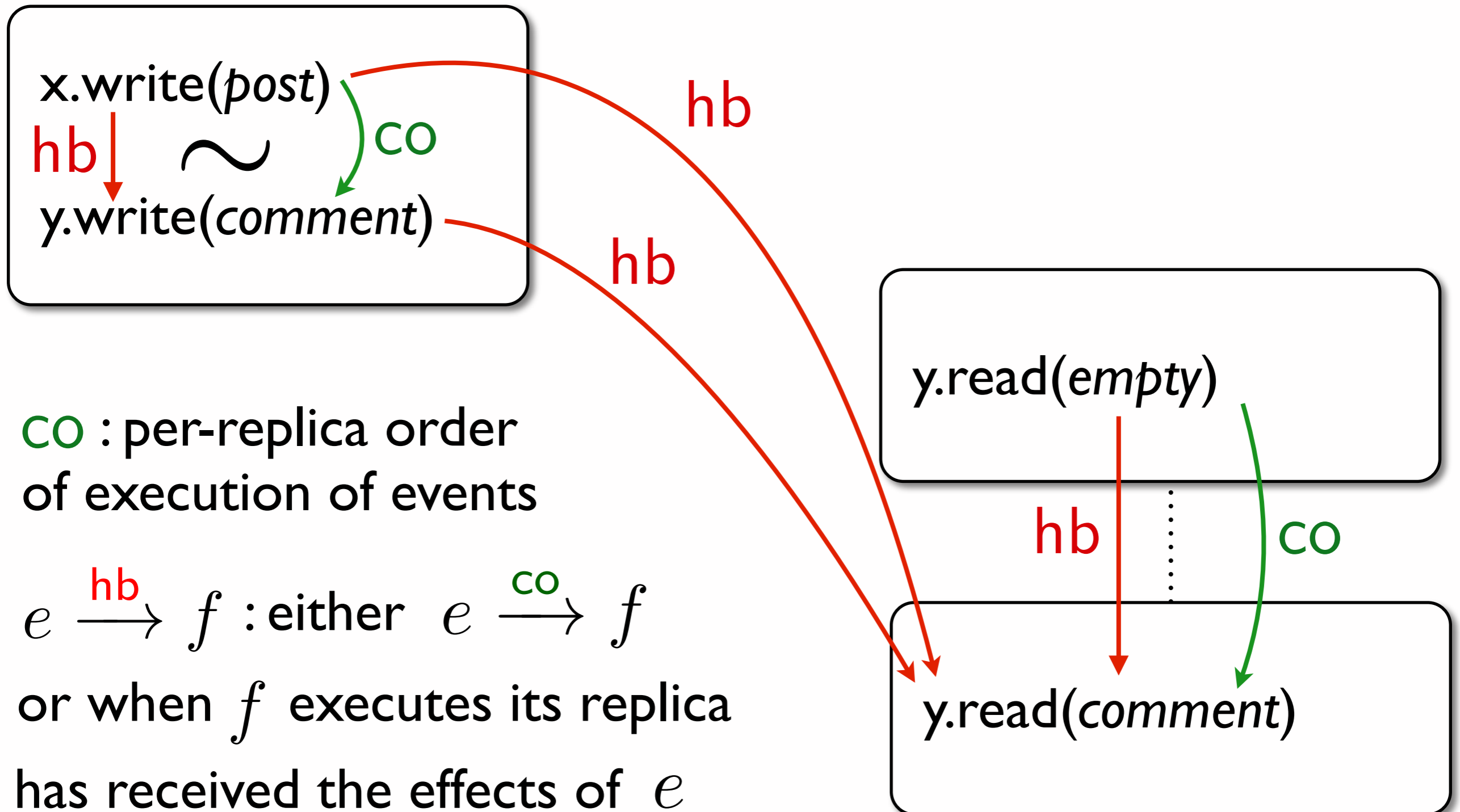
$y.read(empty)$



$y.read(comment)$

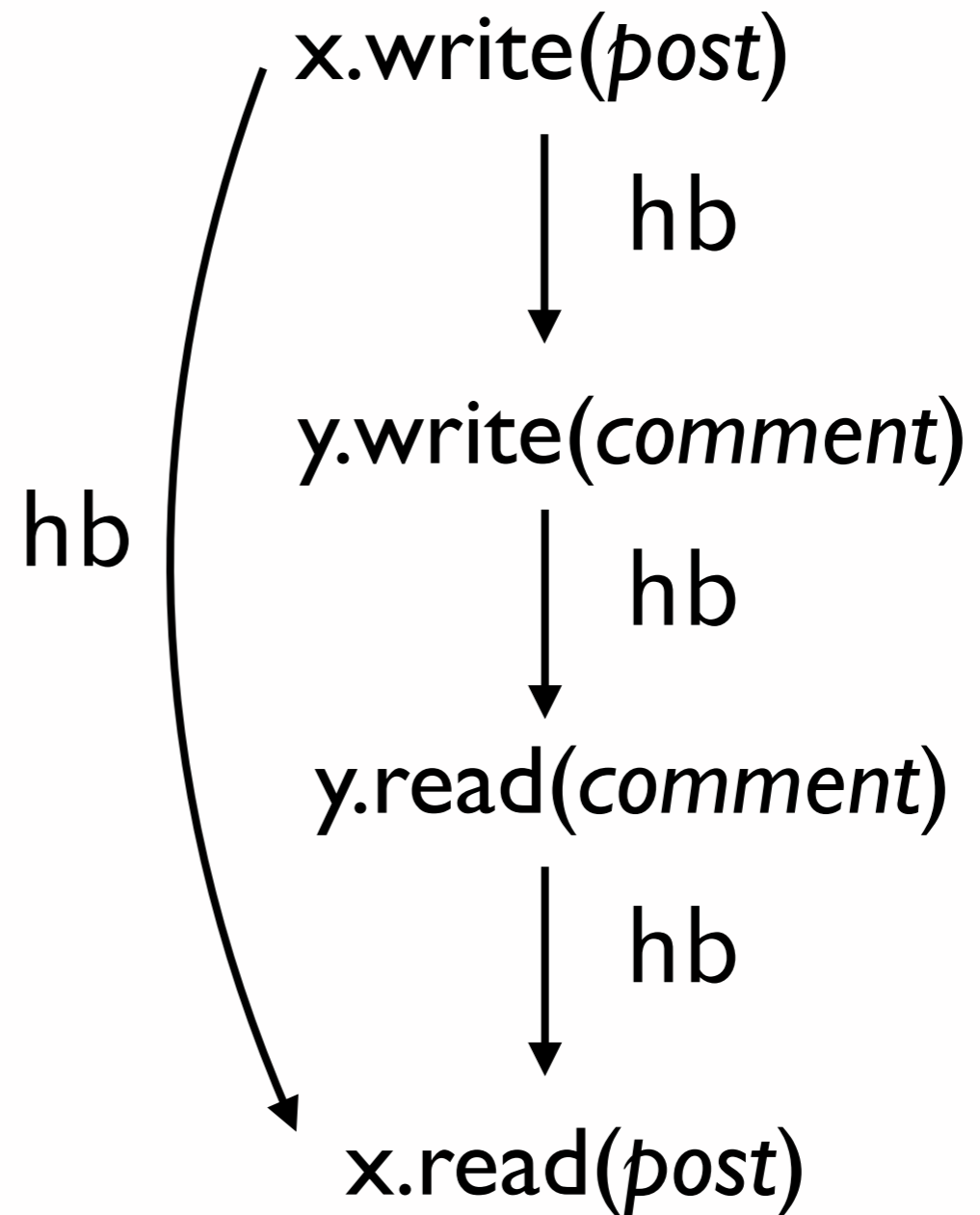


From operational model to abstract executions



Abstract model of PSI: obtained by constraining hb

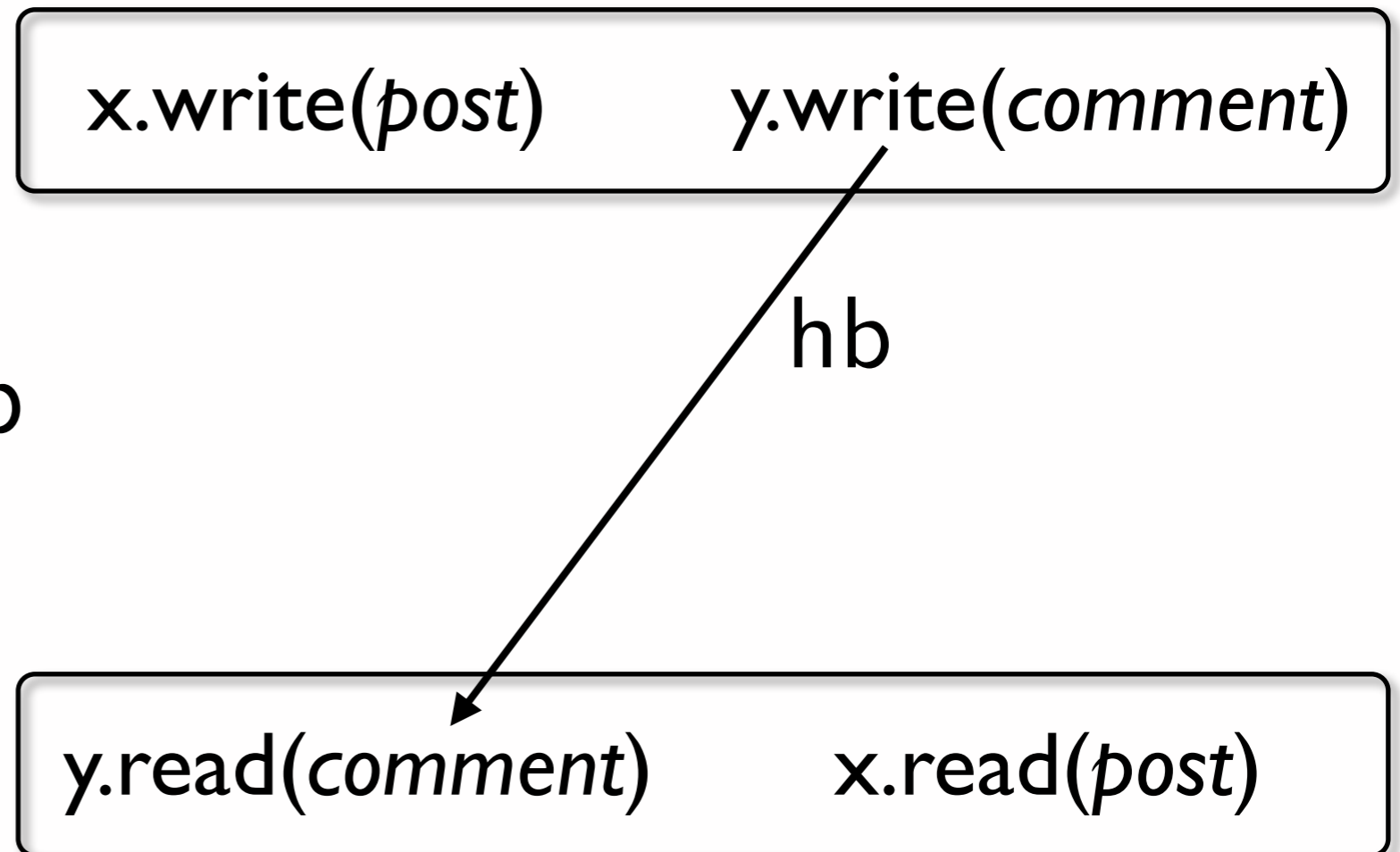
- $co \subseteq hb$
- $hb^+ \subseteq hb$



HB is transitive: causality is preserved

Abstract model of PSI: obtained by constraining hb

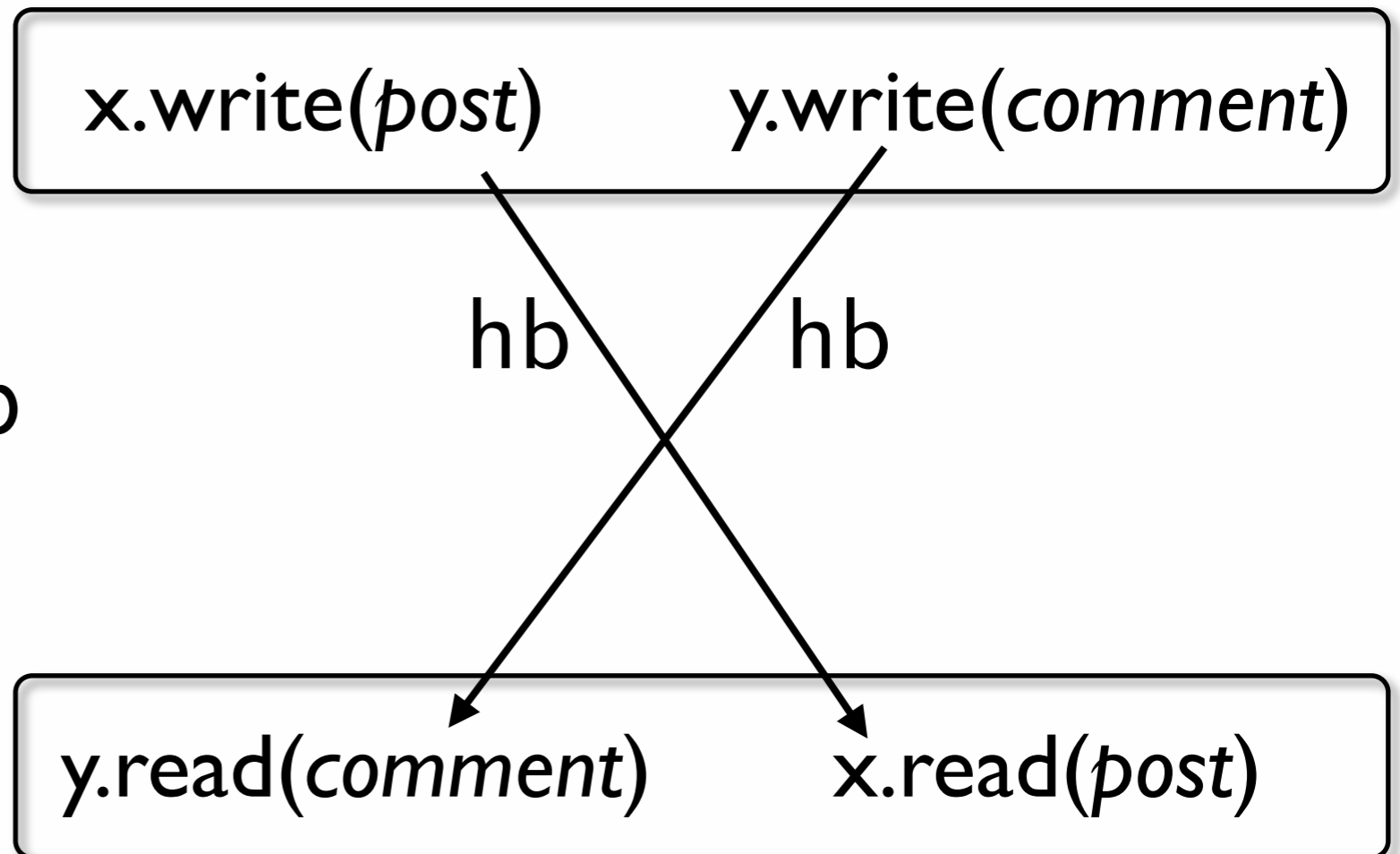
- $co \subseteq hb$
- $hb^+ \subseteq hb$
- $\sim; (hb \setminus \sim); \sim \subseteq hb$



Atomicity: either none or all the events of a transaction are observed by another one

Abstract model of PSI: obtained by constraining hb

- $co \subseteq hb$
- $hb^+ \subseteq hb$
- $\sim; (hb \setminus \sim); \sim \subseteq hb$

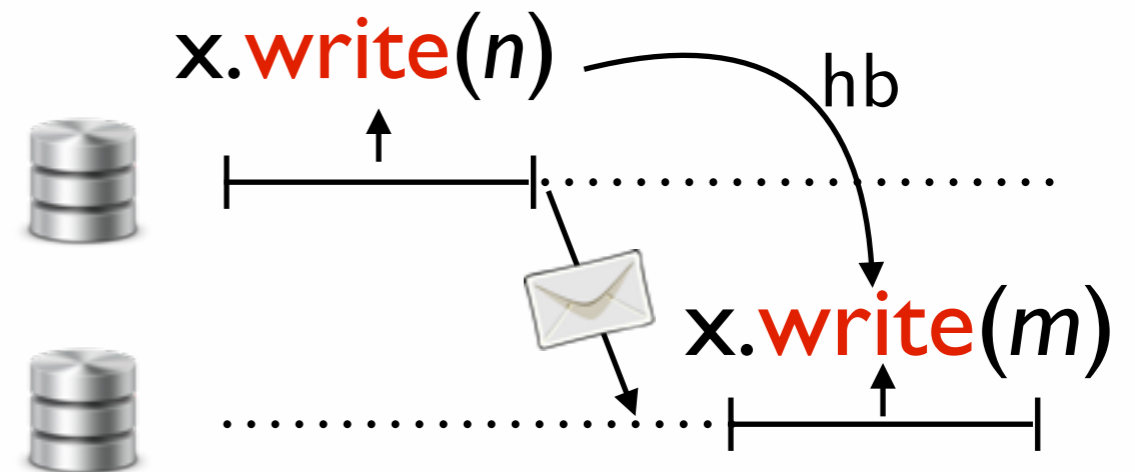
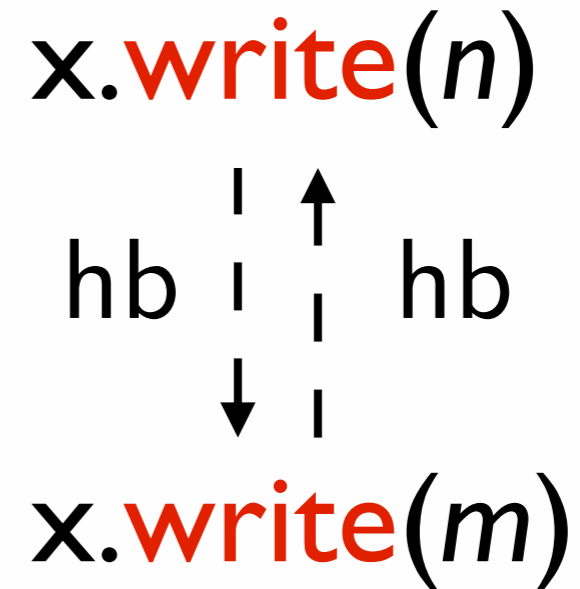


Atomicity: either none or all the events of a transaction are observed by another one

Abstract model of PSI: obtained by constraining hb

- $co \subseteq hb$
- $hb^+ \subseteq hb$
- $\sim; (hb \setminus \sim); \sim \subseteq hb$
- $\forall x. \forall e, f \in Writes_x.$
 $e = f \vee e \xrightarrow{hb} f \vee f \xrightarrow{hb} e$

**Writes on the same object
are related by hb
Write write
conflict detection:**



Abstract model of PSI: obtained by constraining hb

- $co \subseteq hb$

- $hb^+ \subseteq hb$

- $\sim; (hb \setminus \sim); \sim \subseteq hb$

- $\forall x. \forall e, f \in \text{Writes}_x.$

$$e = f \vee e \xrightarrow{hb} f \vee f \xrightarrow{hb} e$$

- $op(e) = x.read(n) \implies op(f) = x.write(n)$

where f is the last write on x happening before e

x.write(1)

↓ hb

x.write(2)

↓ hb

x.write(3)

↗ hb

x.read(3)

**read operations fetch their value from the
most recent write**

Correctness of the specification

Soundness:

every concrete execution is encoded in an abstract one that satisfies the given properties

Correctness of the specification

Soundness:

every concrete execution is encoded in an abstract one that satisfies the given properties

Completeness:

any abstract execution that satisfies the given properties can be obtained from the encoding of a concrete PSI one

Chopping: an example

```
Transaction transfer (int n) {
  TMP1 := read(acct1);
  TMP2 := read(acct2);
  write(acct2, TMP1 + TMP2);
  tryCommit;
}

Transaction lookup1() {
  TMP := read(acct1);
  tryCommit;
}
```

transfer can be **chopped** in two transactions
without introducing new behaviour

```
Transaction withdraw (int n) {
  TMP := read(acct1);
  write(acct1, TMP - n);
  tryCommit;
}

Transaction deposit (int n) {
  TMP := read(acct2);
  write(acct1, TMP + n);
  tryCommit;
}
```

```
Chain transfer' (int n) { withdraw(n); deposit(n); }
```

Chopping: an example


Chopping is not always possible:

```
Transaction transfer (int n) {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  write(acct2, TMP1 + TMP2);  
  tryCommit;  
}
```

```
Transaction Mlookup() {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  tryCommit;  
}
```

```
Chain transfer' (int n) { withdraw(n); deposit(n); }
```

Mlookup can be used to observe
an intermediate state of the database



Chopping: an example

Chopping is not always possible:

```
Transaction transfer (int n) {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  write(acct2, TMP1 + TMP2);  
  tryCommit;  
}
```

```
Transaction Mlookup() {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  tryCommit;  
}
```

```
Chain transfer' (int n) { withdraw(n); deposit(n); }
```

Mlookup can be used to observe an intermediate state of the database

acct1 = 50
acct2 = 0

transfer(20);

acct1 = 30
acct2 = 20

Chopping: an example

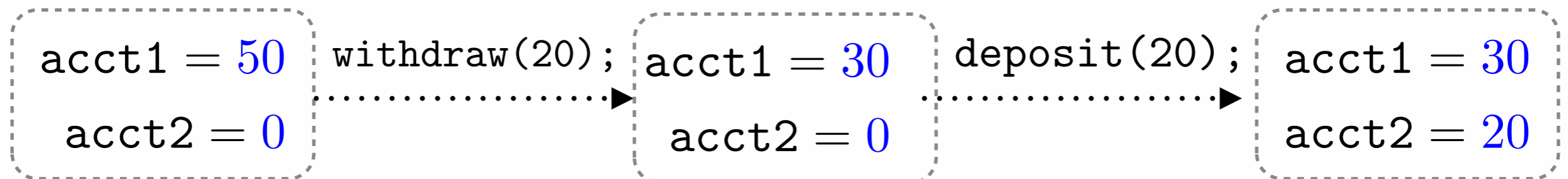
Chopping is not always possible:

```
Transaction transfer (int n) {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  write(acct2, TMP1 + TMP2);  
  tryCommit;  
}
```

```
Transaction Mlookup() {  
  TMP1 := read(acct1);  
  TMP2 := read(acct2);  
  tryCommit;  
}
```

```
Chain transfer' (int n) { withdraw(n); deposit(n); }
```

Mlookup can be used to observe an intermediate state of the database



Chopping graphs

transfer' :

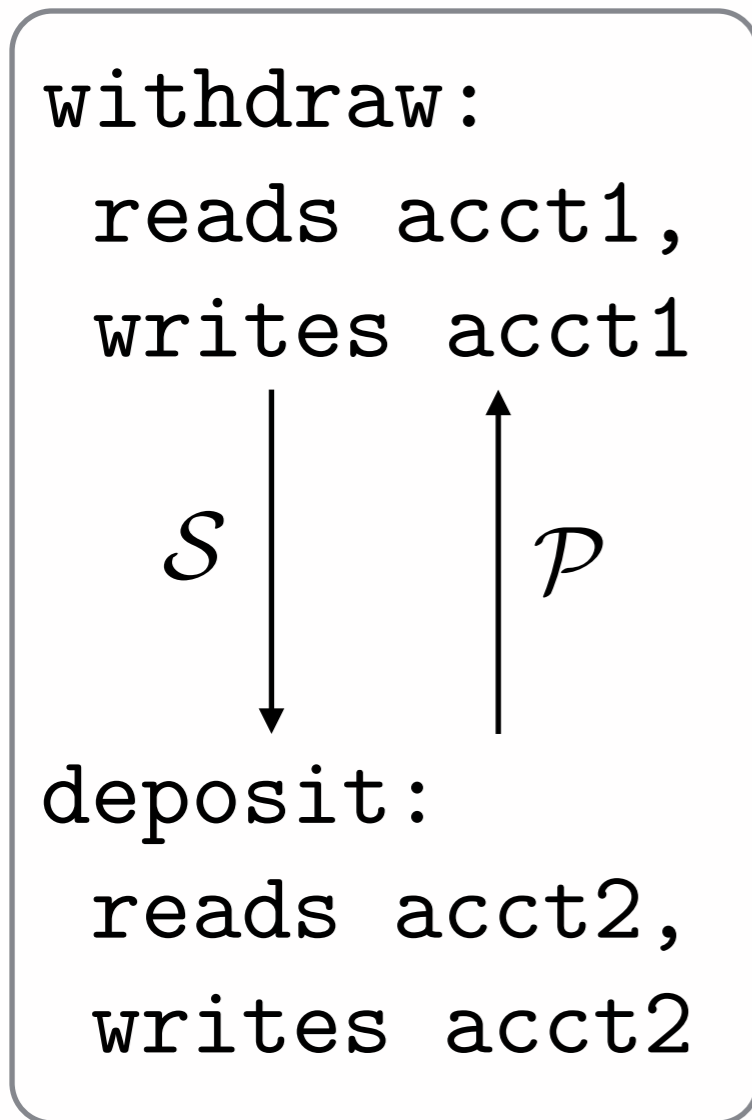
```
withdraw:  
  reads acct1,  
  writes acct1
```

```
deposit:  
  reads acct2,  
  writes acct2
```

```
Mlookup:  
  reads acct1,  
  reads acct2
```

Chopping graphs

transfer' :

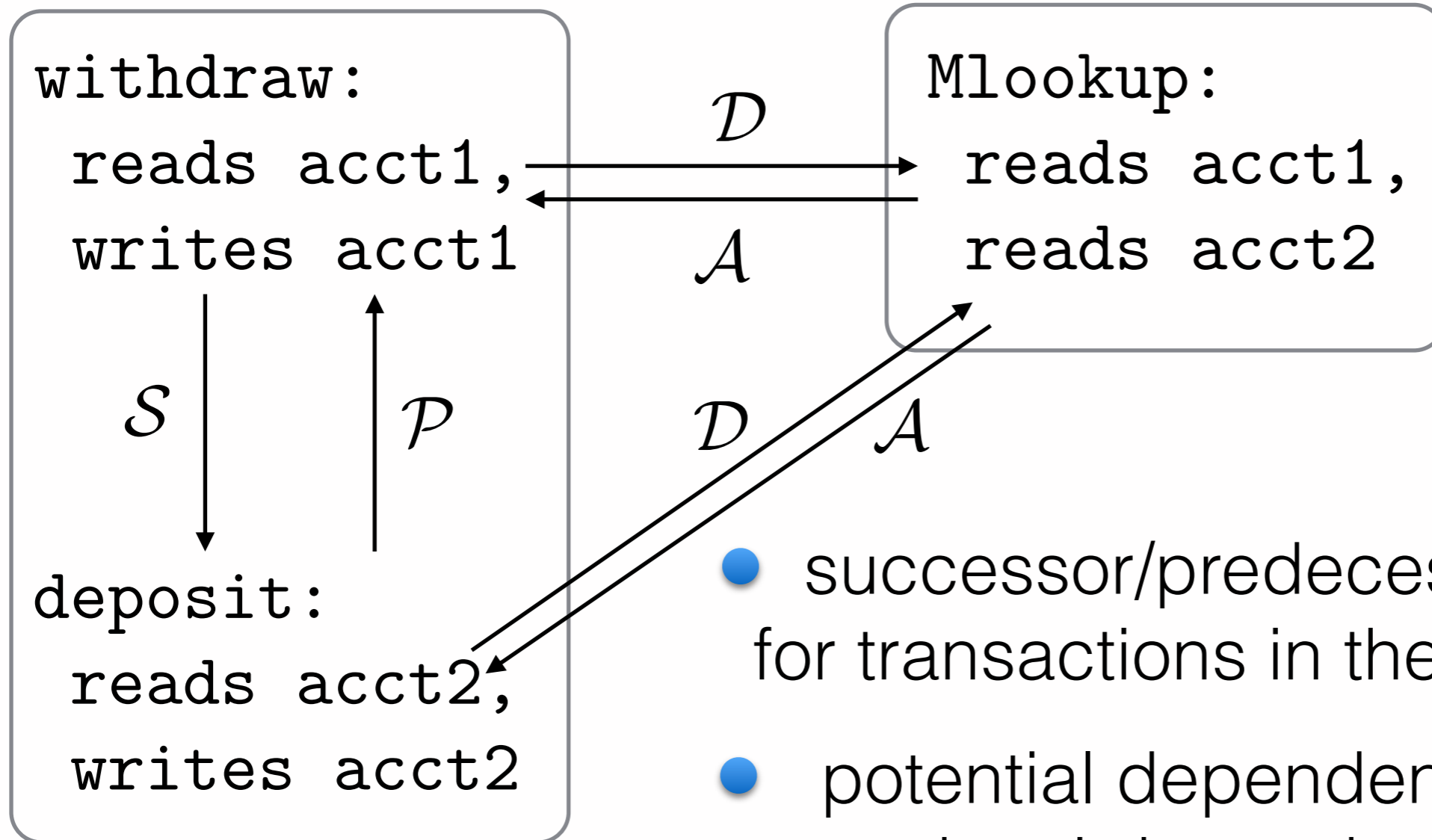


Mlookup:
reads acct1,
reads acct2

- successor/predecessor edges for transactions in the same chain

Chopping graphs

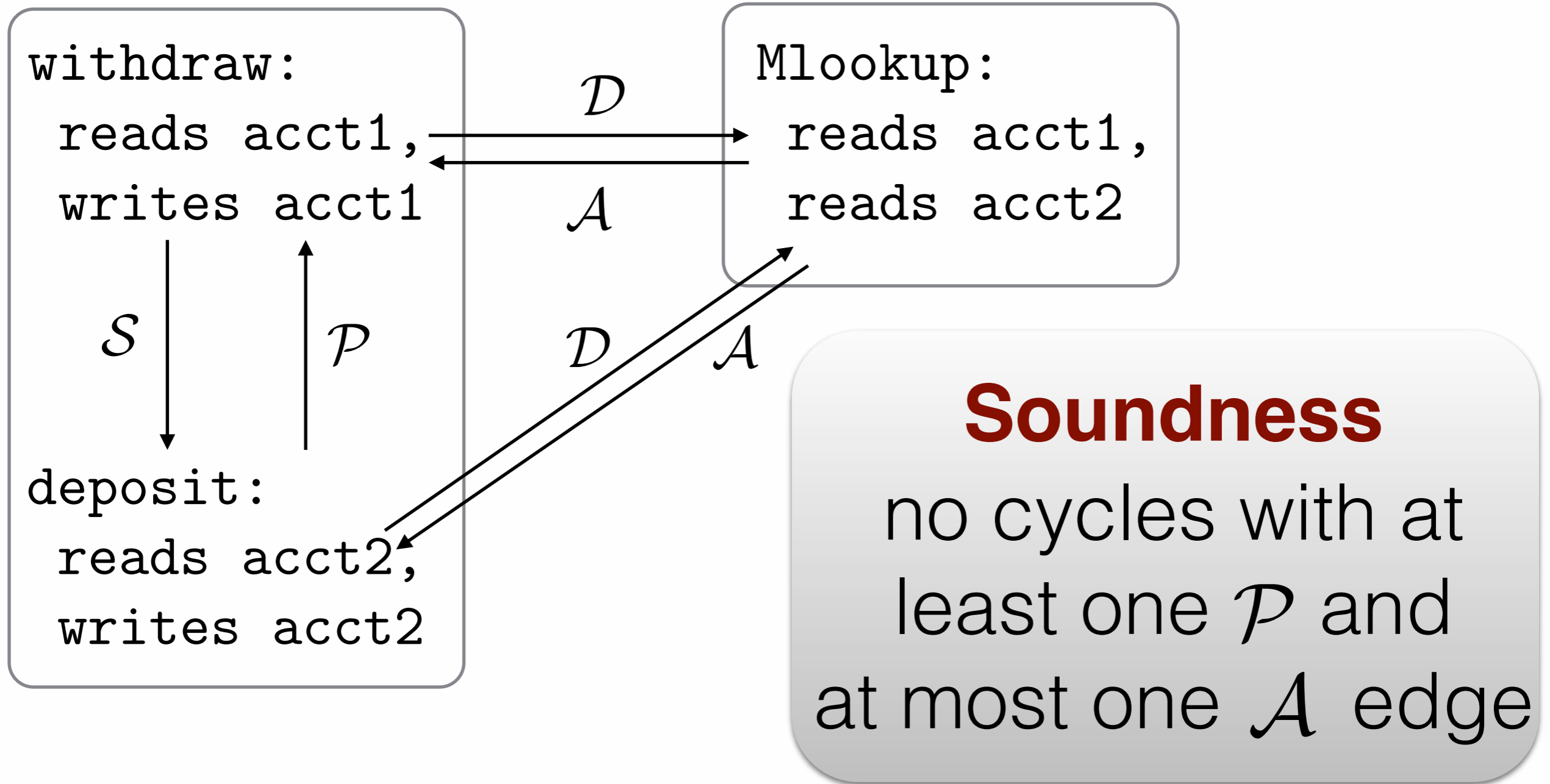
transfer' :



- successor/predecessor edges for transactions in the same chain
- potential dependencies and anti dependencies between transactions in different chains

Chopping criterion for PSI

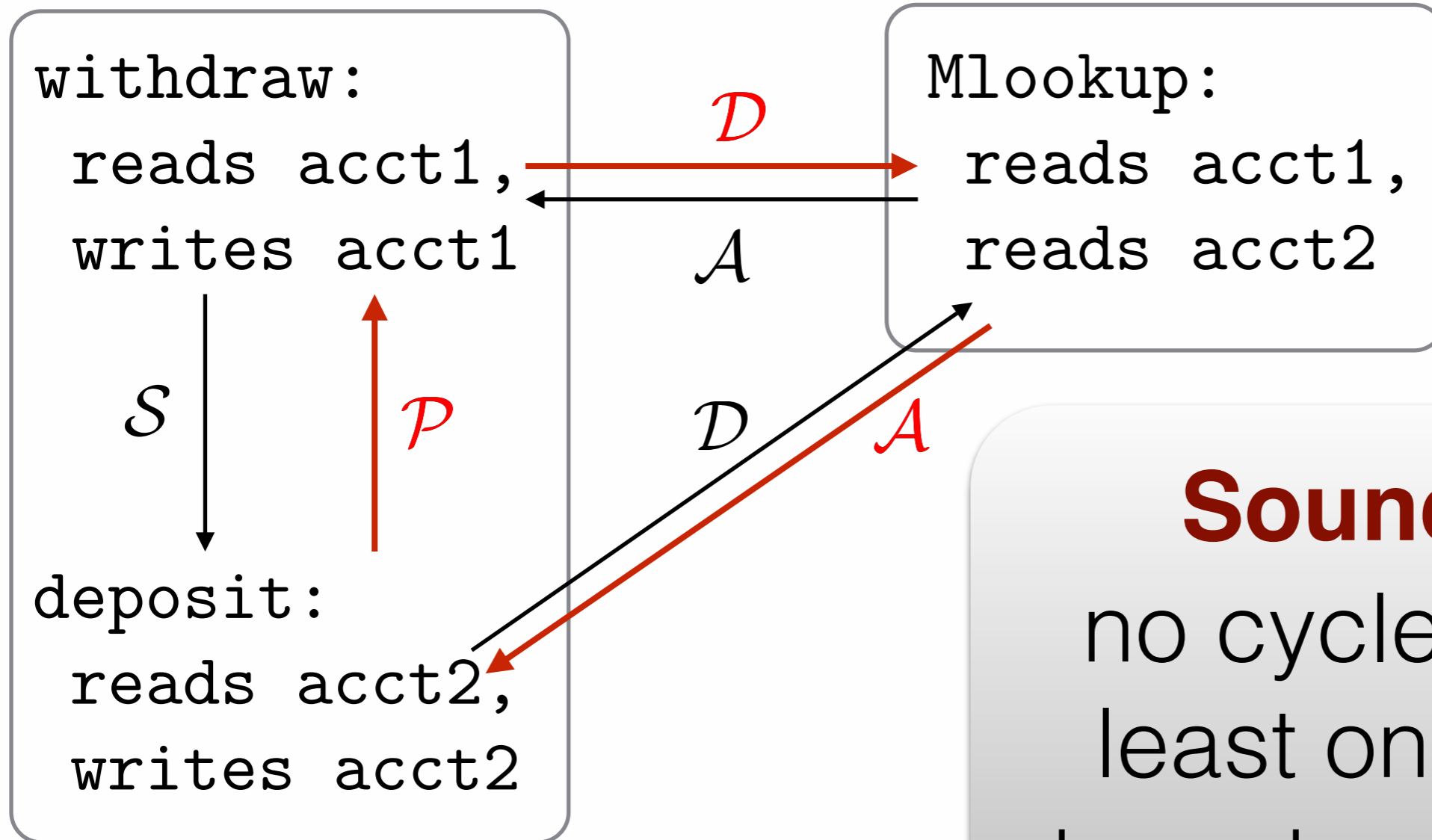
transfer':



Proof: relies heavily on the specification

Chopping criterion for PSI

transfer' :

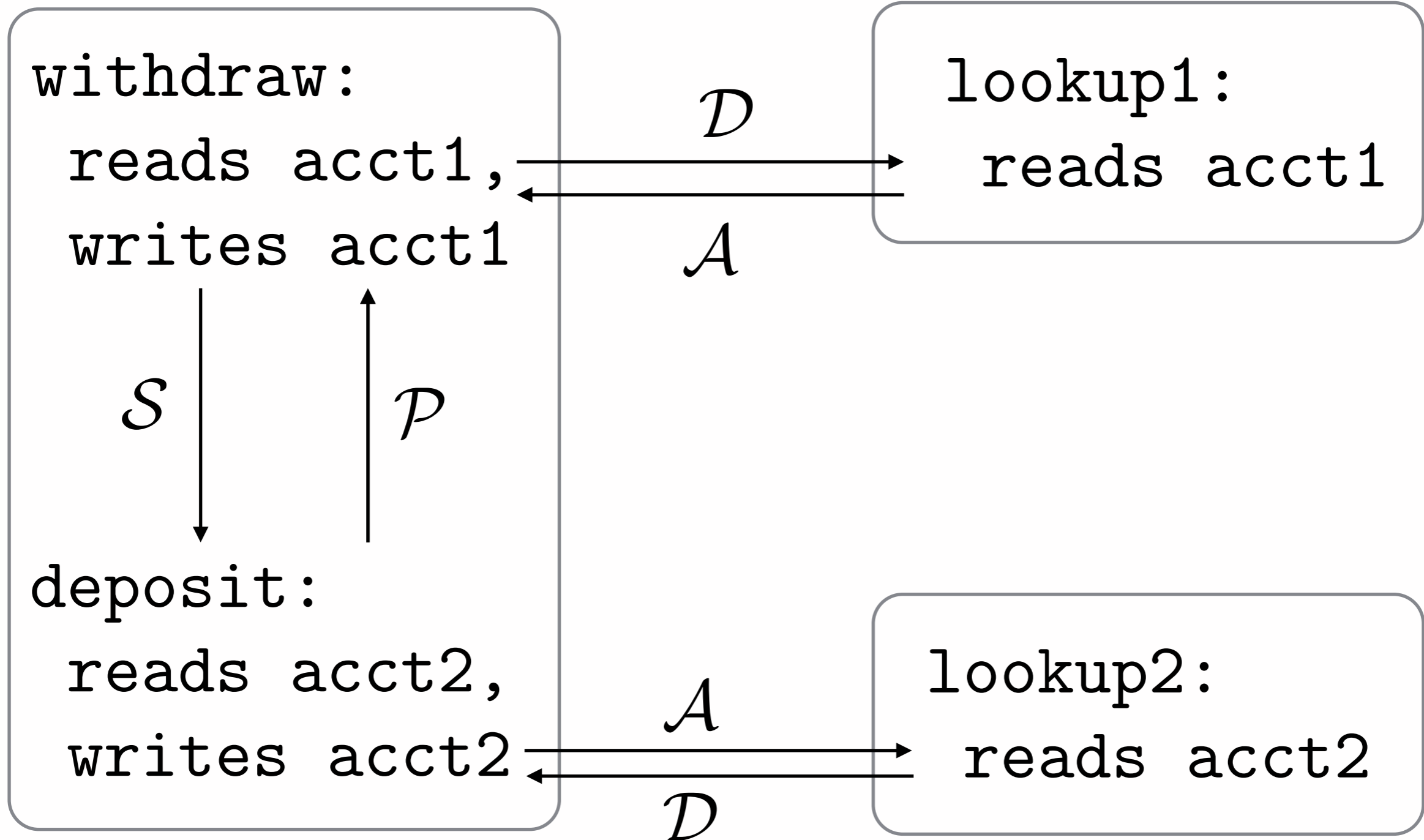


Soundness

no cycles with at least one \mathcal{P} and at most one \mathcal{A} edge

A positive example

transfer' :



Chopping: Serialisability VS. PSI

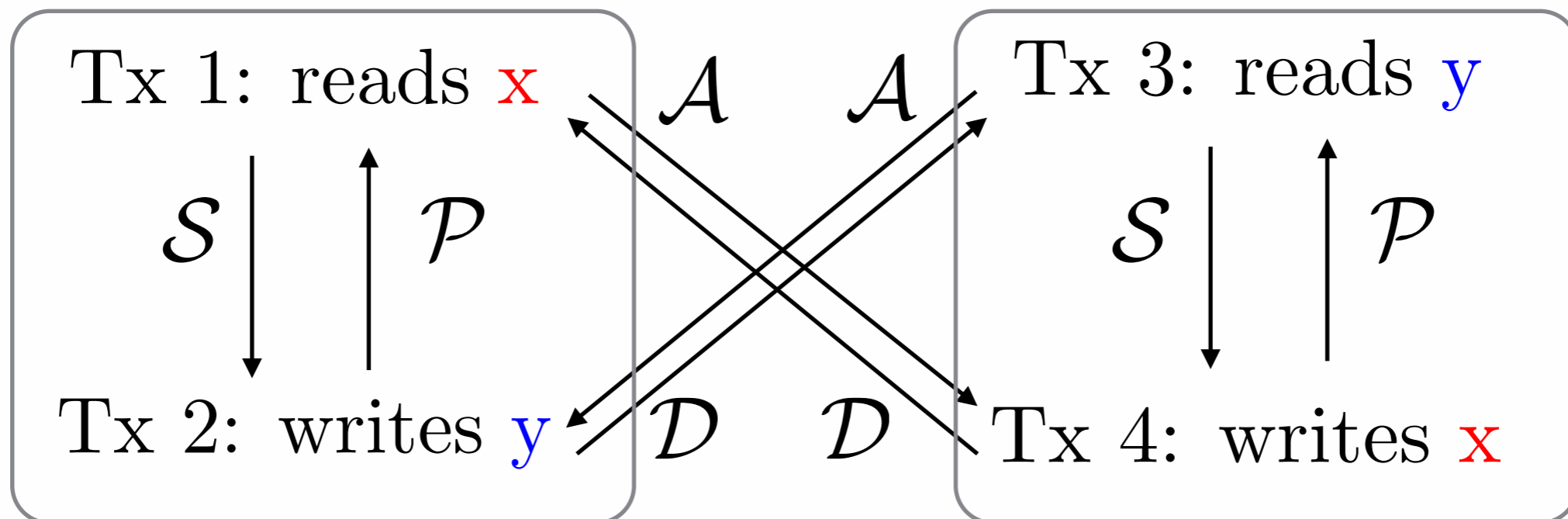
- The existing criterion for serialisability can be applied in PSI databases
proof: show that the criterion for serialisability implies the one for PSI
- But the converse is not true

Chopping: Serialisability VS. PSI

- The existing criterion for serialisability can be applied in PSI databases

proof: show that the criterion for serialisability implies the one for PSI

- But the converse is not true



What to do next (and what we have already done)

- Abstract Specification of Different Consistency Models (**CONCUR 2015**)
- Robustness (**Giovanni Bernardi, work in progress**)
ensure that the behaviour of a program is preserved when the consistency model is weakened
- Chopping for other consistency models
We already have a proposal for SI

Thank you!