

Replicated Data Types: Specification, Verification, Optimality

Hongseok Yang
University of Oxford

Joint work with Sebastian Burckhardt (Microsoft Research),
Alexey Gotsman (IMDEA) and Marek Zawirski (UPMC & INRIA)

What is it like writing programs running on top
of weakly consistent distributed stores?

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

www.amazon.co.uk

Reader

Apple iCloud Facebook Twitter Wikipedia Yahoo! News Popular

amazon.co.uk HONGSEOK's Amazon Today's Deals Gift Cards Sell Help

BOXING DAY DEALS WEEK Great Offers Until December 31 sponsored by NINTENDO 3DS [Shop now](#)

Shop by Department Search All Go Hello, HONGSEOK Your Account Try Prime Basket Wish List

Amazon Mobile Apps Cloud Player for PC and Mac LOVEFiLM Kindle Cloud Drive Appstore for Android Audible Audiobooks

Amazon uses cookies. [What are cookies?](#)

Hundreds of titles from £0.99 each in our biggest ever sale [Shop now](#)

Buy 2 or more, save 20% on Kindle Accessories [Shop now](#)

Digital Deals [Learn more](#)

Amazon Family Prize Draw Up to 60% off Fashion Boxing Day Deals Amazon Prime Subscribe & Save A Space Adventure

BOXING DAY DEALS WEEK Great Offers Until December 31 [Shop now](#) sponsored by NINTENDO 3DS

DISCOVER The Amazon Clothing Store

Shop Levi's, Ben Sherman, French Connection, 7 For All Mankind and more. [See more](#)

Advertisement

Information about each customer

The screenshot shows the Amazon.co.uk website interface. At the top, there's a browser header with tabs for 'www.amazon.co.uk' and links to 'Apple', 'iCloud', 'Facebook', 'Twitter', 'Wikipedia', 'Yahoo', 'News', and 'Popular'. A blue speech bubble highlights the 'Hello, HONGSEOK Your Account' text in the top right corner, which is enclosed in a red rounded rectangle. Below the header, the Amazon logo is on the left, followed by 'HONGSEOK's Amazon', 'Today's Deals', 'Gift Cards', 'Sell', and 'Help'. To the right, there's a 'BOXING DAY Great Offers Until December 31' banner, a search bar with 'Search All' and 'Go' button, and a 'Basket' icon with '4' items. The main content area features sections for Kindle, Cloud Drive, Appstore for Android, and Audible Audiobooks. Promotional banners include '12 Days of KINDLE' (Hundreds of titles from £0.99 each) and 'Buy 2 or more, save 20% on Kindle Accessories'. At the bottom, there are links for 'Amazon Family Prize Draw', 'Up to 60% off Fashion', 'Boxing Day Deals', 'Amazon Prime', 'Subscribe & Save', and 'A Space Adventure'. A large 'BOXING DAY DEALS WEEK' banner with 'Great Offers Until December 31' and a 'Shop now' link is prominently displayed at the bottom center.

Shopping cart (Shapiro et al.)

Bob
in UK

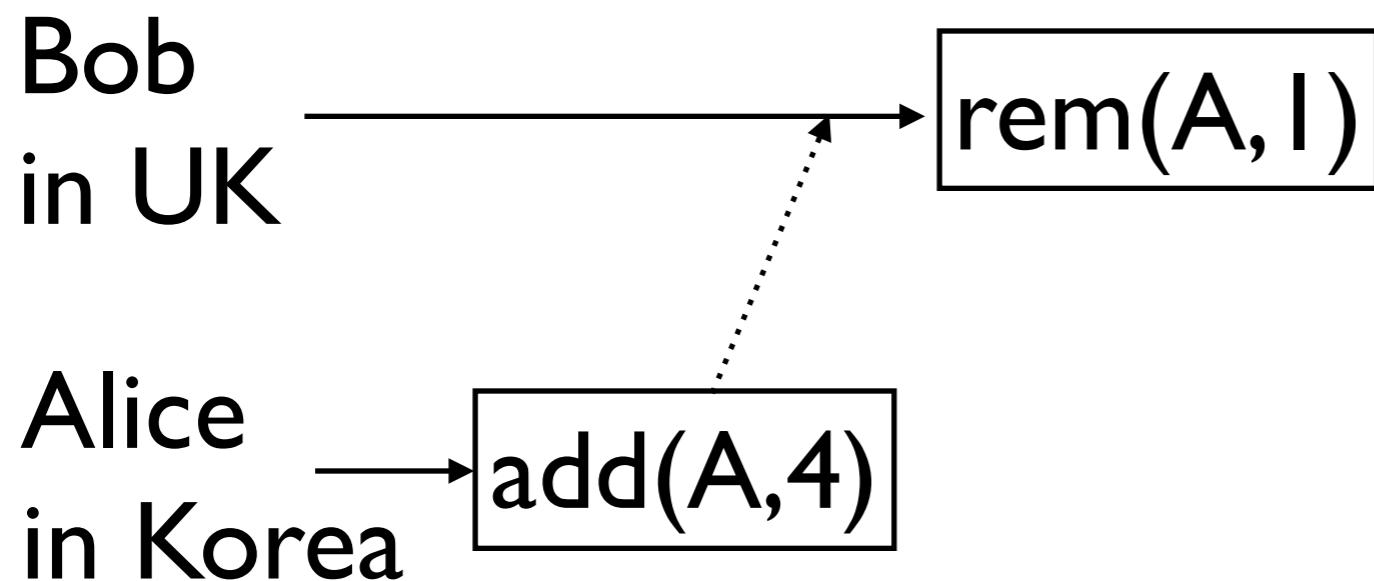
Alice
in Korea

Shopping cart (Shapiro et al.)

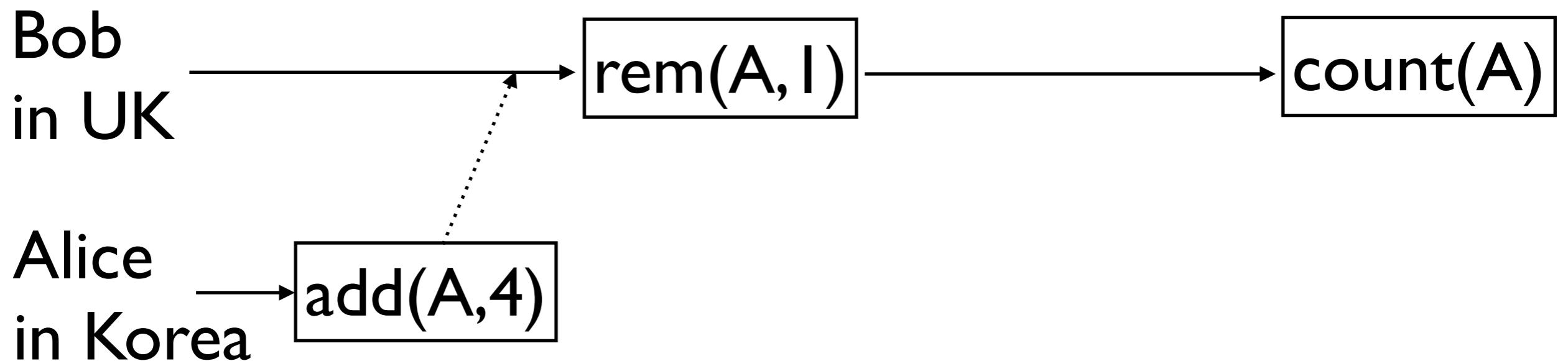
Bob
in UK

Alice
in Korea → add(A,4)

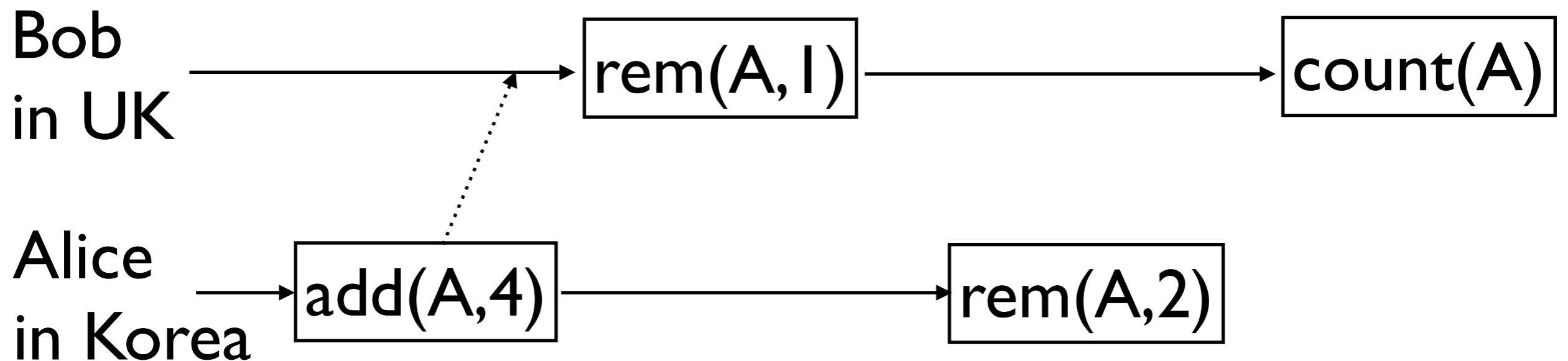
Shopping cart (Shapiro et al.)



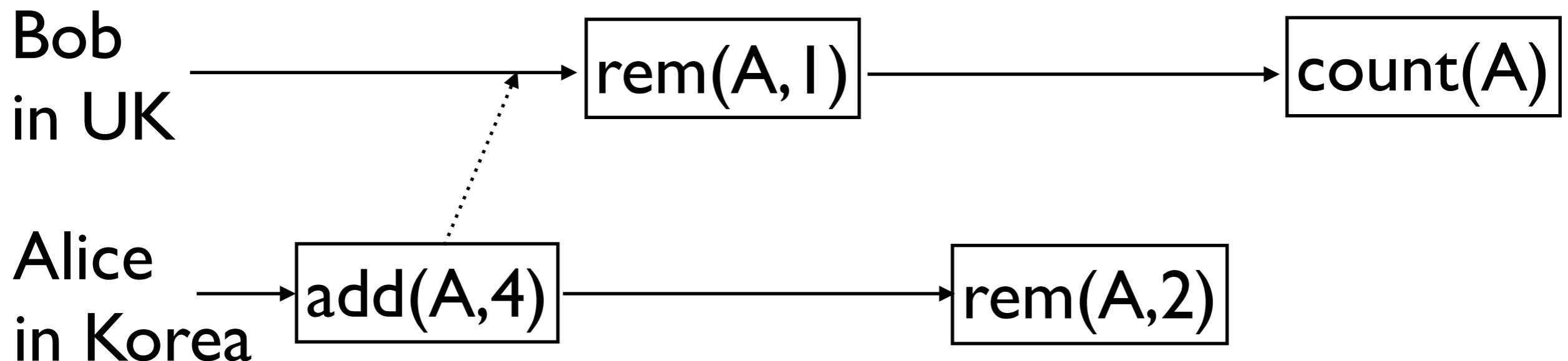
Shopping cart (Shapiro et al.)



Shopping cart (Shapiro et al.)



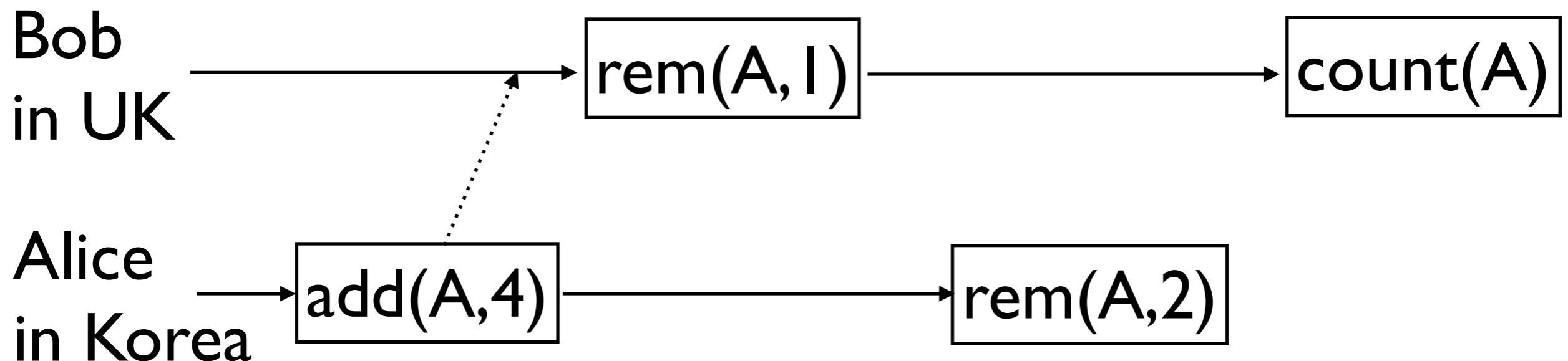
Shopping cart (Shapiro et al.)



[Q] What cannot be the result of `count(A)`?

- (a) 1
- (b) 3
- (c) 5
- (d) all possible

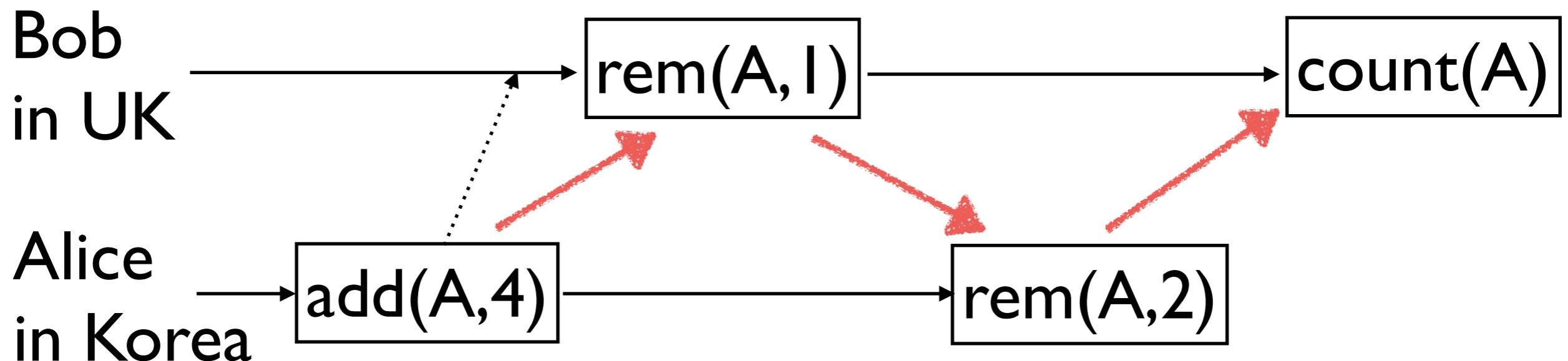
Shopping cart (Shapiro et al.)



[Q] What cannot be the result of `count(A)`?

- (a) 1
- (b) 3
- (c) 5
- (d) all possible

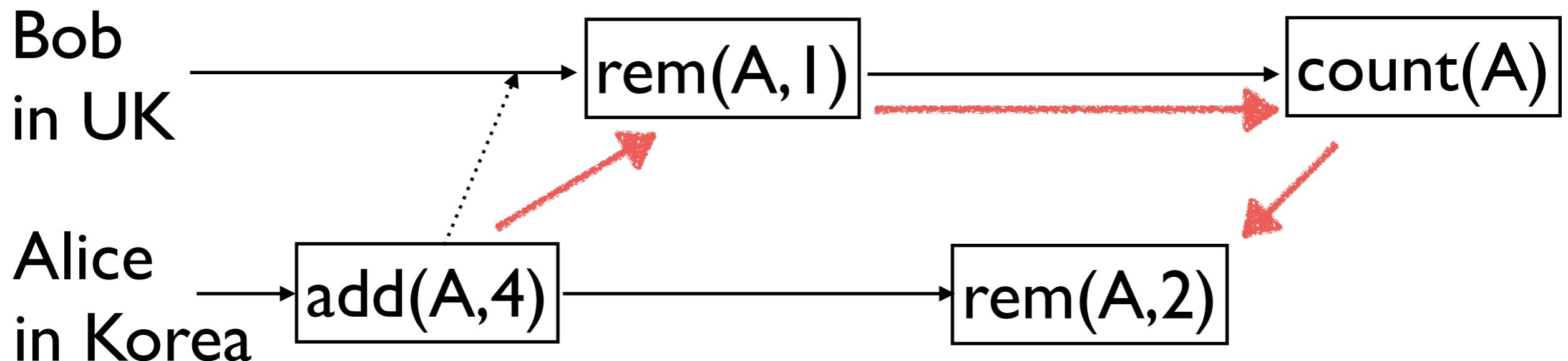
Shopping cart (Shapiro et al.)



[Q] What cannot be the result of $\text{count}(A)$?

- (a) 1
- (b) 3
- (c) 5
- (d) all possible

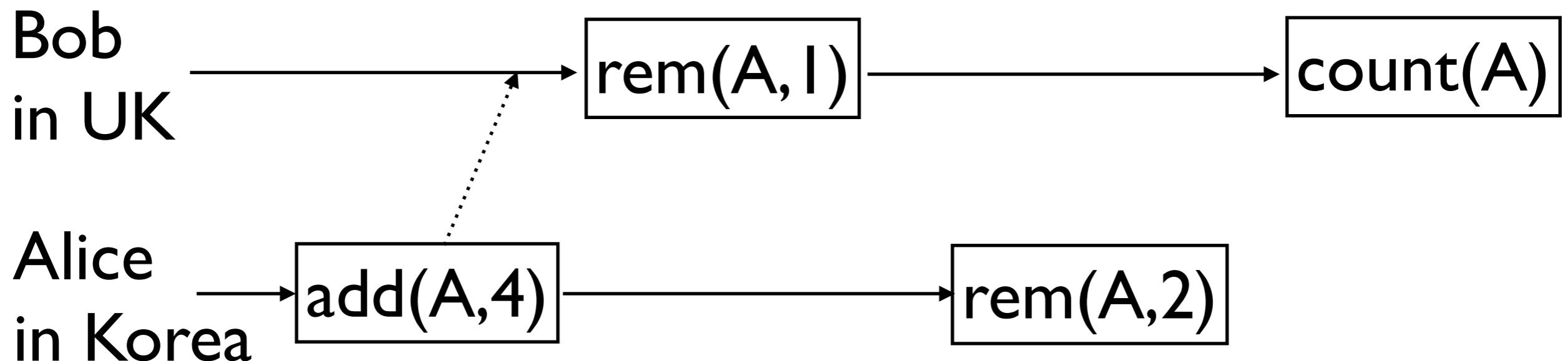
Shopping cart (Shapiro et al.)



[Q] What cannot be the result of $\text{count}(A)$?

- (a) 1
- (b) 3
- (c) 5
- (d) all possible

Shopping cart (Shapiro et al.)

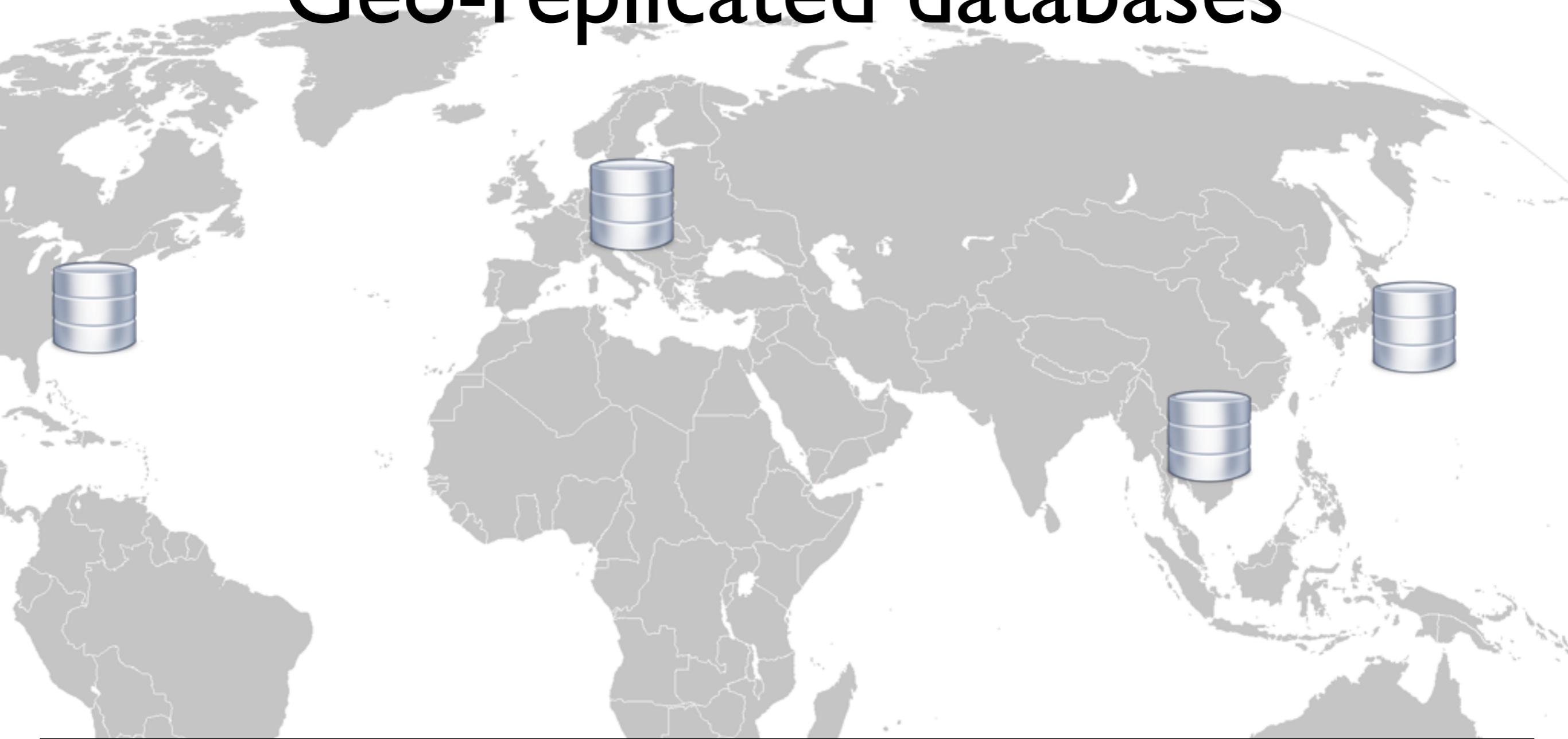


[Q] What cannot be the result of $\text{count}(A)$?

- (a) 1
- (b) 3
- (c) 5
- (d) all possible

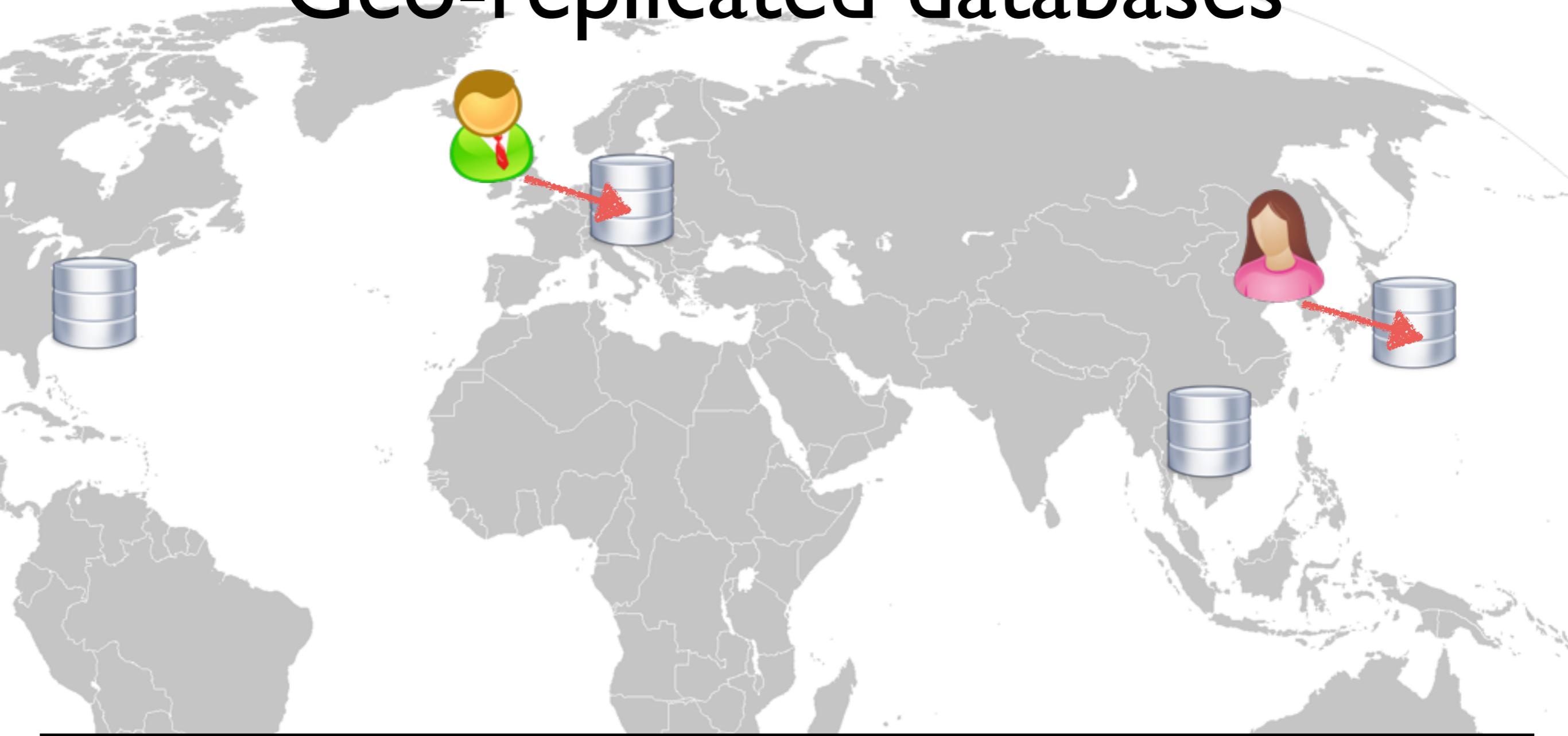
Because of weak consistency of the store.

Geo-replicated databases



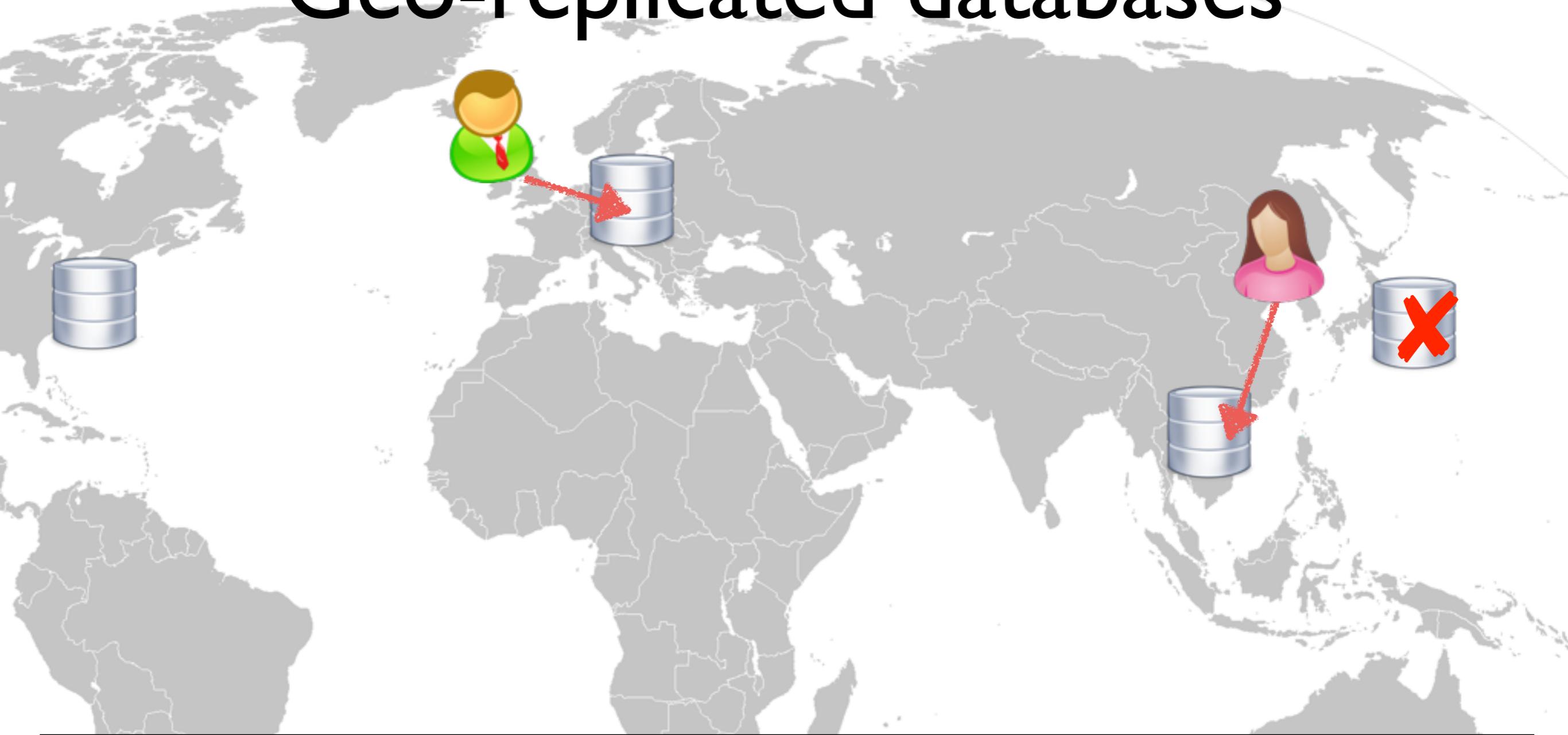
- Every data centre stores a complete replica of data
- Purpose: Minimising latency. Fault tolerance.

Geo-replicated databases



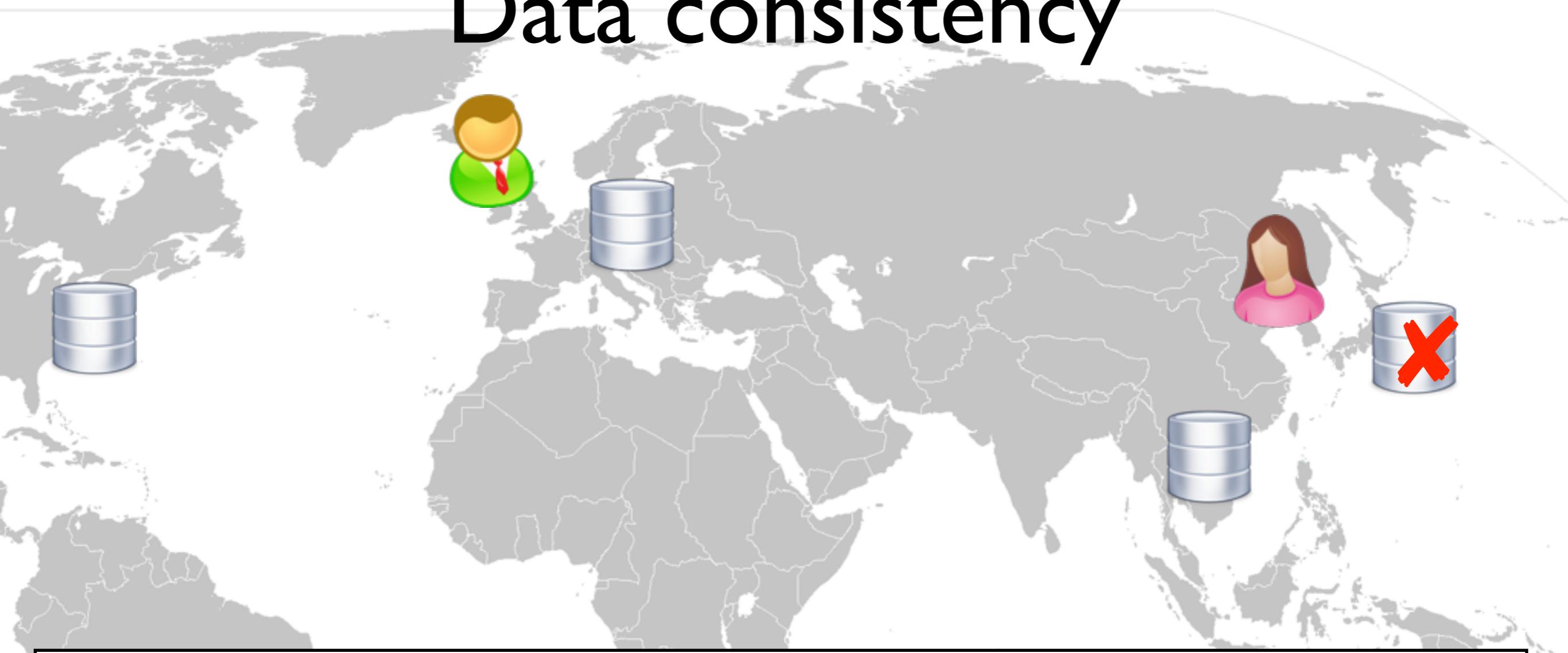
- Every data centre stores a complete replica of data
- Purpose: **Minimising latency.** Fault tolerance.

Geo-replicated databases



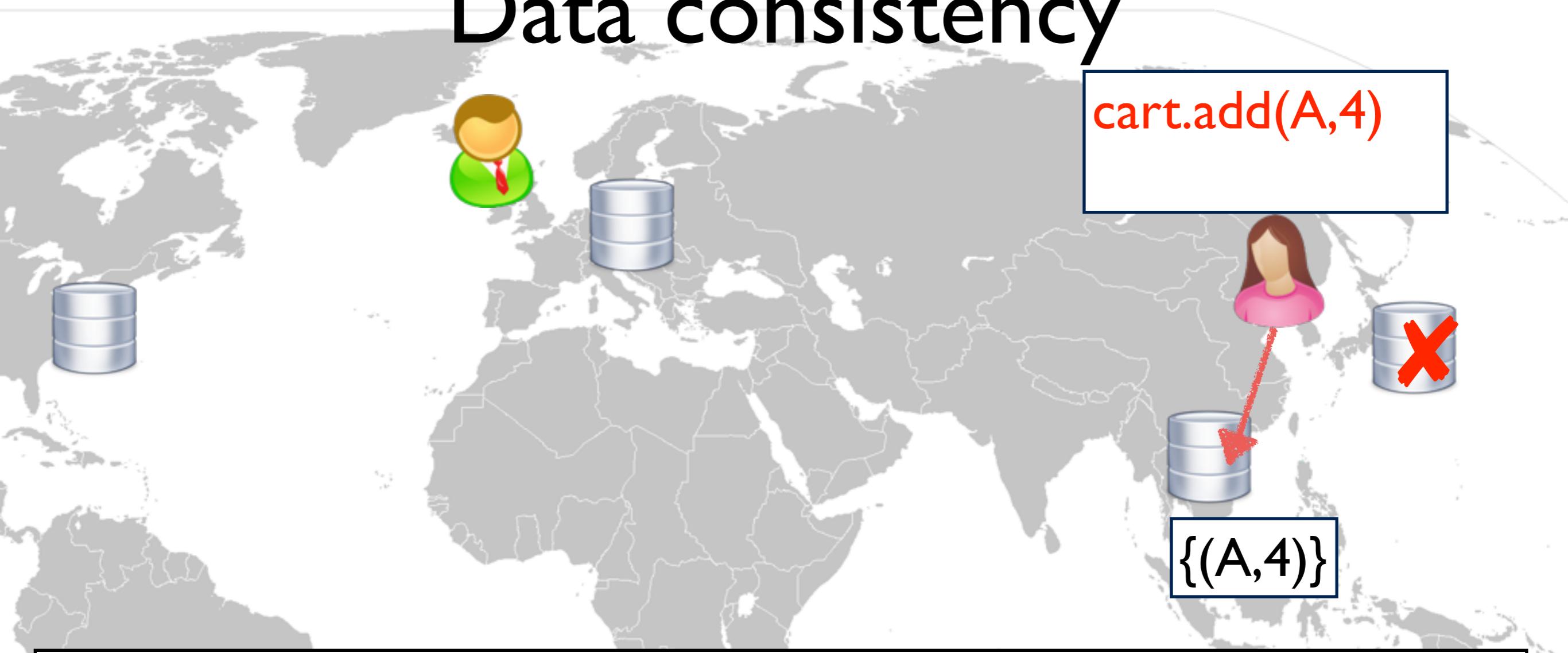
- Every data centre stores a complete replica of data
- Purpose: Minimising latency. **Fault tolerance.**

Data consistency



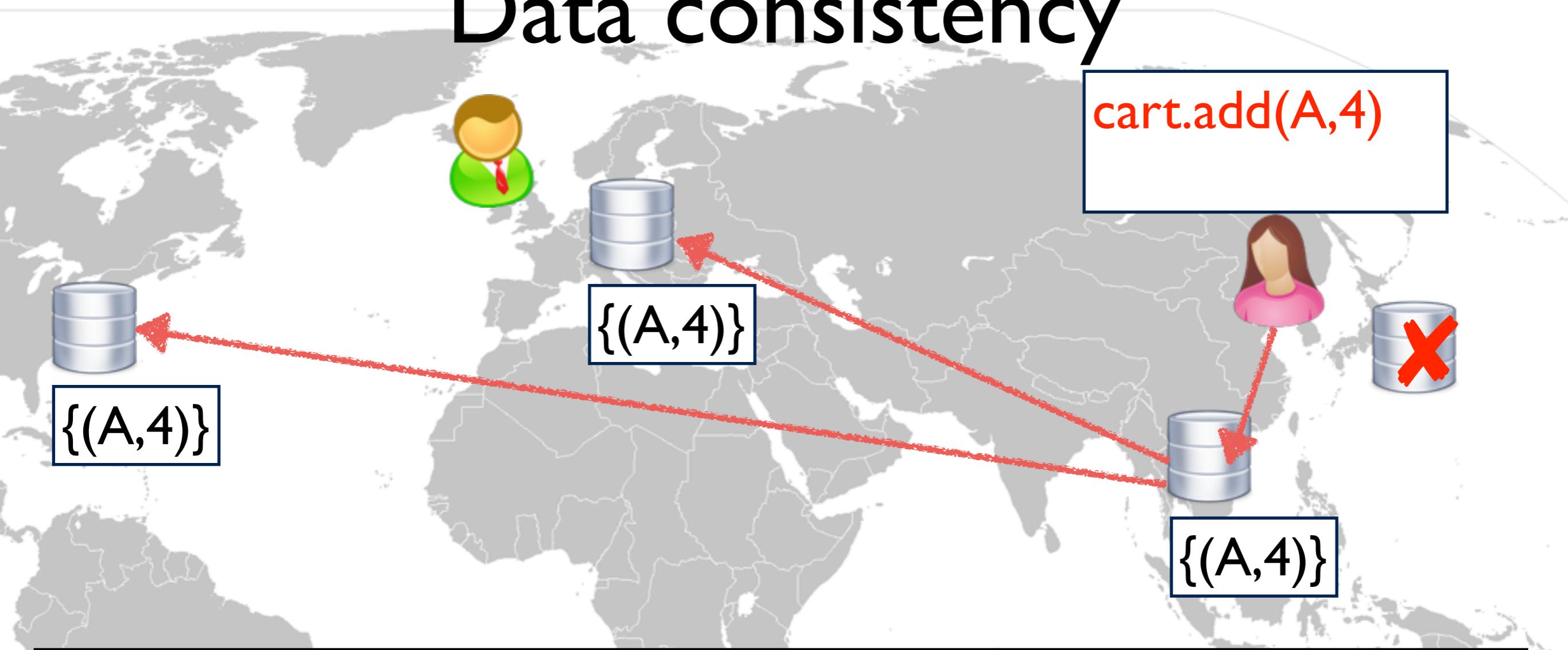
- Do data centre stores behave as if a single store?
- Strong consistency: Yes.
- Weak consistency: No.

Data consistency



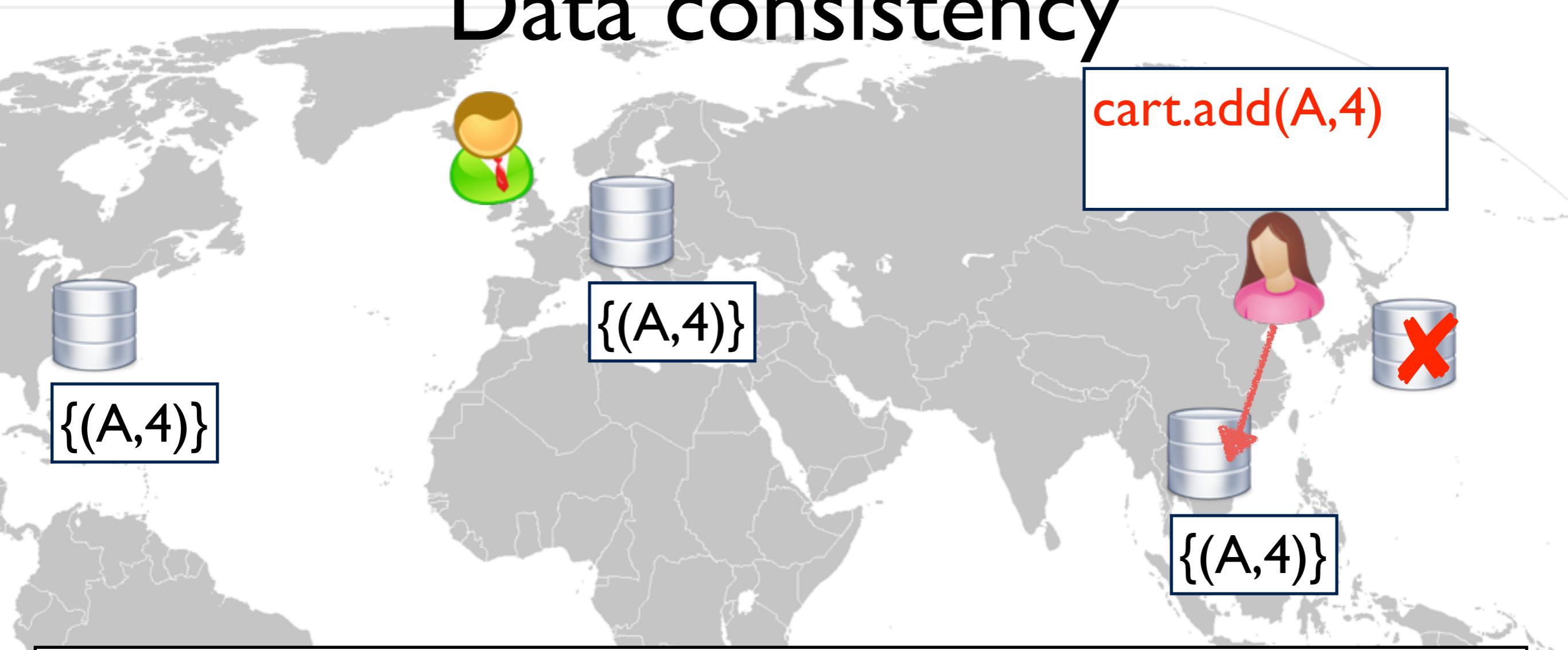
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No.

Data consistency



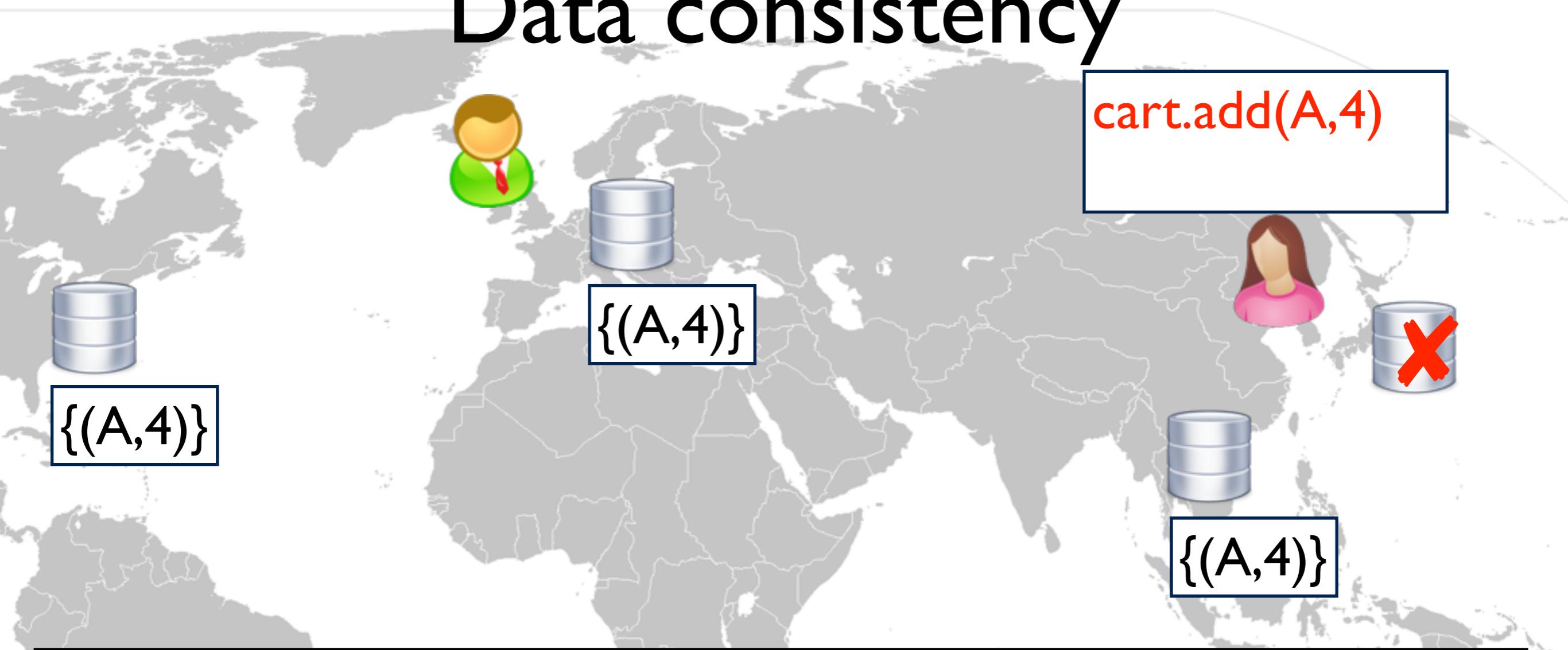
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No.

Data consistency



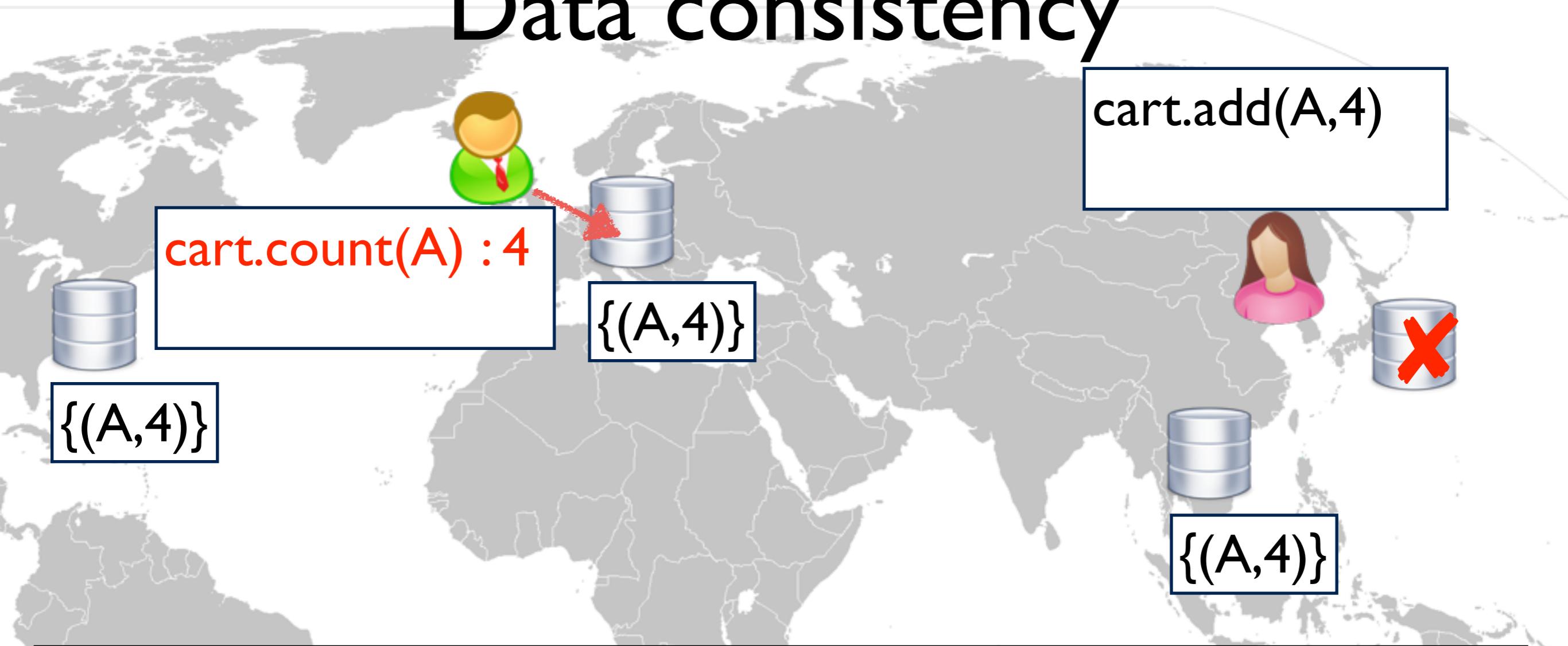
- Do data centre stores behave as if a single store?
- **Strong consistency:** Yes. Block until all get updated.
- **Weak consistency:** No.

Data consistency



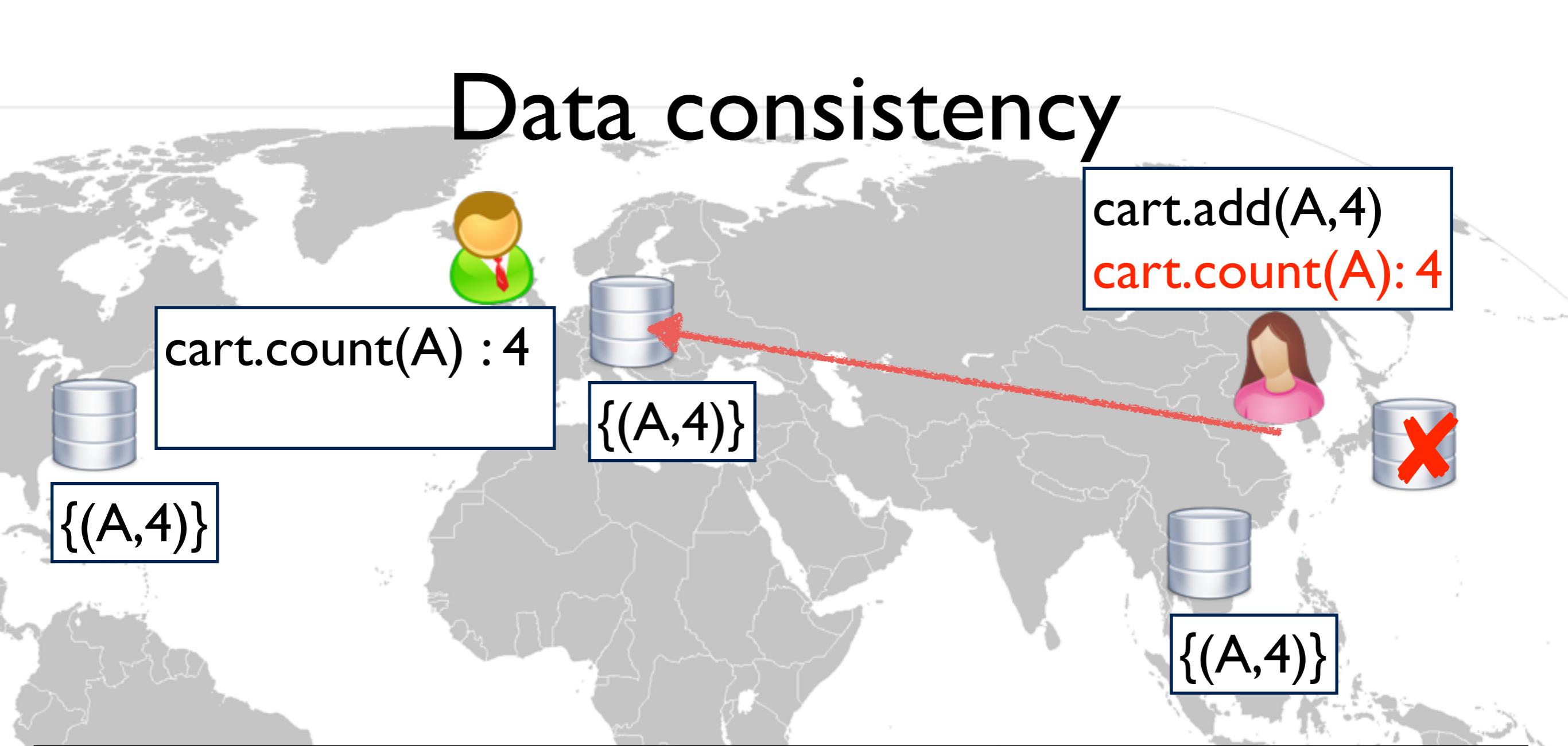
- Do data centre stores behave as if a single store?
- **Strong consistency:** Yes. Block until all get updated.
- **Weak consistency:** No.

Data consistency



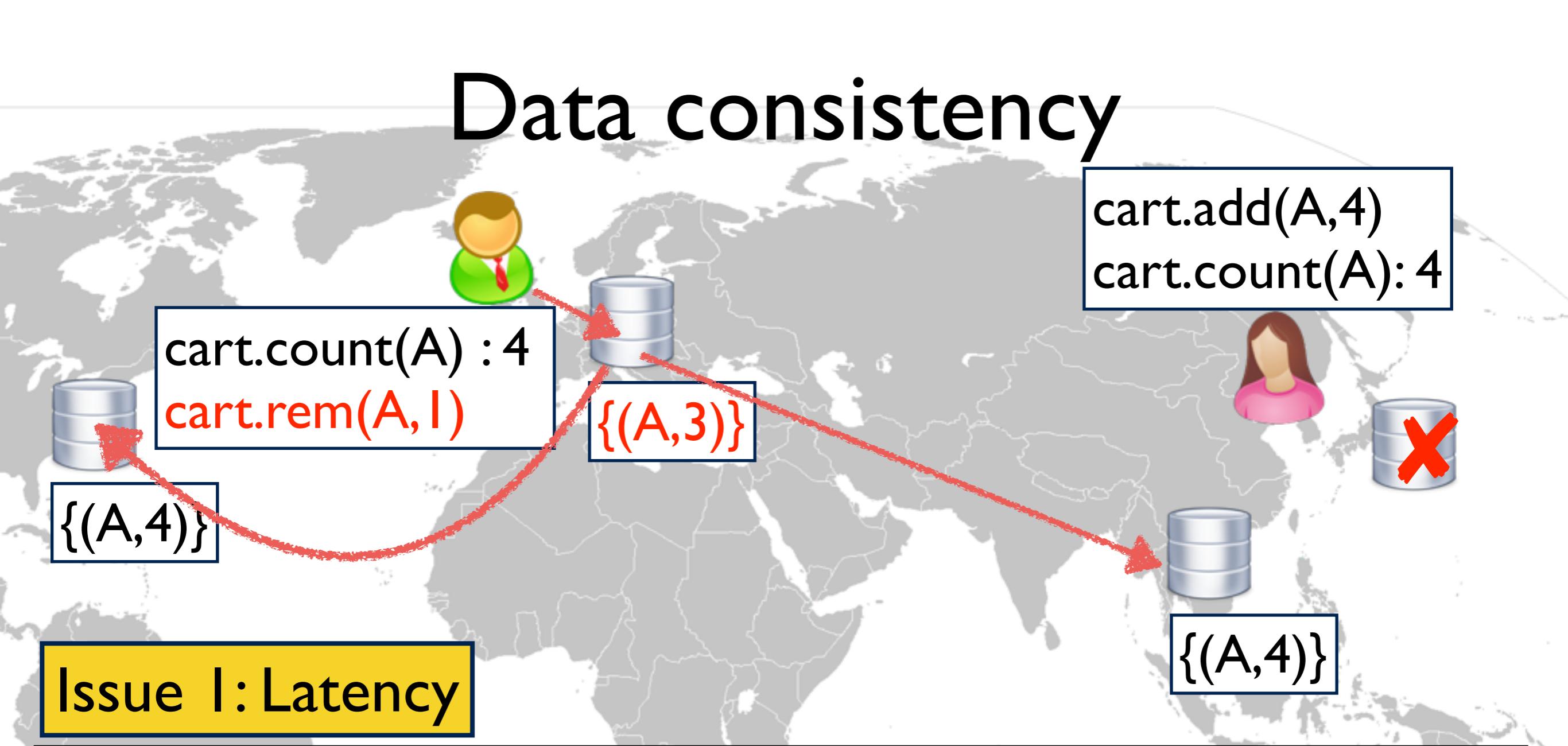
- Do data centre stores behave as if a single store?
- **Strong consistency:** Yes. Block until all get updated.
- **Weak consistency:** No.

Data consistency



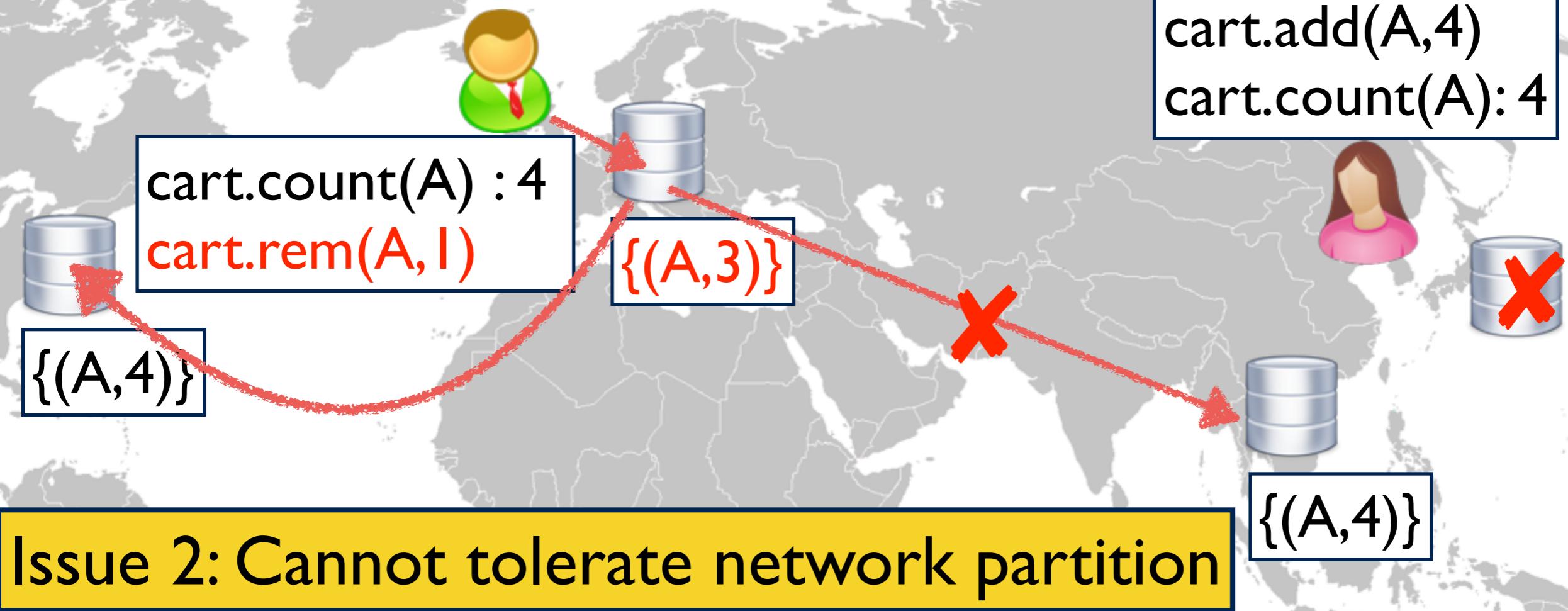
- Do data centre stores behave as if a single store?
- **Strong consistency:** Yes. Block until all get updated.
- **Weak consistency:** No.

Data consistency



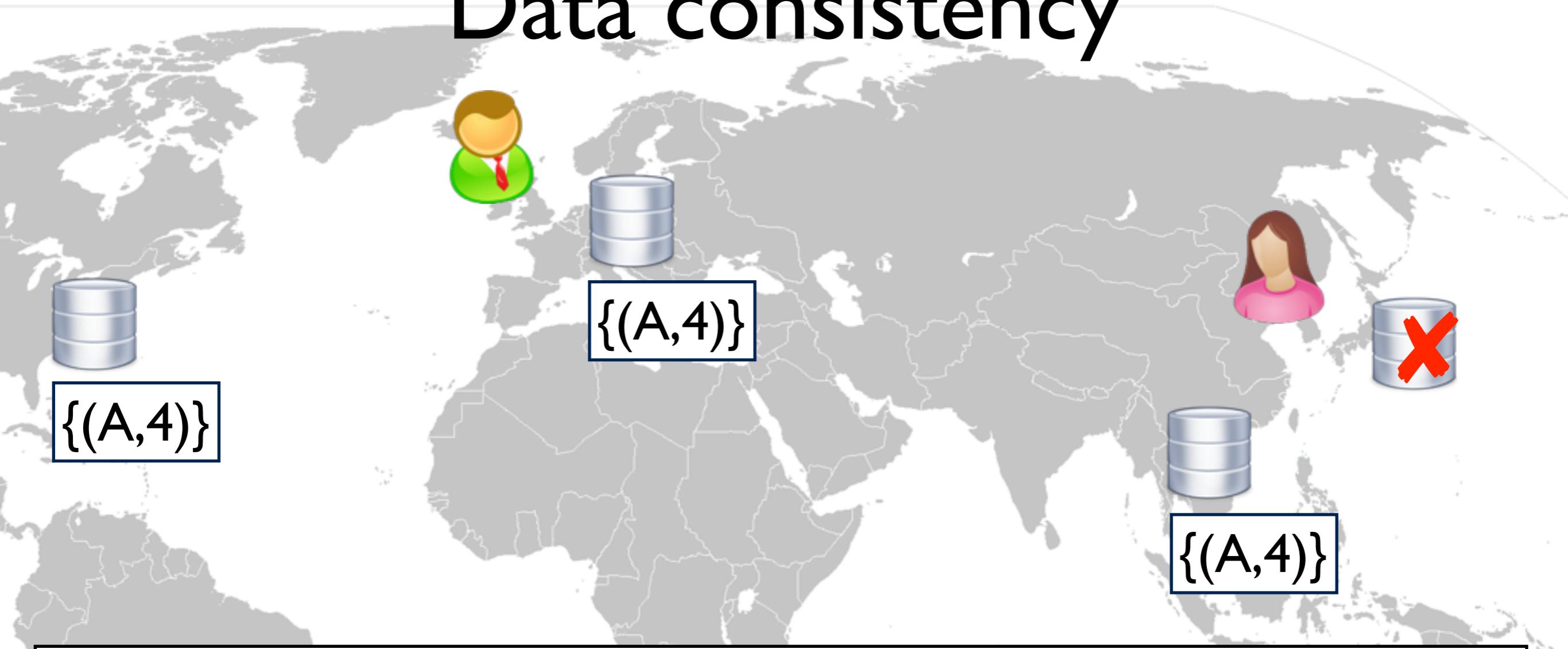
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No.

Data consistency



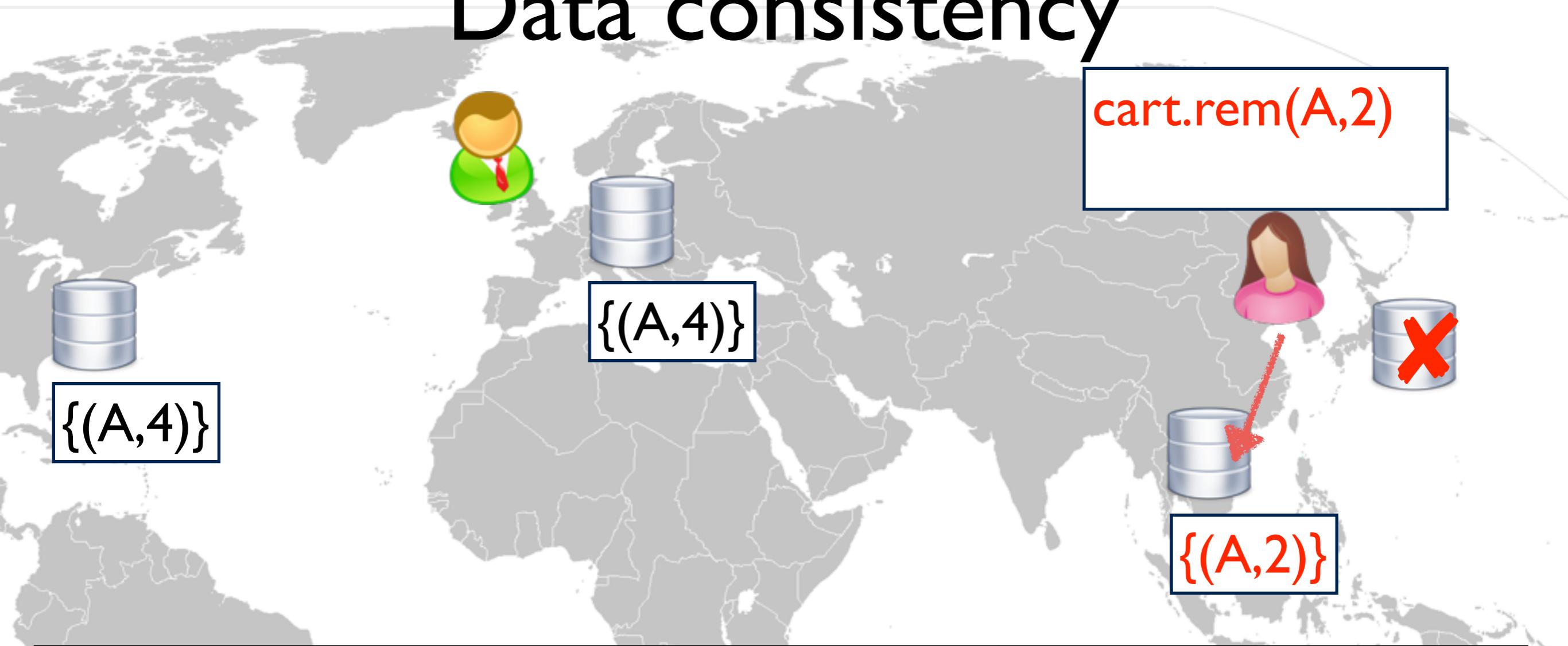
- Do data centre stores behave as if a single store?
- **Strong consistency:** Yes. Block until all get updated.
- **Weak consistency:** No.

Data consistency



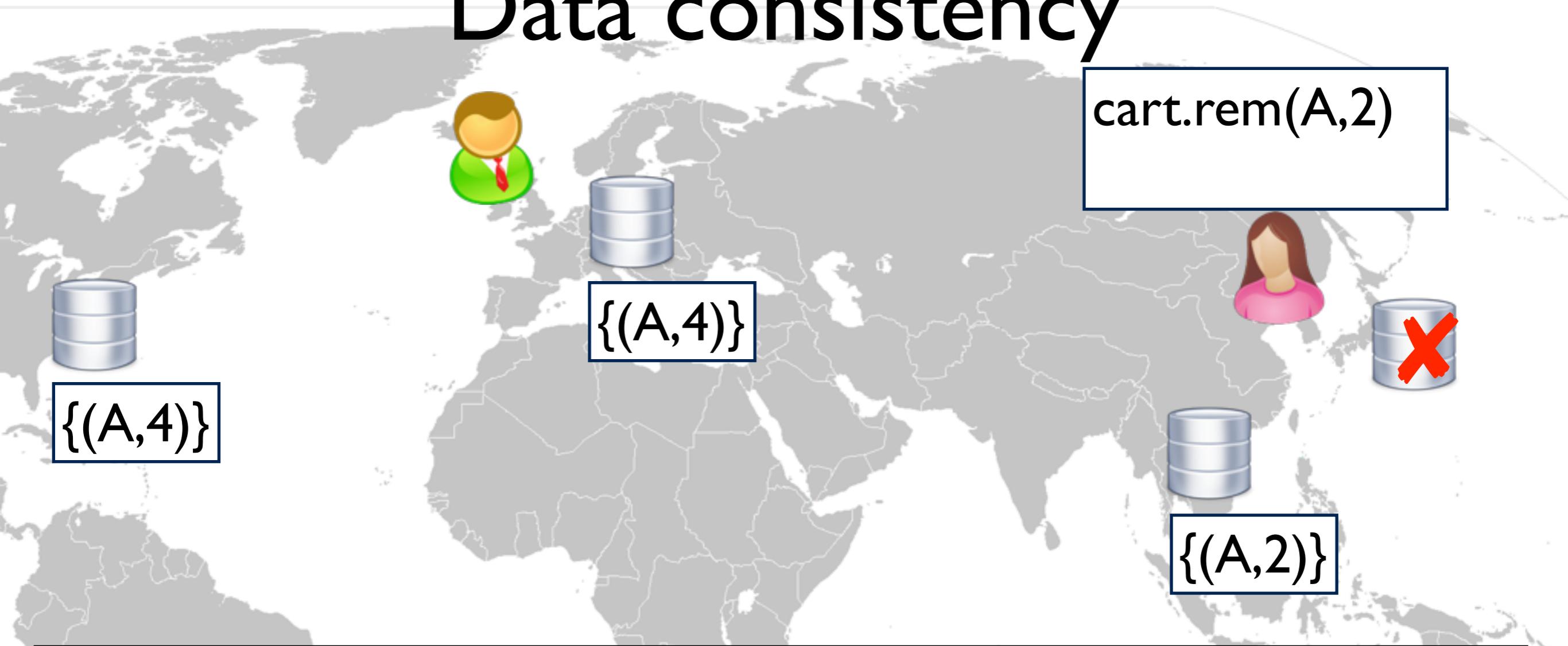
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

Data consistency



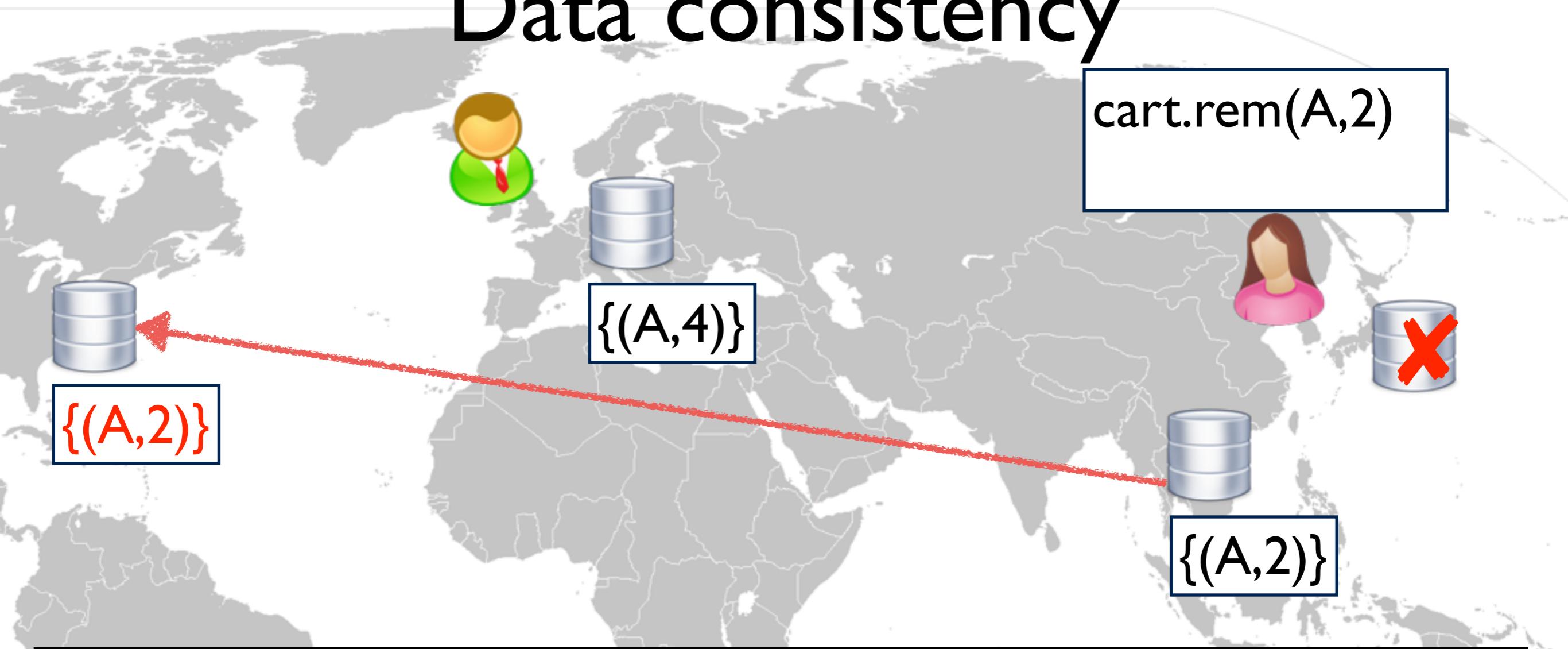
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

Data consistency



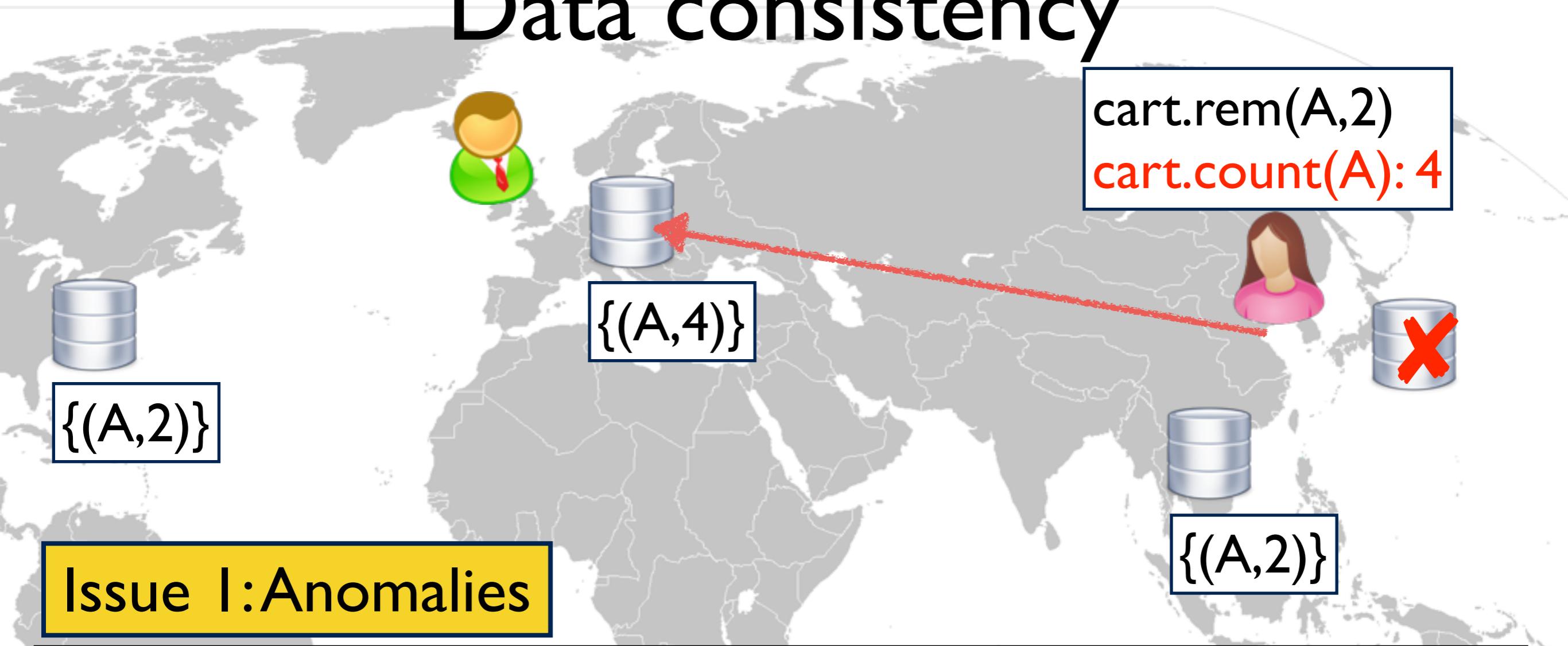
- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

Data consistency



- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

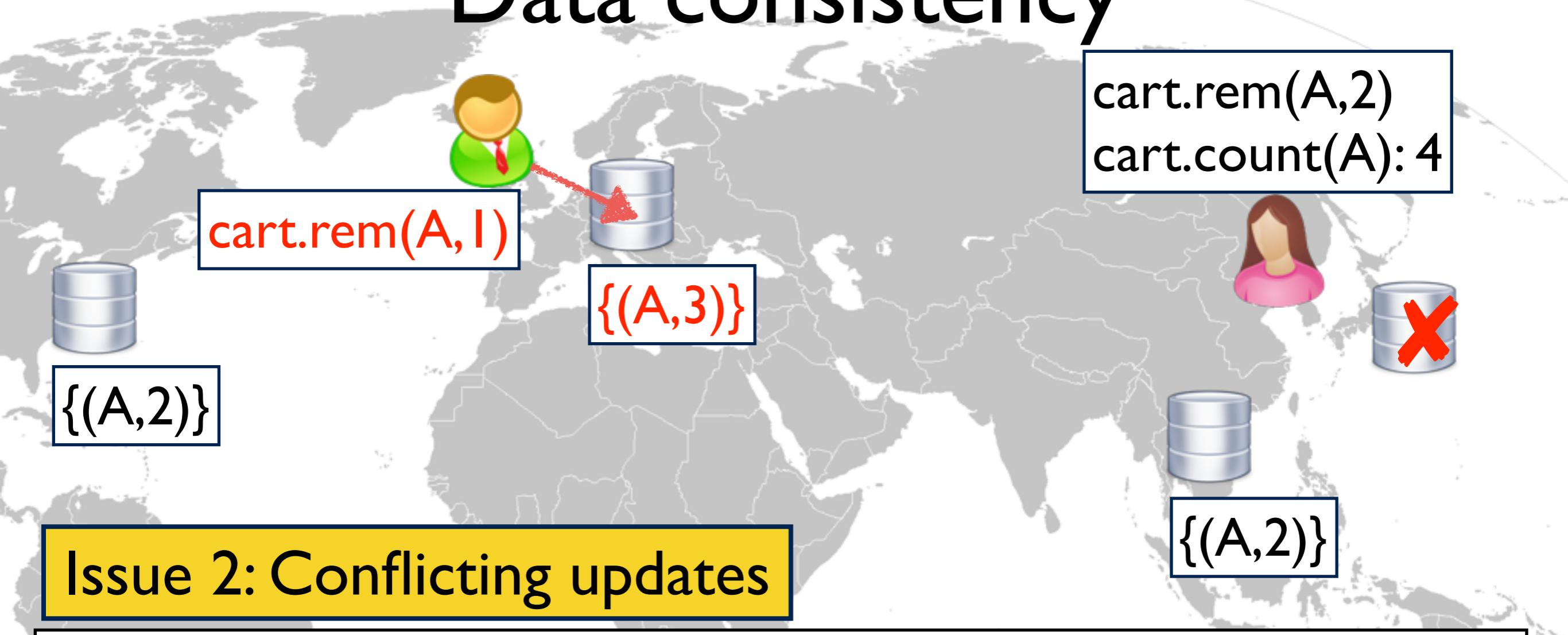
Data consistency



Issue I: Anomalies

- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

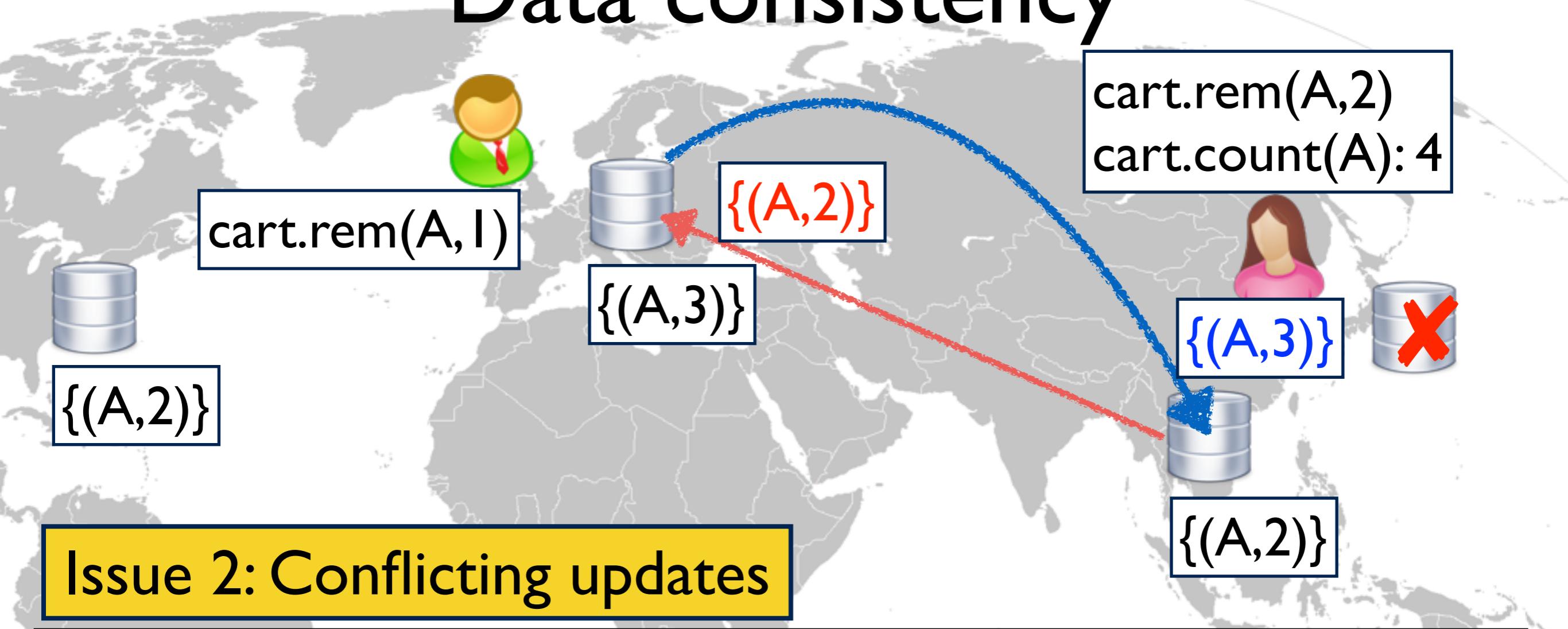
Data consistency



Issue 2: Conflicting updates

- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

Data consistency



Issue 2: Conflicting updates

- Do data centre stores behave as if a single store?
- Strong consistency: Yes. Block until all get updated.
- Weak consistency: No. First update. Then propagate.

Popularity of weak consistency

- Dynamo/Simple DB (Amazon), Riak, Cassandra (Facebook/Twitter), CouchDB, Google Doc,etc
- Due to better responsiveness and availability.
- But created a programming challenge.

Shopping cart

- Map from book ids to natural numbers.
- Initially, constant-0 function.
- $\text{add} : \{0, \dots, N-1\} \times \text{Nat} \rightarrow \text{Unit}$
- $\text{rem} : \{0, \dots, N-1\} \times \text{Nat} \rightarrow \text{Unit}$
- $\text{count} : \{0, \dots, N-1\} \rightarrow \text{Nat}$

Sequential data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```

Sequential data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```

Sequential data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```

Sequential data type

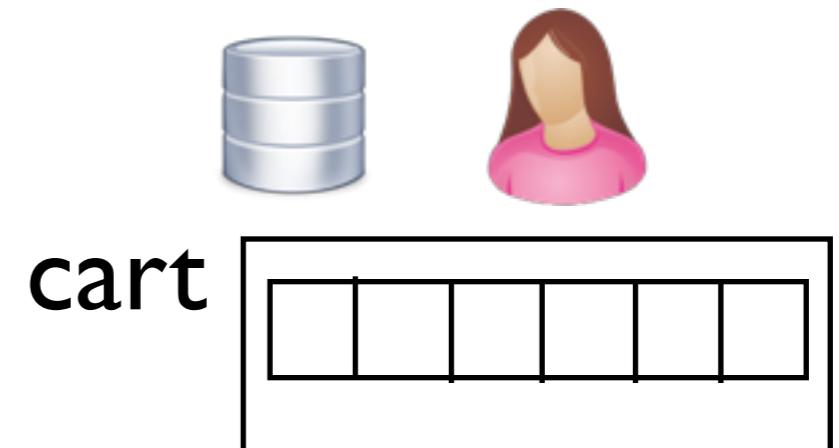
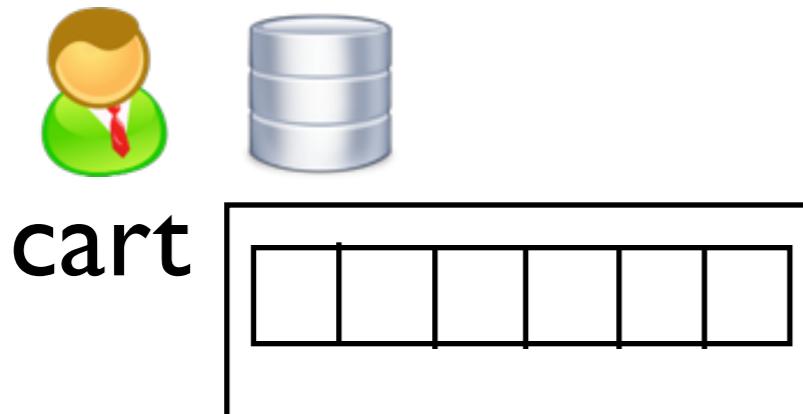
```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```

Sequential data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```

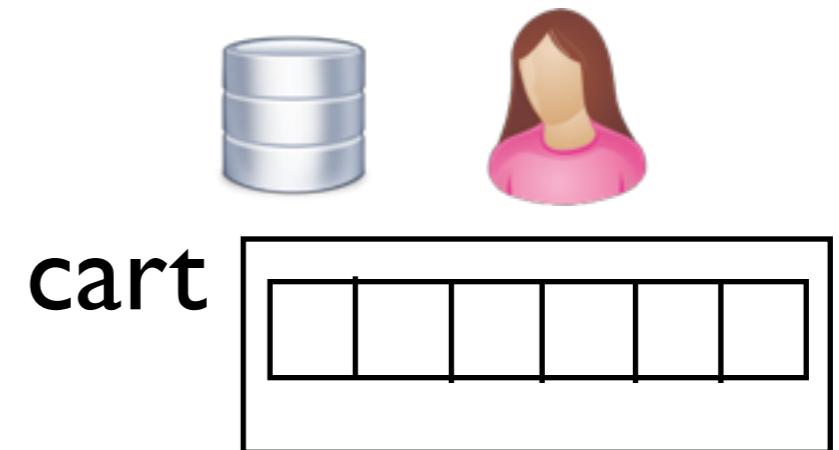
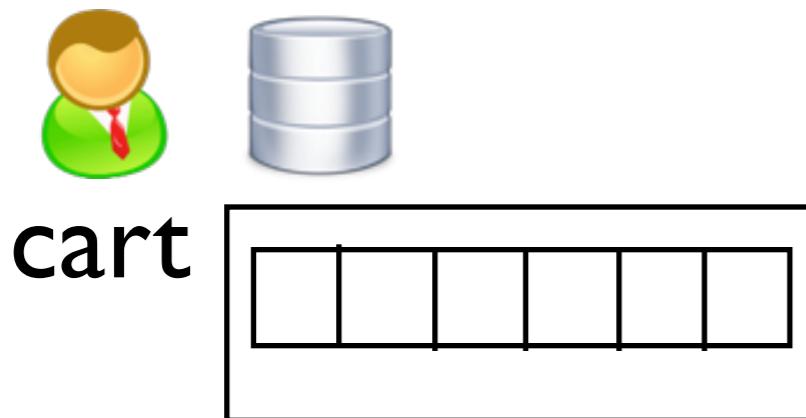
Replicated ~~Sequential~~✓ data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```



Replicated ~~Sequential~~ data type

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b,n) = { val m = map(b); map(b) = m+n }  
    def rem(b,n) = { val m = map(b); map(b) = max(m-n,0) }  
    def count(b) = { val m = map(b); return m }  
}
```



Things to be added:

1. Asynchronous message sending & receiving.
2. Conflict detection and resolution.

Primitive replicated data types (AKA CRDT) [Shapiro+ 2011]

- Designed for weakly consistent replicated stores.
- Send & receive update messages asynchronously.
- Detect and resolve conflicts.
- Register, set, dag, graph, etc.
- Conceptually, implement a weak shared memory on top of a replicated store.

Multi-valued register

$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

- Implements a memory cell of type T .

Multi-valued register

$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

- Implements a memory cell of type T .
- Resolves conflicts by taking union.

Multi-valued register

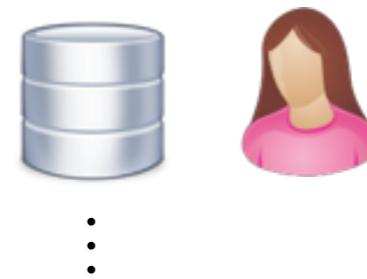
$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

- Implements a memory cell of type T .
- Resolves conflicts by taking union.



:

x.set(1)



:

x.set(2)

Multi-valued register

$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

- Implements a memory cell of type T .
- Resolves conflicts by taking union.



:

$x.\text{set}(1)$

:

$x.\text{get}() // \{1,2\}$



:

$x.\text{set}(2)$



Multi-valued register

`set : T → Unit` `get : Unit → Set[T]`

- Implements a memory cell of type T.
- Resolves conflicts by taking union.



⋮

`x.set(1)`

⋮

`x.get() // {1,2}`

⋮

`x.set(3)`



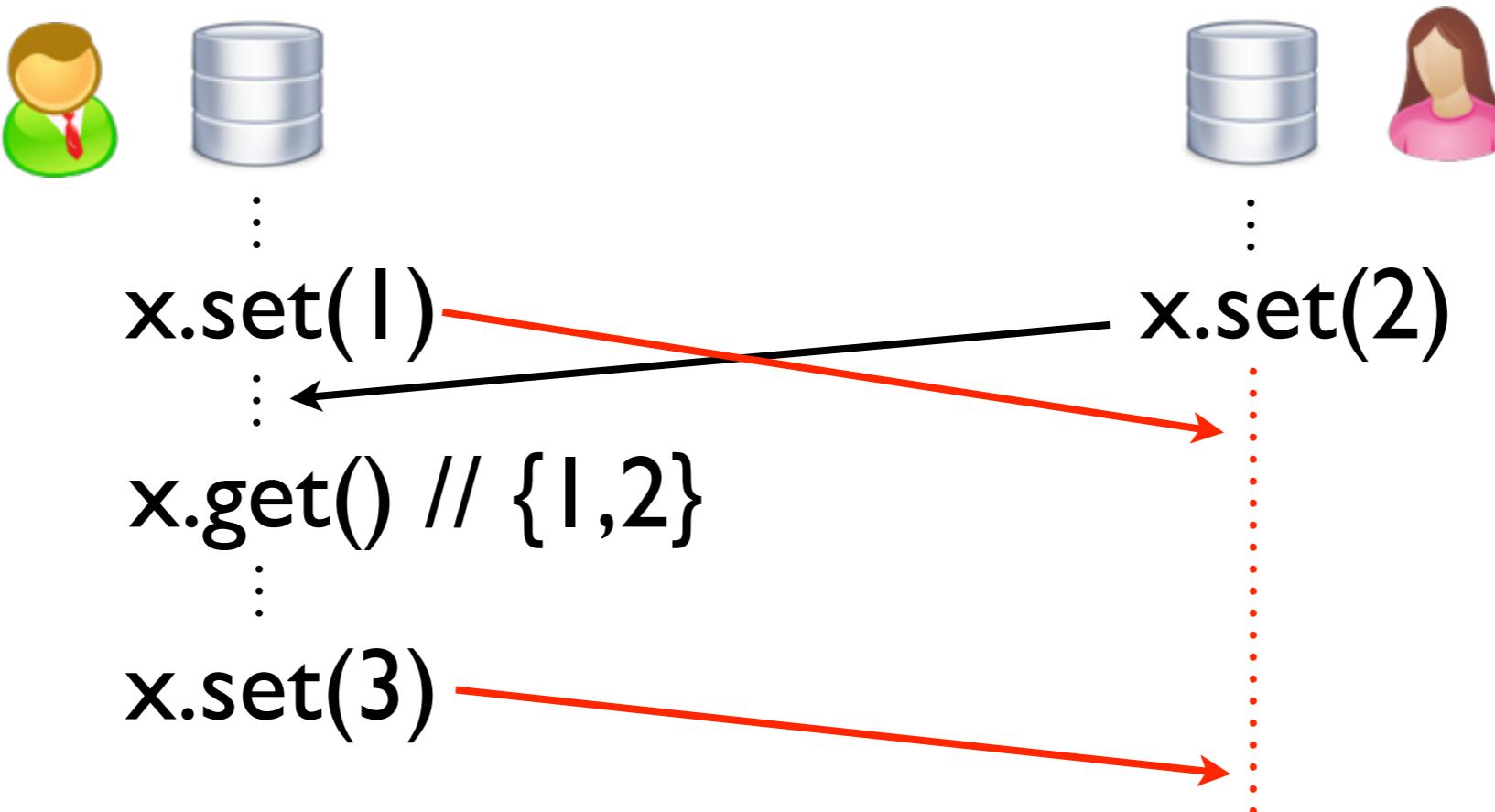
⋮

`x.set(2)`

Multi-valued register

$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

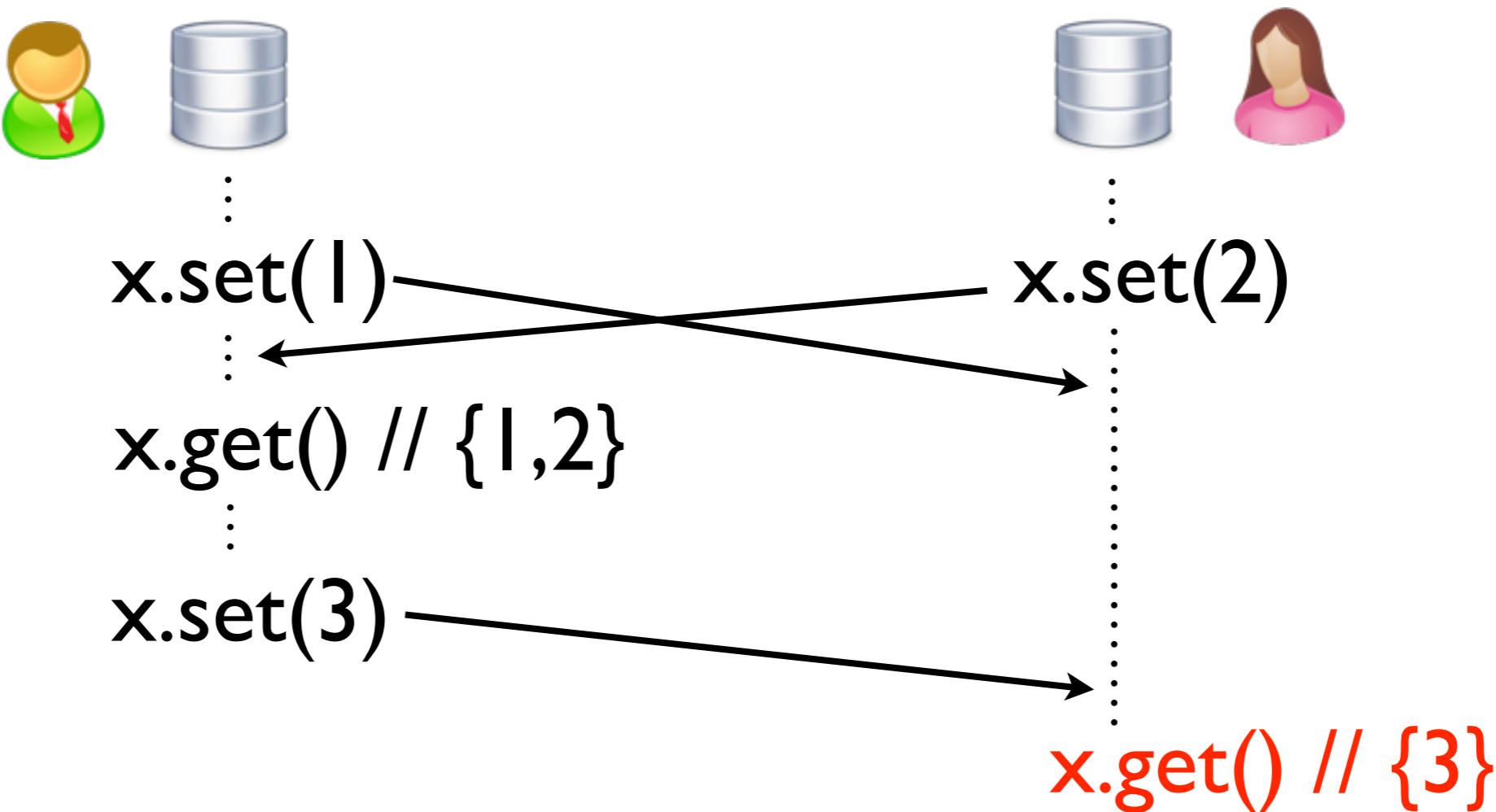
- Implements a memory cell of type T .
- Resolves conflicts by taking union.



Multi-valued register

$\text{set} : T \rightarrow \text{Unit}$ $\text{get} : \text{Unit} \rightarrow \text{Set}[T]$

- Implements a memory cell of type T .
- Resolves conflicts by taking union.

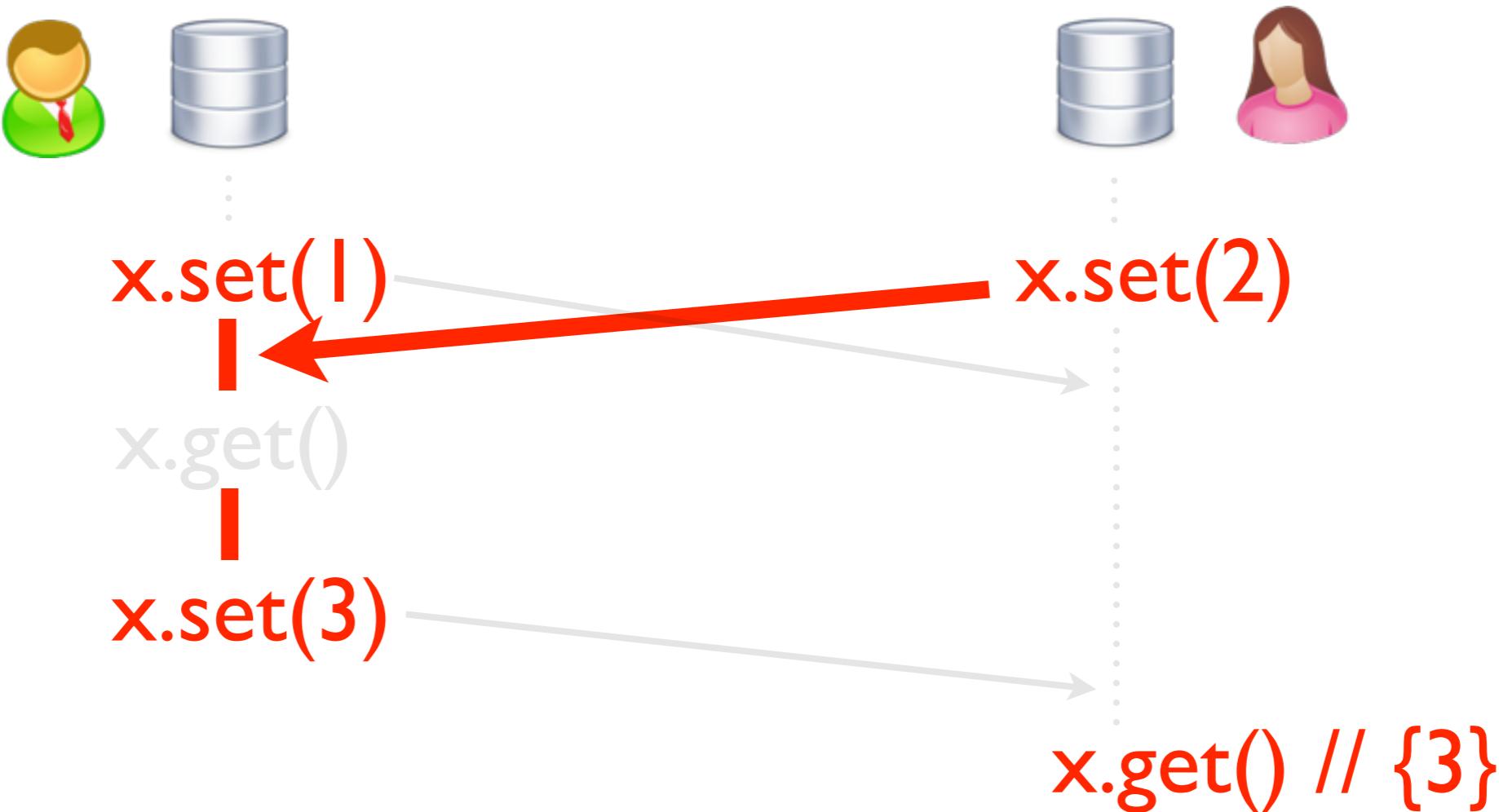


Multi-valued register

`set : T → Unit`

`get : Unit → Set[T]`

- Implements a memory cell of type T.
- Resolves conflicts by taking union.



Shopping cart for replicated stores

```
class Cart {  
    val map = new Array[Nat](N)  
  
    def add(b, n) = {  
        val m = map(b)  
        map(b) = m+n  
    }  
    def rem(b, n) = {  
        val m = map(b)  
        map(b) = max(m-n, 0)  
    }  
    def count(b) = {  
        val m = map(b)  
        return m  
    }  
}
```

Shopping cart for replicated stores

```
class Cart {  
    val map = new Array[MVRegister[Nat]](N)  
    for(i <- 0 until N) map(i) = new MVRegister[Nat]  
  
    def add(b, n) = {  
        val m = map(b).get()  
        map(b).set(m+n)  
    }  
    def rem(b, n) = {  
        val m = map(b).get()  
        map(b).set(max(m-n, 0))  
    }  
    def count(b) = {  
        val m = map(b).get()  
        return m  
    }  
}
```

I. Use MVRegister.

Shopping cart for replicated stores

```
class Cart {  
    val map = new array[MVRegister[Nat]](N)  
    for(i <- 0 until N) map(i) = new MVRegister[Nat]  
  
    def add(b,n) = {  
        val m = sum(map(b).get())  
        map(b).set(m+n)  
    }  
    def rem(b,n) = {  
        val m = sum(map(b).get())  
        map(b).set(max(m-n,0))  
    }  
    def count(b) = {  
        val m = sum(map(b).get())  
        return m  
    }  
}
```

- I. Use MVRegister.
2. Add conflict resolution.

Shopping cart for replicated stores

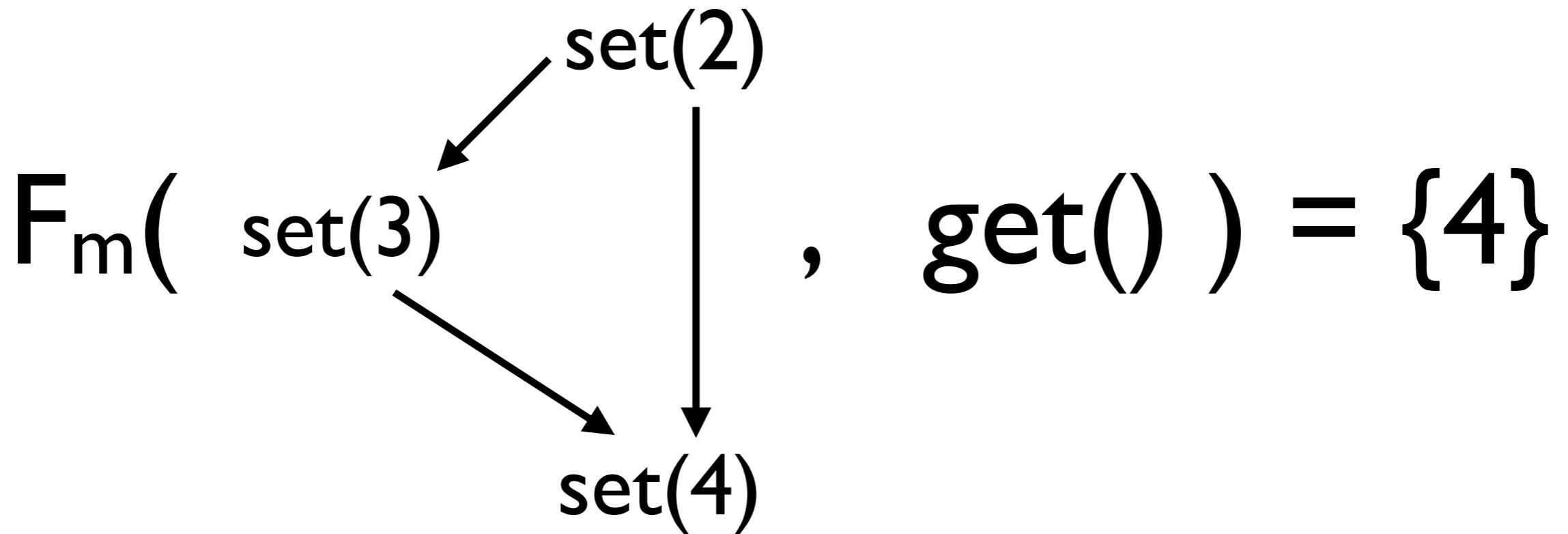
```
class Cart {  
    val map = new array[MVRegister[Nat]](N)  
    for(i <- 0 until N) map(i) = new MVRegister[Nat]  
  
    def add(b,n) = {  
        val m = max(map(b).get())  
        map(b).set(m+n)  
    }  
    def rem(b,n) = {  
        val m = max(map(b).get())  
        map(b).set(max(m-n,0))  
    }  
    def count(b) = {  
        val m = max(map(b).get())  
        return m  
    }  
}
```

- I. Use MVRegister.
2. Add conflict resolution.

Semantics of replicated data types

$\text{Dag}[\text{Op}] = \text{Set of Op-labelled finite dags}$

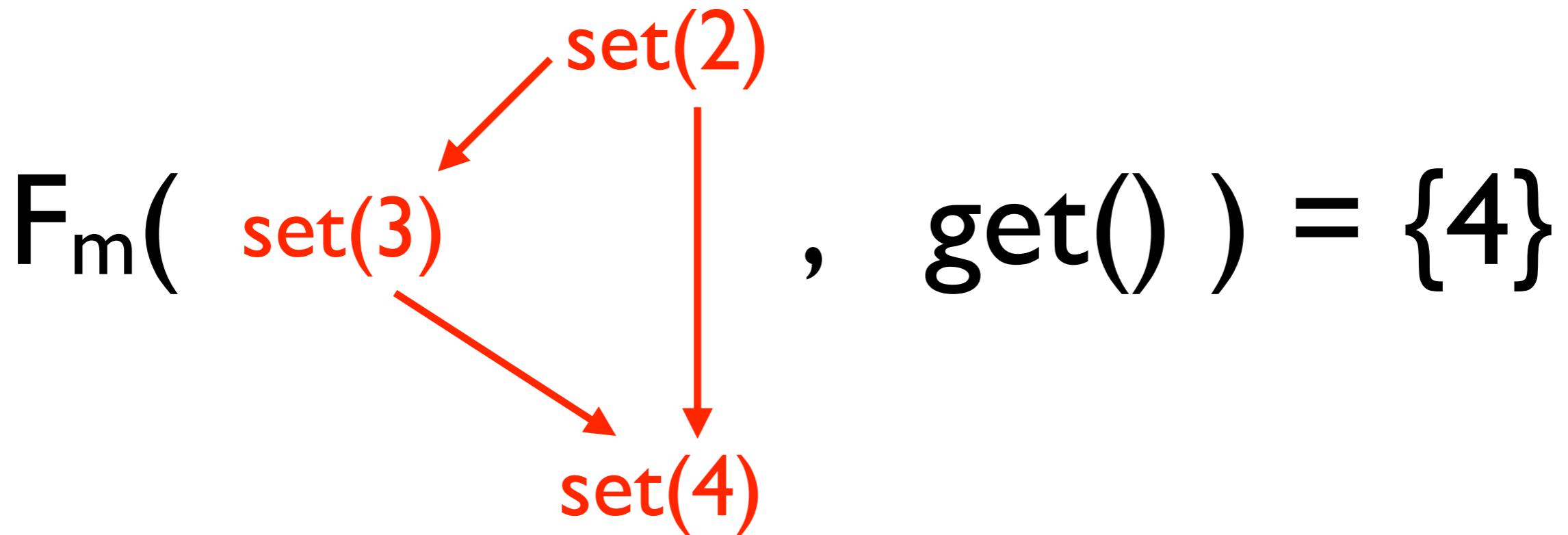
$F : \text{Dag}[\text{Op}] \times \text{Op} \rightarrow_{\text{partial}} \text{Value}$



Semantics of replicated data types

$\text{Dag}[\text{Op}] = \text{Set of Op-labelled finite dags}$

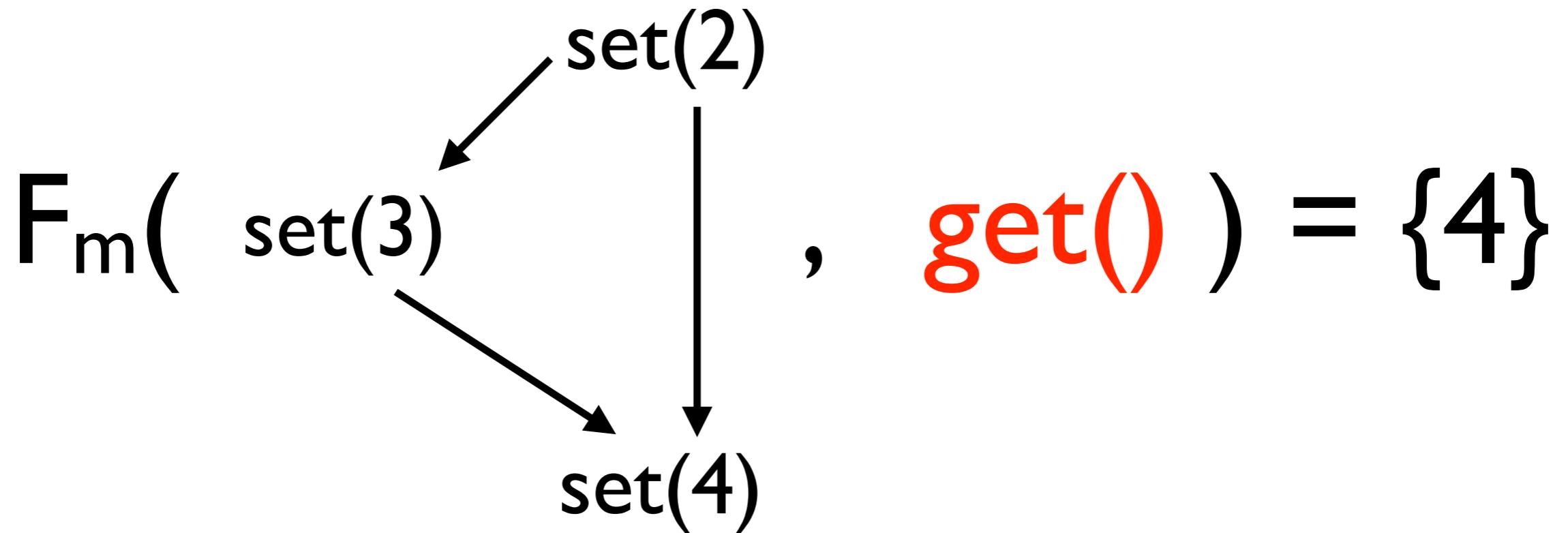
$F : \text{Dag}[\text{Op}] \times \text{Op} \rightarrow_{\text{partial}} \text{Value}$



Semantics of replicated data types

$\text{Dag}[\text{Op}] = \text{Set of Op-labelled finite dags}$

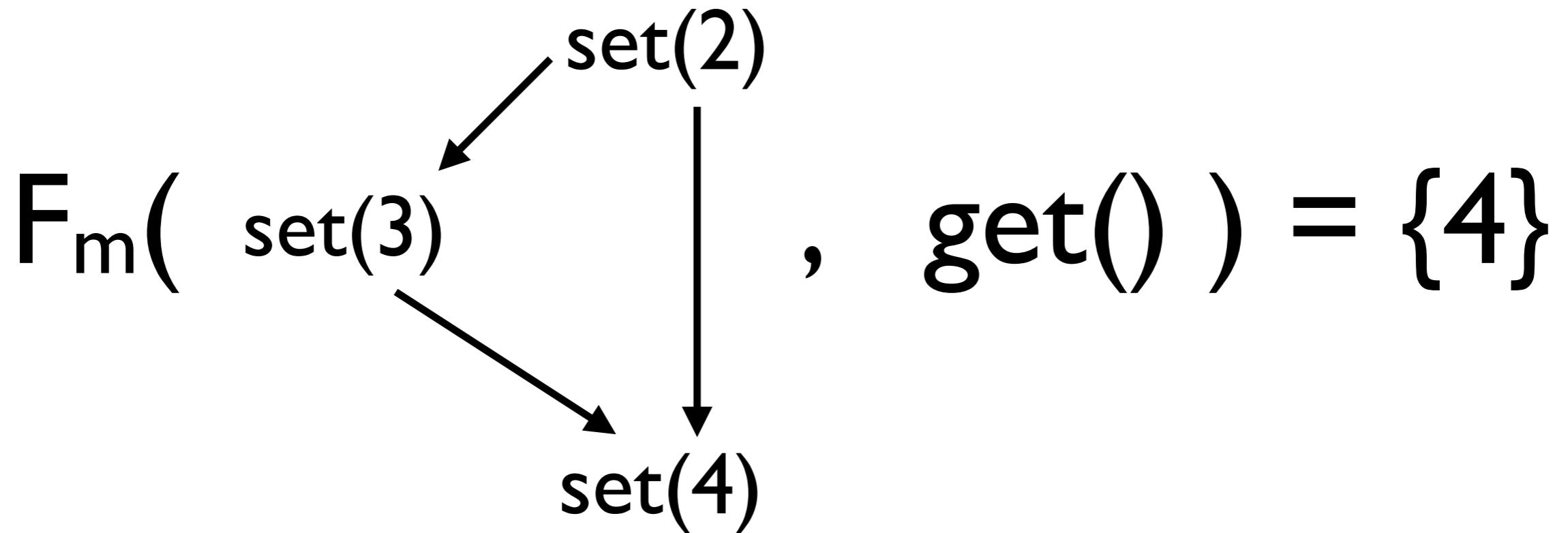
$F : \text{Dag}[\text{Op}] \times \text{Op} \rightarrow_{\text{partial}} \text{Value}$



Semantics of replicated data types

$\text{Dag}[\text{Op}] = \text{Set of Op-labelled finite dags}$

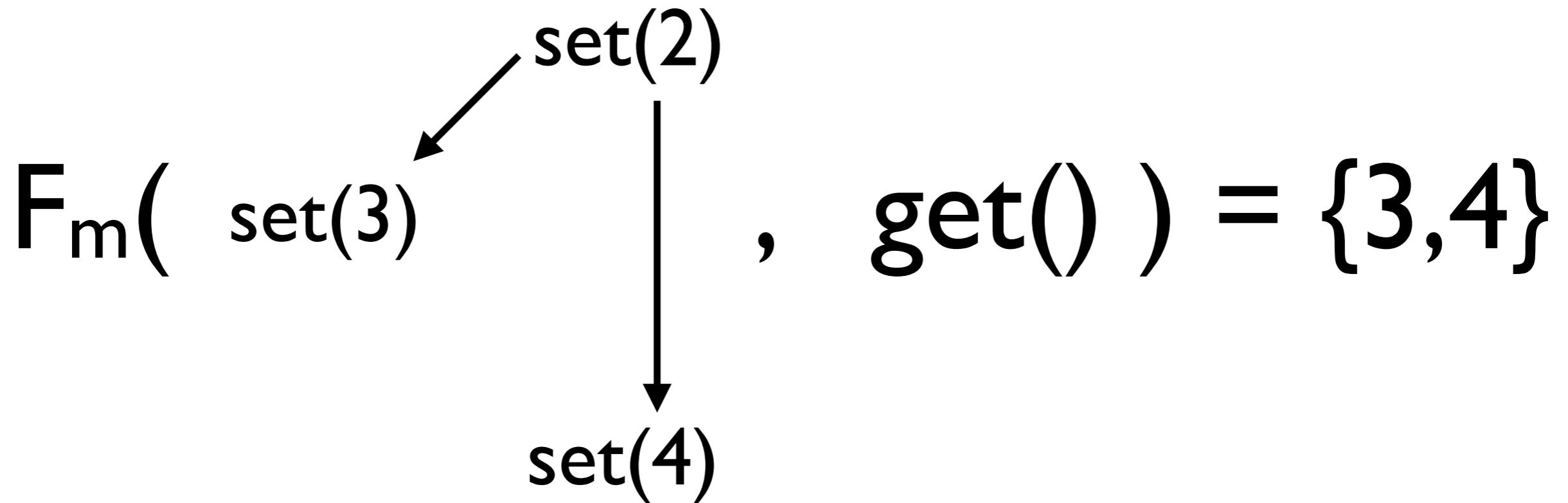
$F : \text{Dag}[\text{Op}] \times \text{Op} \rightarrow_{\text{partial}} \text{Value}$



Semantics of replicated data types

$\text{Dag}[\text{Op}] = \text{Set of Op-labelled finite dags}$

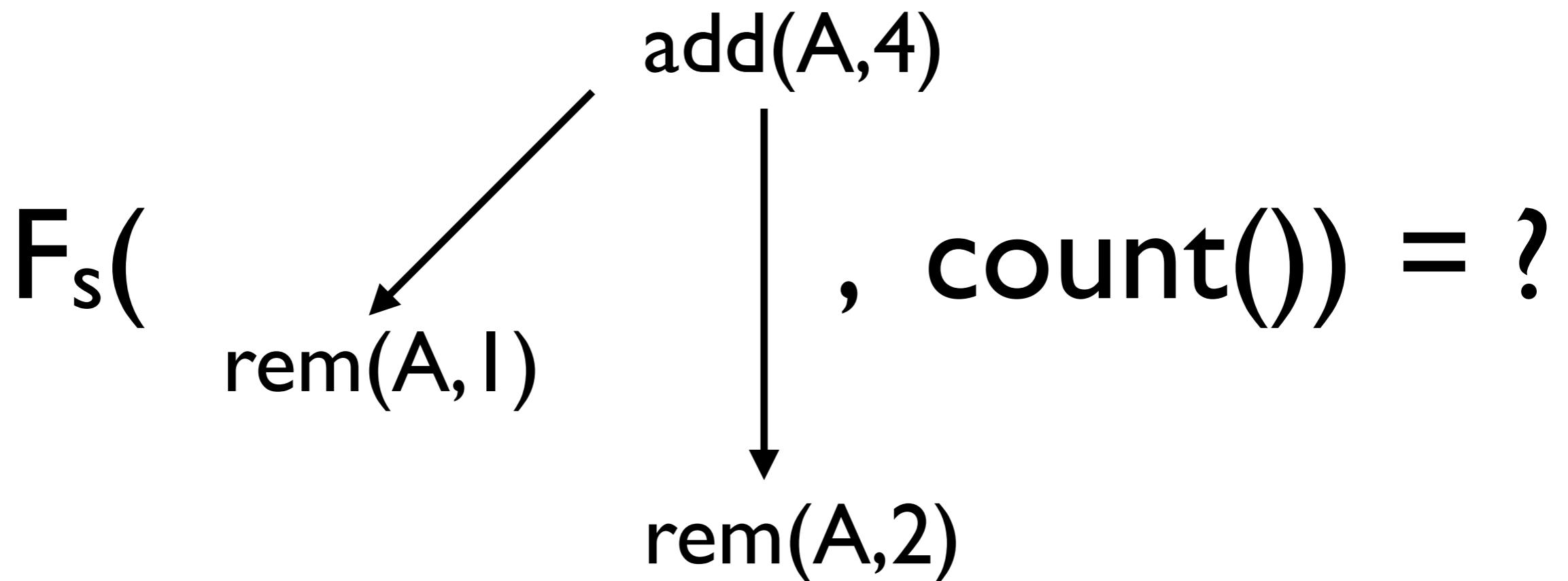
$F : \text{Dag}[\text{Op}] \times \text{Op} \rightarrow_{\text{partial}} \text{Value}$



```

val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

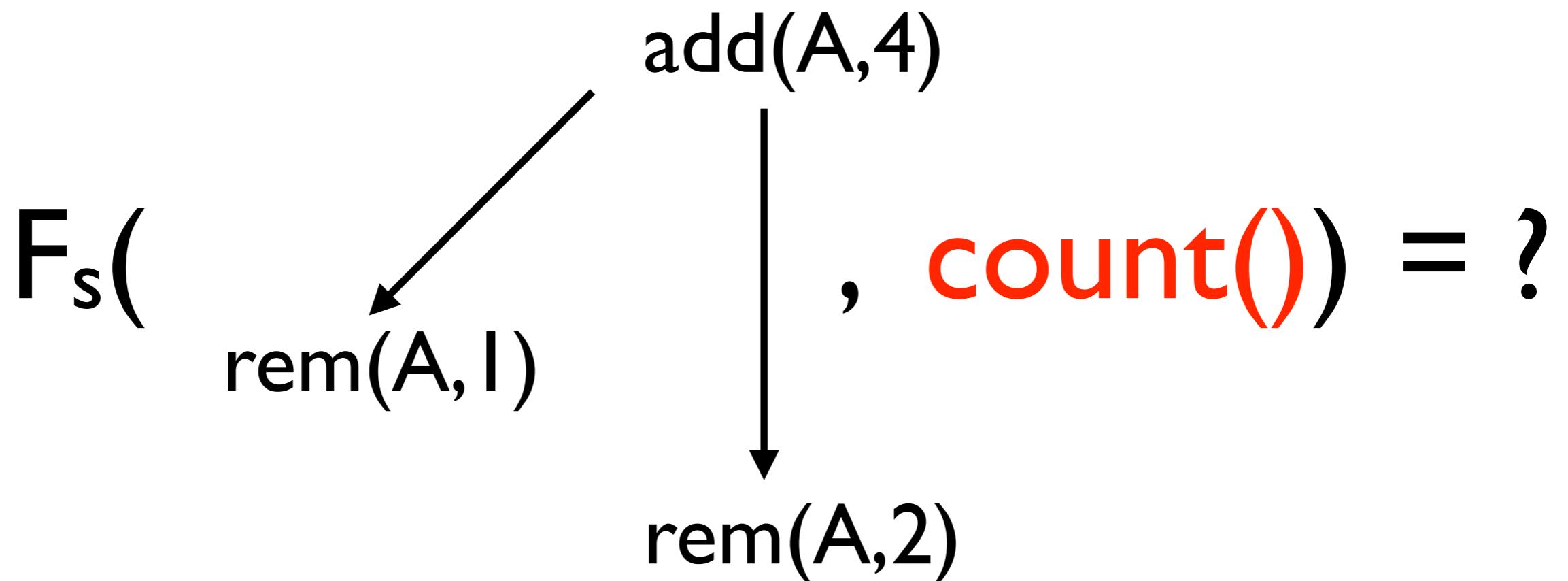
```



```

val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

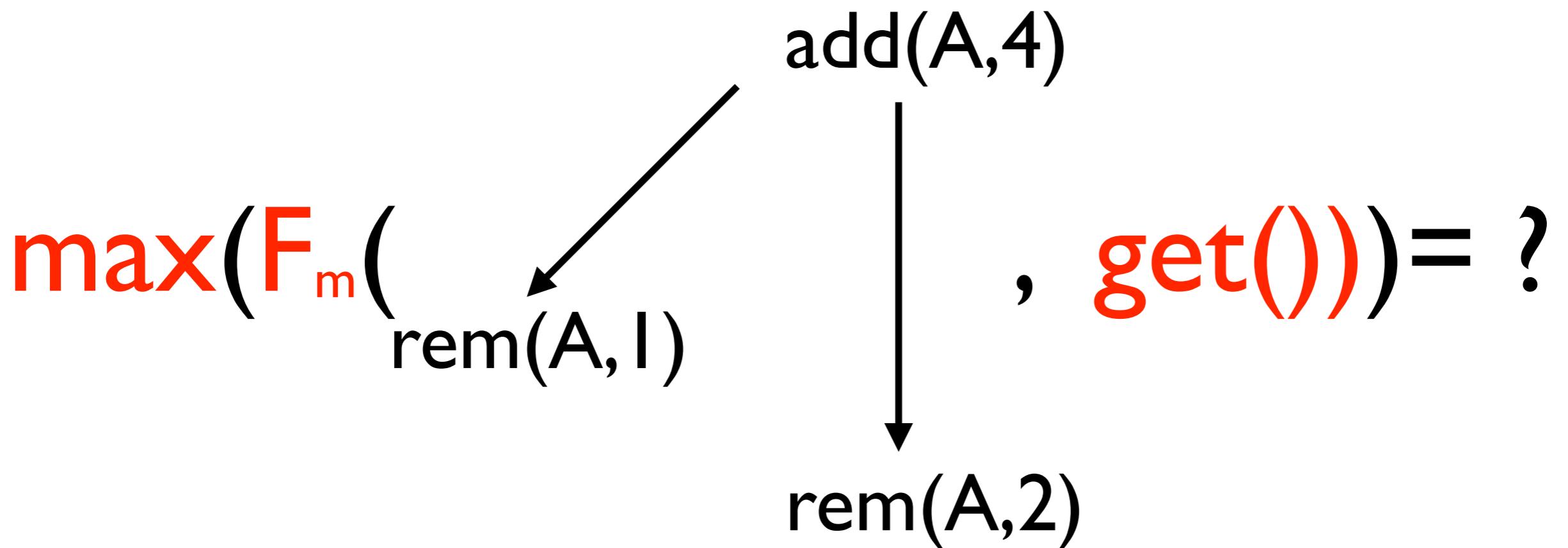
```



```

val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

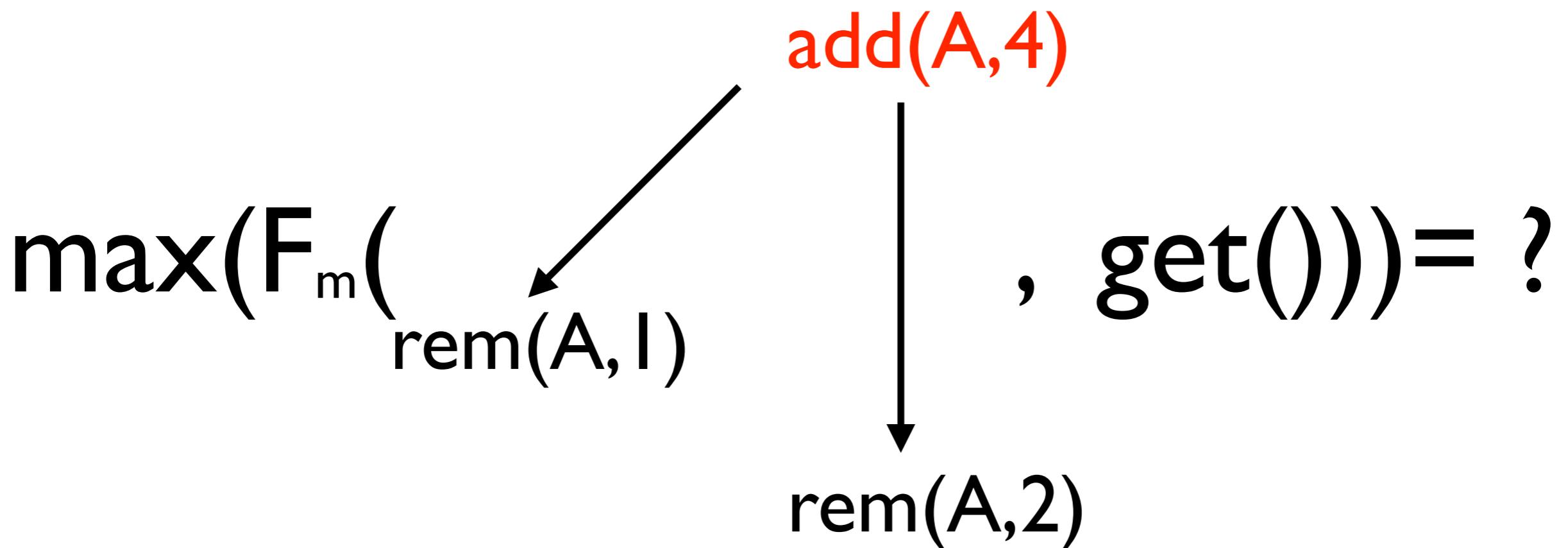
```



```

val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

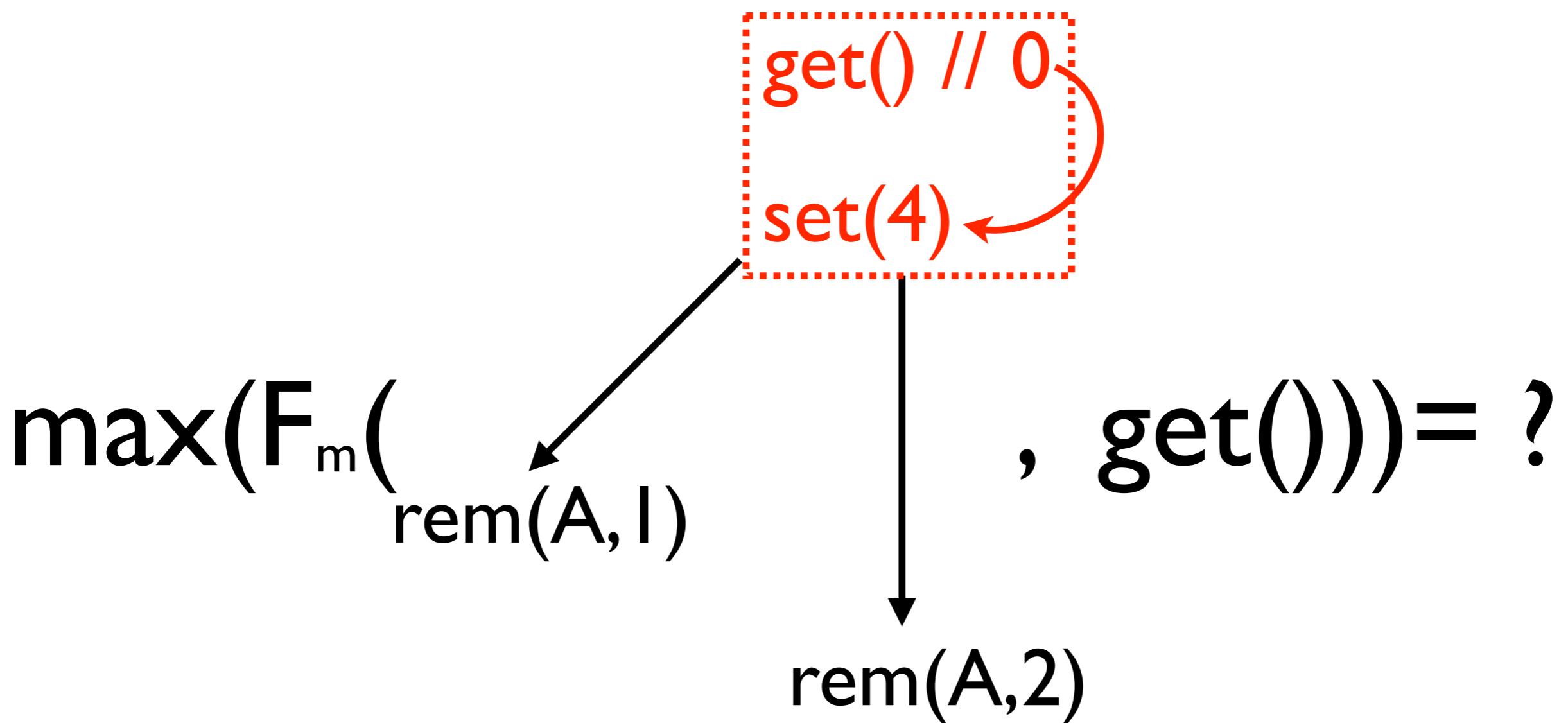
```



```

val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

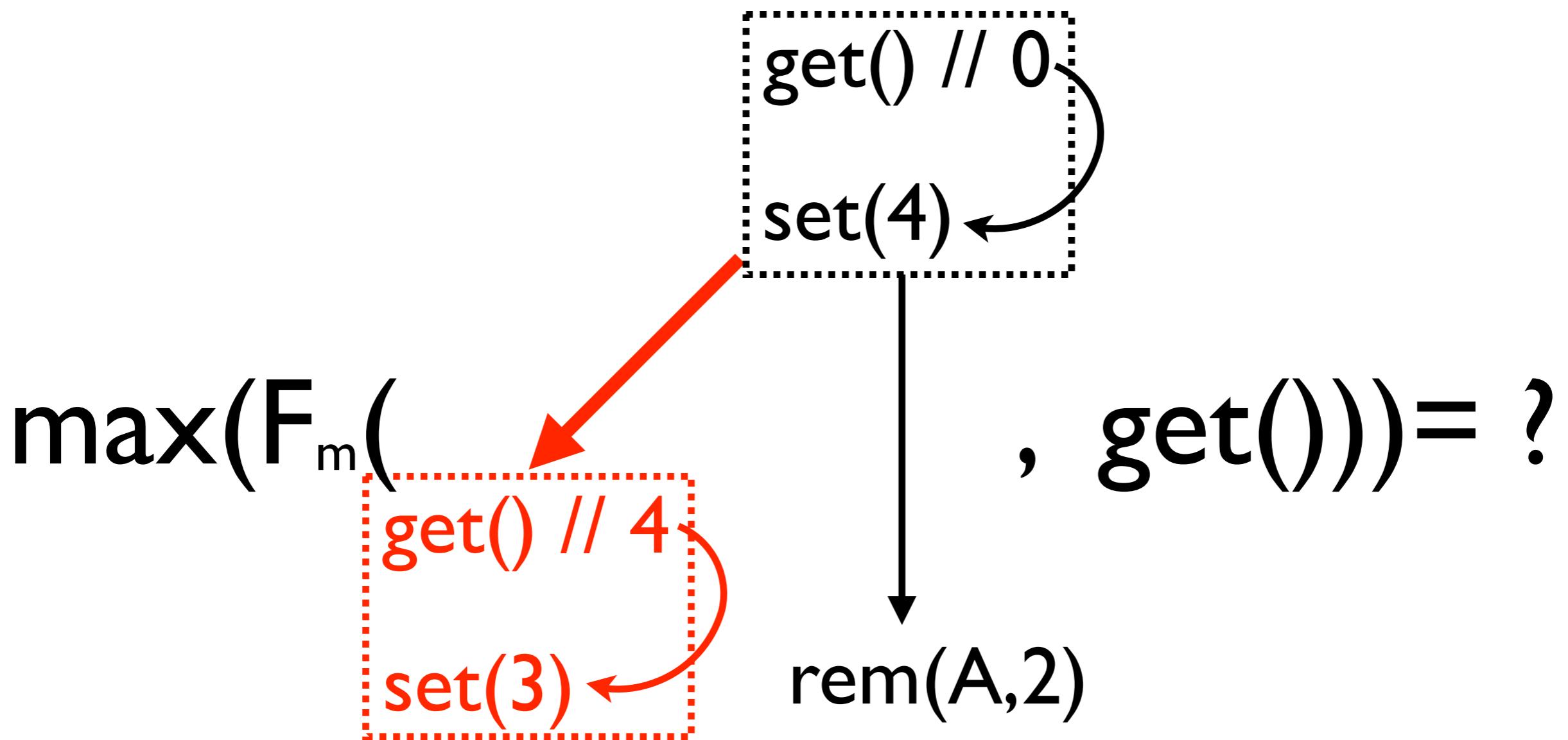
```



```

val map = new array[MVRegister[Nat]](N)
.
.
.
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

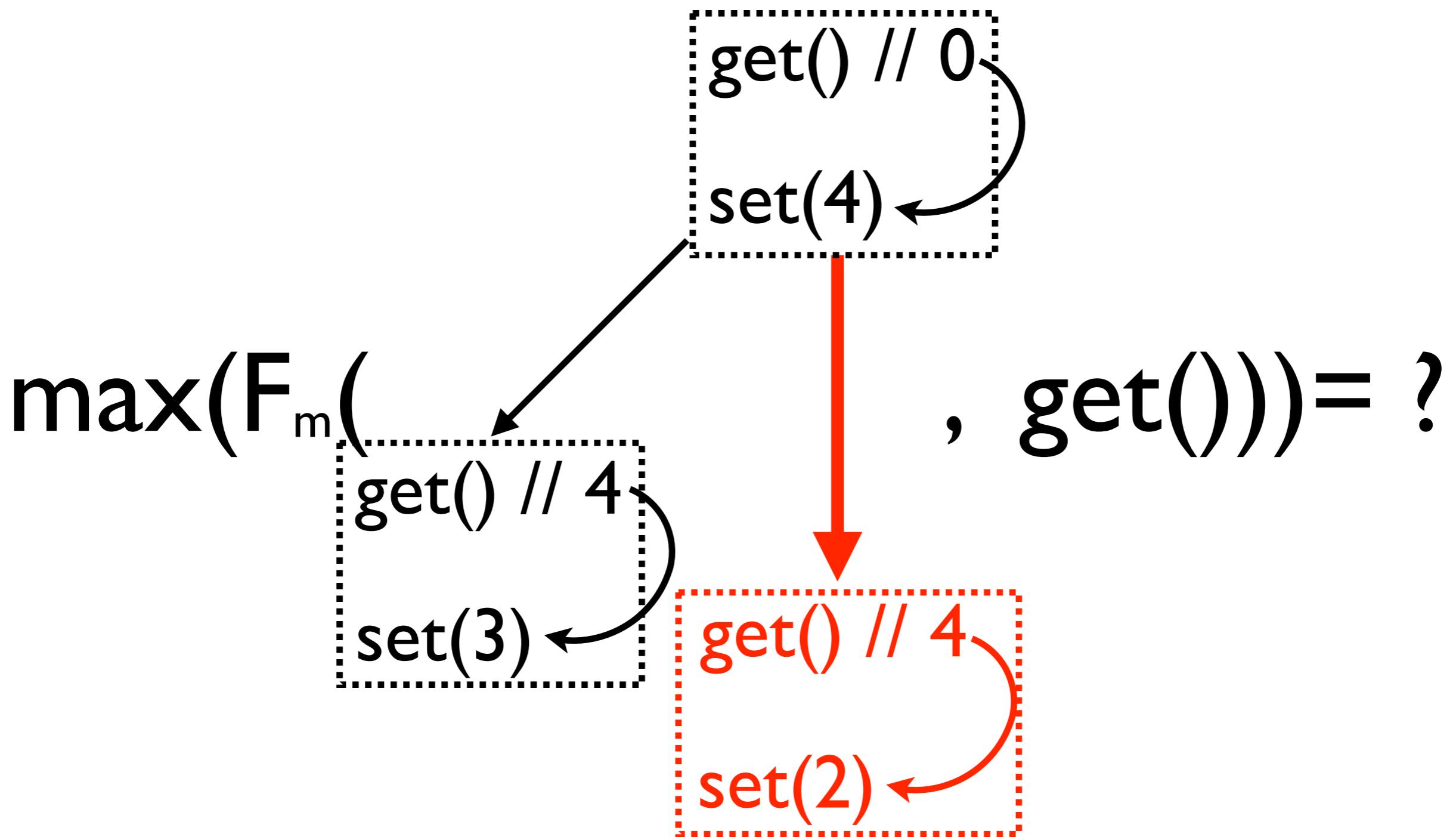
```



```

val map = new array[MVRegister[Nat]](N)
.
.
.
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

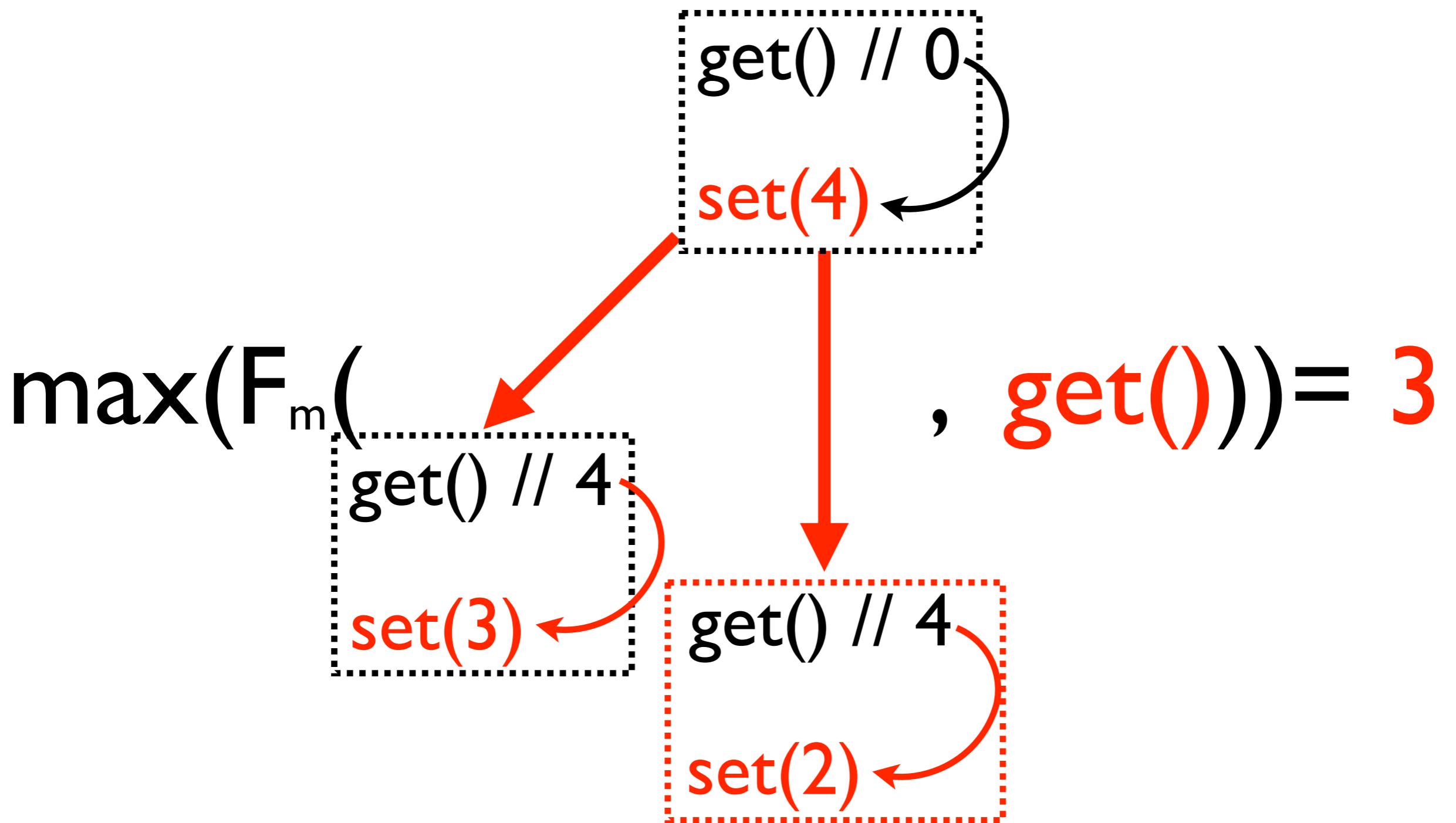
```



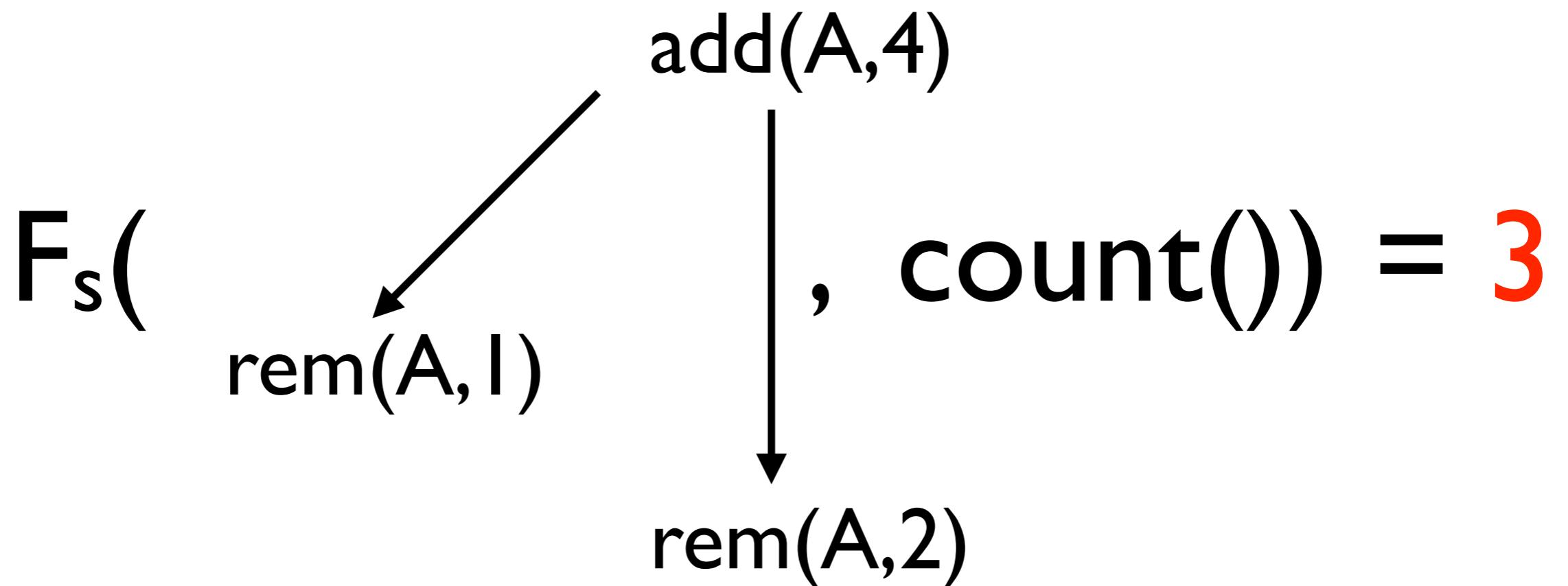
```

val map = new array[MVRegister[Nat]](N)
.
.
.
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

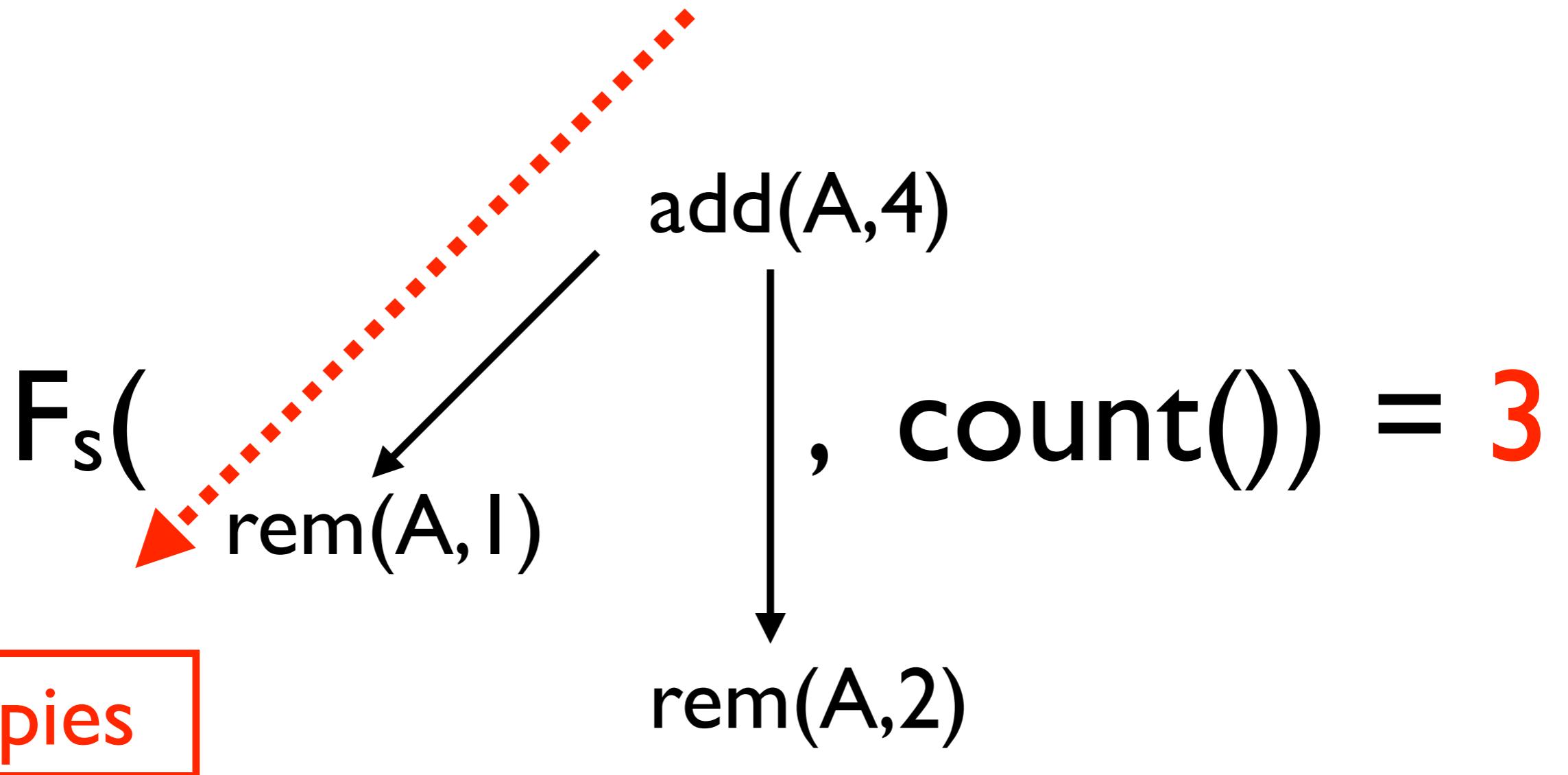
```



```
val map = new array[MVRegister[Nat]](N)
. . .
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }
```



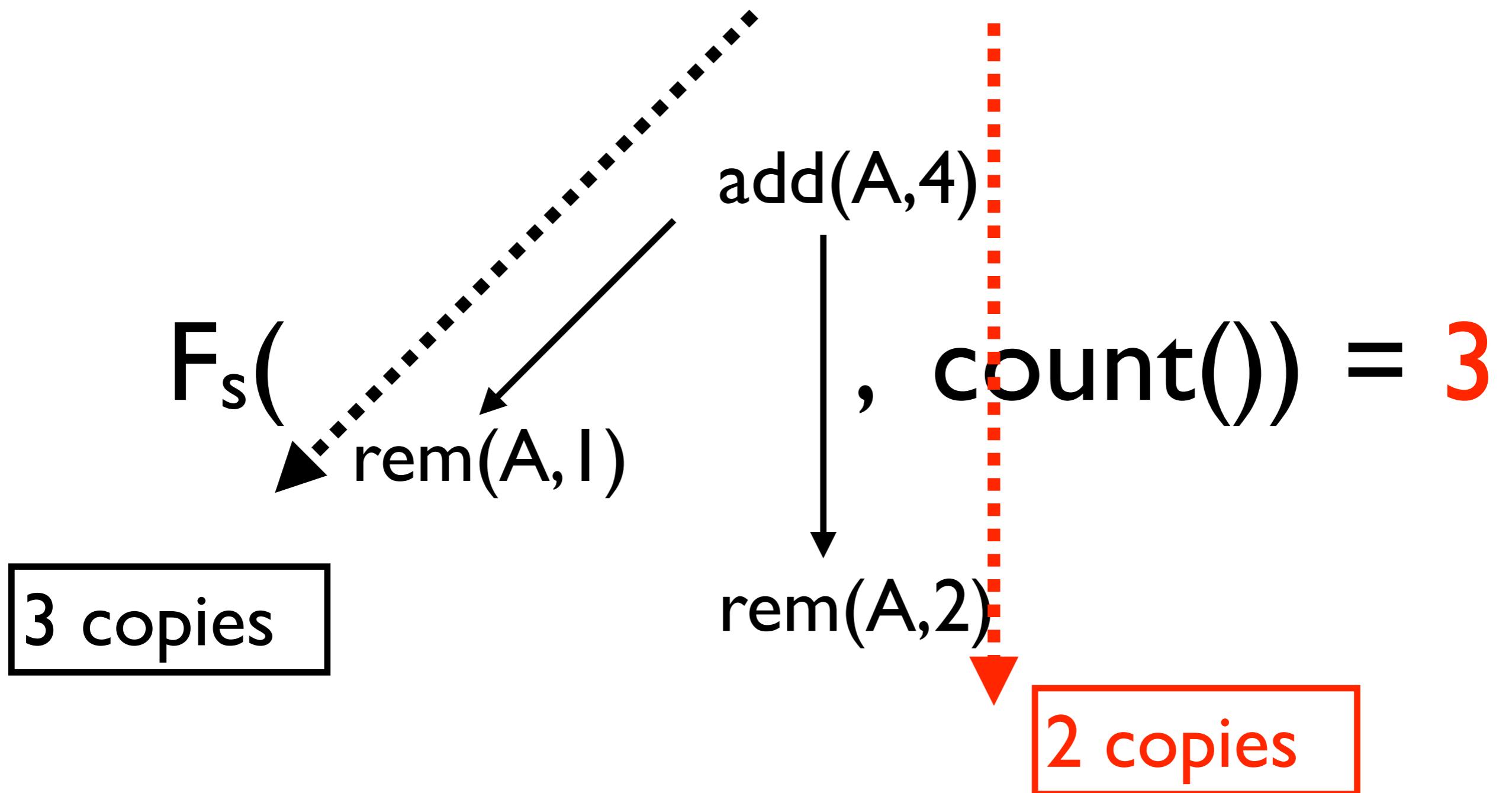
```
val map = new array[MVRegister[Nat]](N)
.
.
.
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }
```



```

val map = new array[MVRegister[Nat]](N)
.
.
.
def add(b,n)={ val m = max(map(b).get()); map(b).set(m+n) }
def rem(b,n)={ val m = max(map(b).get()); map(b).set(max(m-n,0)) }
def count(b)={ val m = max(map(b).get()); return m }

```



Semantics

- Developed F-based semantics of distributed programs with weak transactions.
- It is equivalent to a more standard semantics.
- The details are in our tech report. Email me if you are interested.

Take-home message

Developing a good programming model for weakly consistent distributed stores is a big challenge.