

# Composite Replicated Data Types

Alexey Gotsman

IMDEA Software Institute, Madrid, Spain

*Joint work with Hongseok Yang (Oxford)*

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

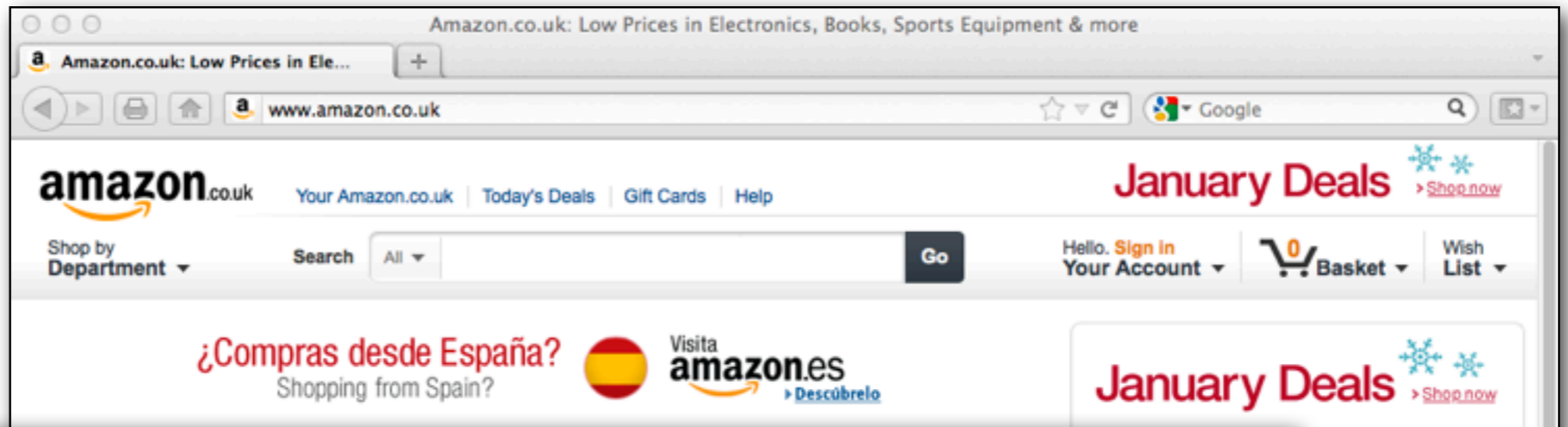
Birthday

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



Google

Google Search I'm Feeling Lucky

facebook

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.

Email or Phone  
Keep me logged in

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



Load-balancing, fault-tolerance

# Replicas on mobile devices



Offline use

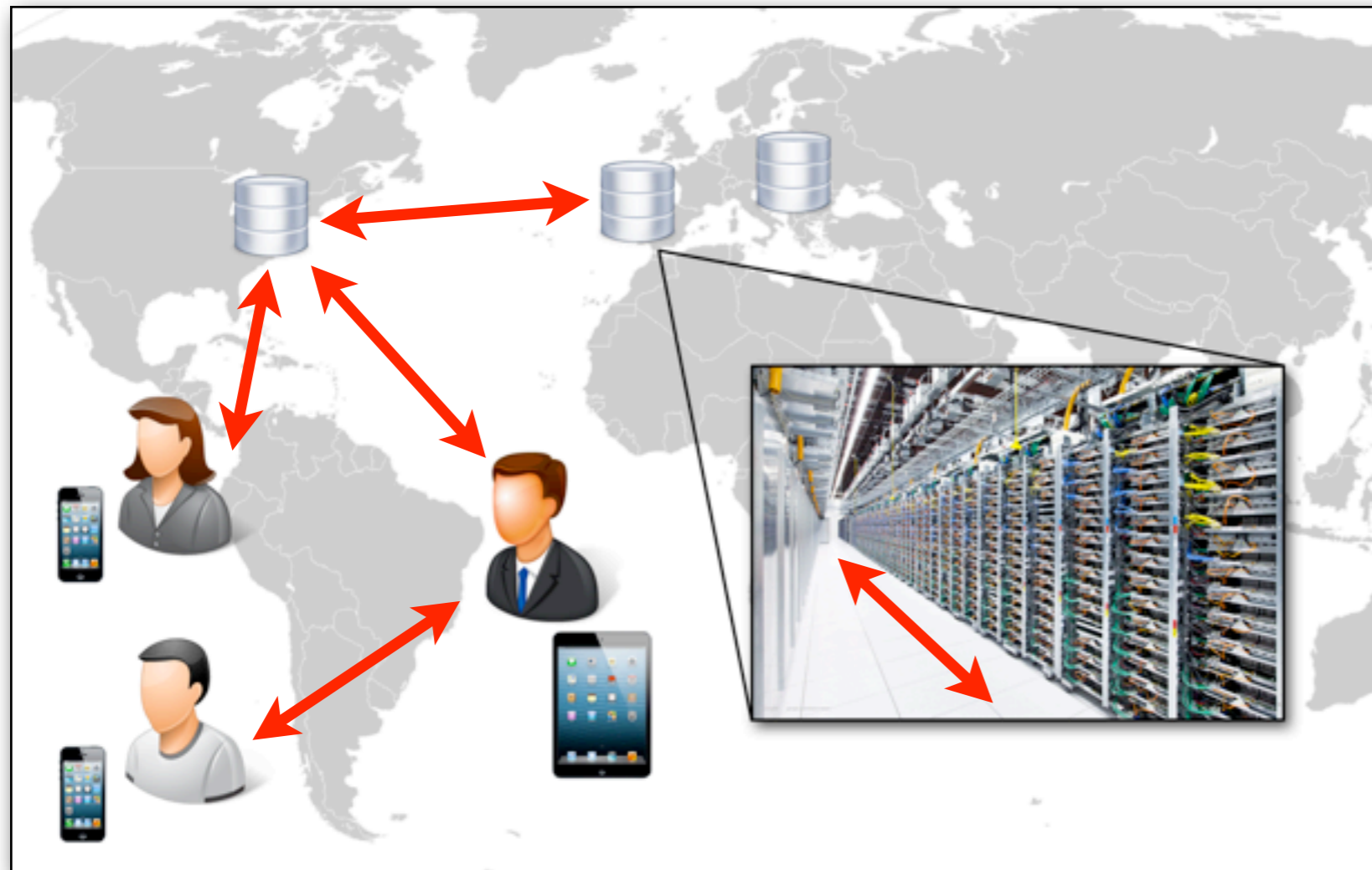




$\approx$



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



≈



- **Serialisability**: the system behaves like a serial processor of transactions on a centralised database
- Requires **synchronisation**: contact other replicas when processing a request
- Expensive or impossible



# Eventually consistent databases



```
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```



No synchronisation: update your replica now,  
propagate to others later

# Eventually consistent databases



```
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```

```
friends[Alice].get: {Bob}  
friends[Bob].get: {Alice}
```

- All updates by the transaction delivered together
- Can preserve invariants, but isn't serialisability:  
asynchronous communication leads to **anomalies**



```
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```



```
wall[Bob].add(post)
```





```
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```



```
wall[Bob].add(post)
```



```
wall[Bob].get: post
```



```
friends[Bob].get: ⊘
```



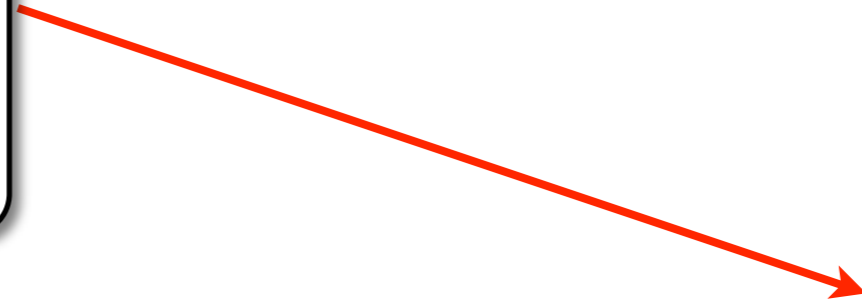
```
friends[Alice].add(Bob)
friends[Bob].add(Alice)
```



```
wall[Bob].add(post)
```



```
wall[Bob].get: post
```



```
friends[Bob].get: ∅
```



This talk - causal consistency model: causality is preserved



friends[Alice].add(Bob)  
friends[Bob].add(Alice)

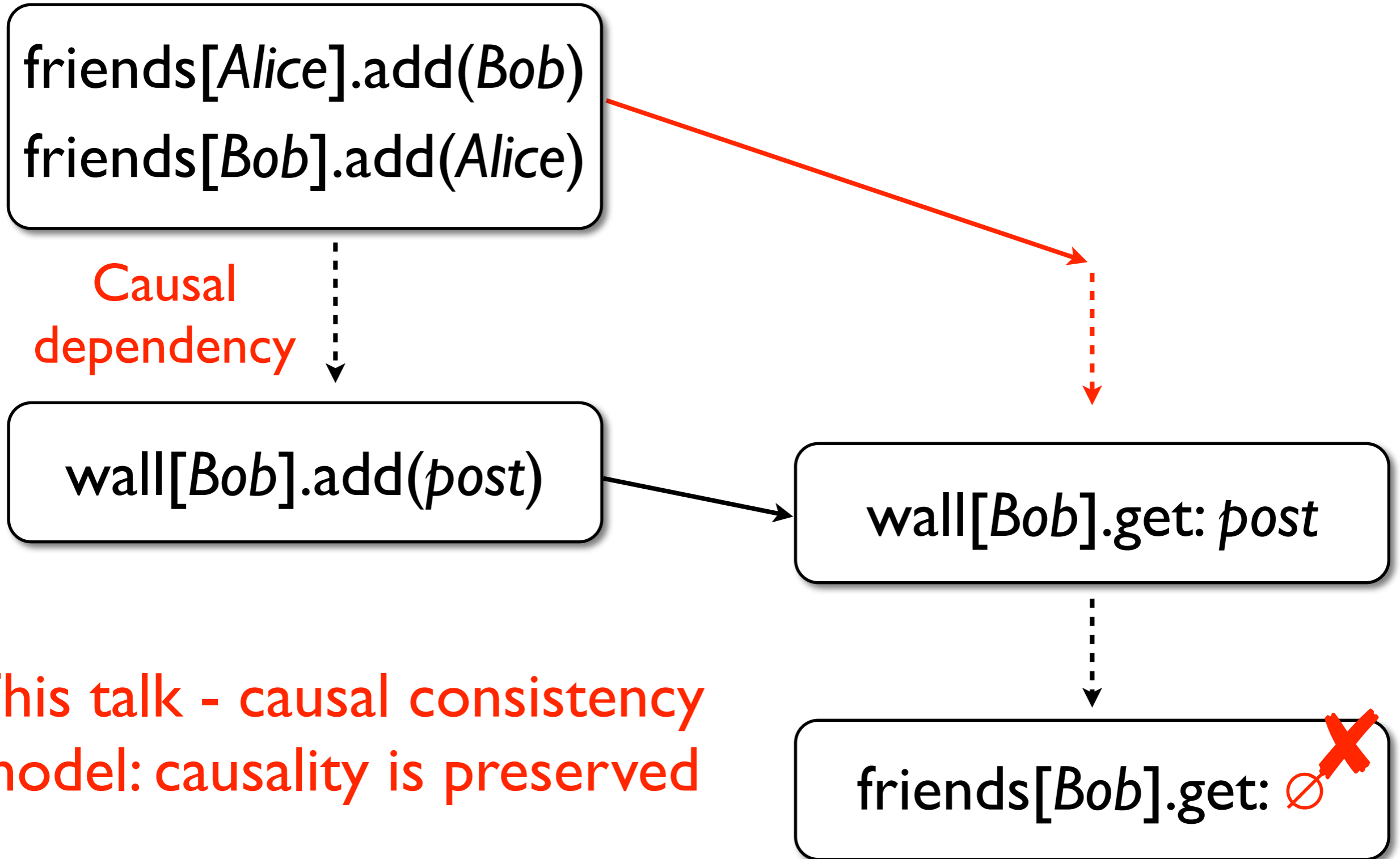
Causal  
dependency

wall[Bob].add(post)

wall[Bob].get: post

friends[Bob].get: ~~∅~~

This talk - causal consistency  
model: causality is preserved







```
requests[Alice].  
add(Bob)
```

Lack of synchronisation  
leads to conflicting  
updates



```
requests[Alice].  
add(Bob)
```



```
requests[Alice].  
add(Bob)
```

Lack of synchronisation  
leads to conflicting  
updates



```
requests[Alice].  
add(Bob)
```



```
requests[Alice].  
add(Bob)
```

accept



Lack of synchronisation  
leads to conflicting  
updates

```
requests[Alice].remove(Bob)  
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```





```
requests[Alice].  
add(Bob)
```



```
requests[Alice].  
add(Bob)
```

accept



Should remove cancel  
the concurrent add?

```
requests[Alice].remove(Bob)  
friends[Alice].add(Bob)  
friends[Bob].add(Alice)
```



`requests[Alice].  
add(Bob)`

Conflict-resolution policies:

**Remove wins:** `requests[Alice] = ∅`

**Add wins:** `requests[Alice] = {Bob}`

**accept**



Should remove cancel  
the concurrent add?

`requests[Alice].remove(Bob)`  
`friends[Alice].add(Bob)`  
`friends[Bob].add(Alice)`



`requests[Alice].  
add(Bob)`

Conflict-resolution policies:

**Remove wins:** `requests[Alice] = ∅`

**Add wins:** `requests[Alice] = {Bob}`

**accept**



Should remove cancel  
the concurrent add?

`requests[Alice].remove(Bob)`  
`friends[Alice].add(Bob)`  
`friends[Bob].add(Alice)`

## Conflict-resolution policies:

**Remove wins:**  $\text{requests}[Alice] = \emptyset$

**Add wins:**  $\text{requests}[Alice] = \{Bob\}$

- Encapsulated in replicated data types (aka CRDTs) [Shapiro+ 2011]
- Object  $\rightarrow$  Type  $\rightarrow$  Conflict-resolution policy
- Remove-wins set, add-wins set, counters, registers...

# The brave new world of eventual consistency

Features to maintain correctness despite weak consistency:

- Consistency models restricting anomalies:  
causal consistency
- Programming concepts:  
replicated data types, transactions

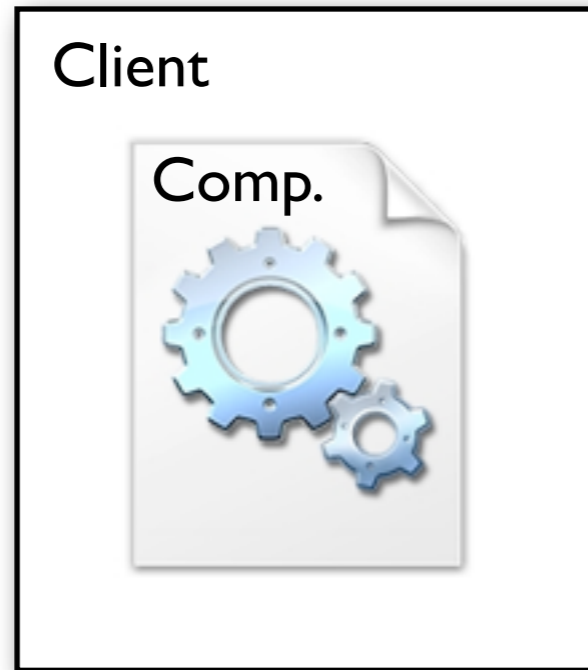
Subtle semantics → programming is difficult

# Goal

- Want to help programmers: specification, verification, programming models
- Need a better theoretical understanding of programming on eventual consistency
- Groundwork: formal semantics specification [POPL'14]
  - ▶ Executions represented using structures of events and partial orders
  - ▶ Replicated data type specifications

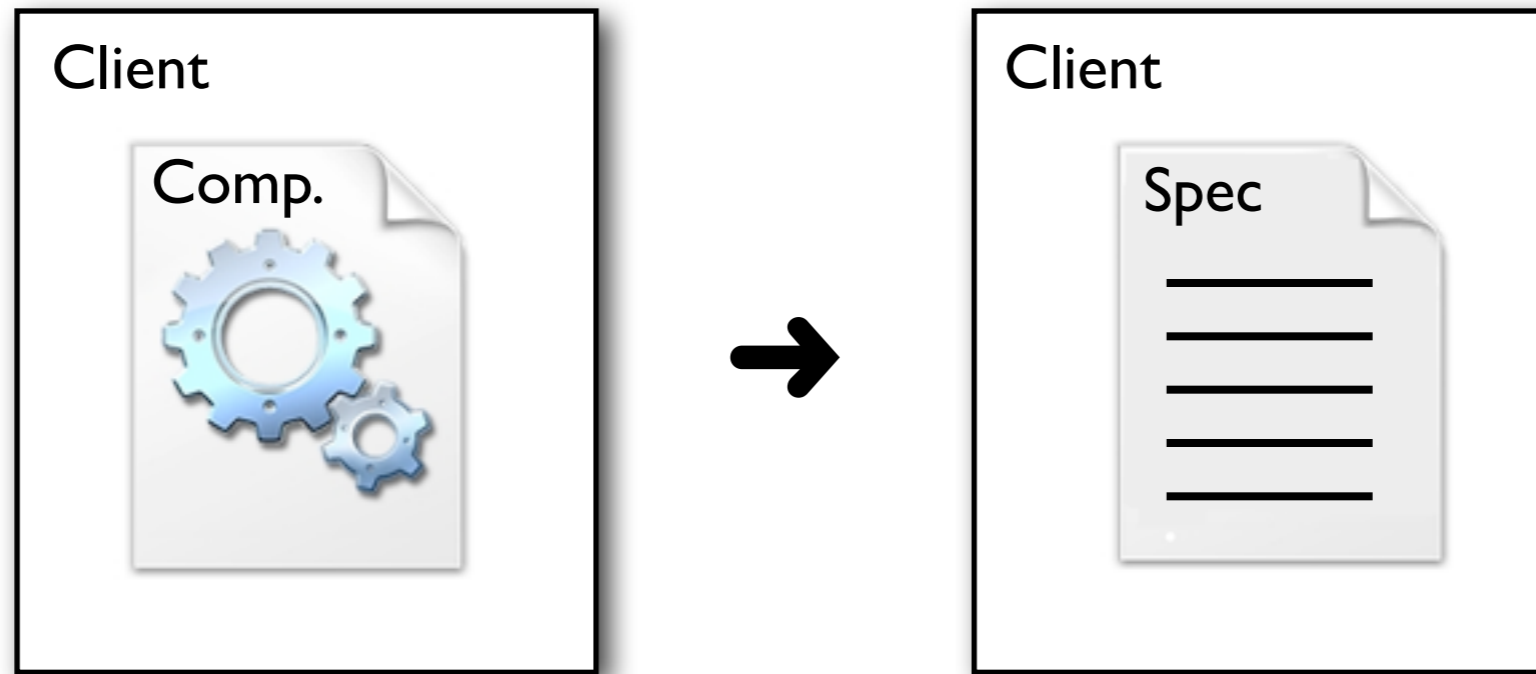


# This talk: modularity



- Programmers encapsulate functionality into components
- Use a component without knowing its implementation

# This talk: modularity



- Programmers encapsulate functionality into components
- Use a component without knowing its implementation
- Want to abstract from component internals when reasoning about its client → replace it by a spec

# Composite replicated data types

≈ ADT, module

```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].  
        remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```

# Composite replicated data types

≈ ADT, module

```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].  
        remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```

- Objects of replicated data types

# Composite replicated data types

≈ ADT, module

```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].  
        remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```

- Objects of replicated data types
- Information hiding: no direct access to the objects
- Methods implemented by causally consistent transactions

# Composite replicated data types

≈ ADT, module

```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].  
                                remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```

- Replicated data types ensure conflict resolution
- Transactions maintain integrity invariants
- Want to abstract from the implementation when reasoning about client programs



```
datatype SocialGraph {
  RemoveWinsSet friends[];
  RemoveWinsSet requests[];

  tx accept(Bob → Alice) {
    ...
    requests[Alice].remove(Bob);
    friends[Alice].add(Bob);
    friends[Bob].add(Alice);
    ...
  }
}
```

```
primitive datatype SocialGraph
operation accept(Bob → Alice)
```



Existing spec mechanism  
[POPL'14]

Object of a composite  
data type D



Object of a primitive data  
type with a spec  $\llbracket D \rrbracket$

**Denotation**  $\llbracket D \rrbracket$  of D  
= best spec

C

graph.accept

D

● friends.add

●

requests.  
remove

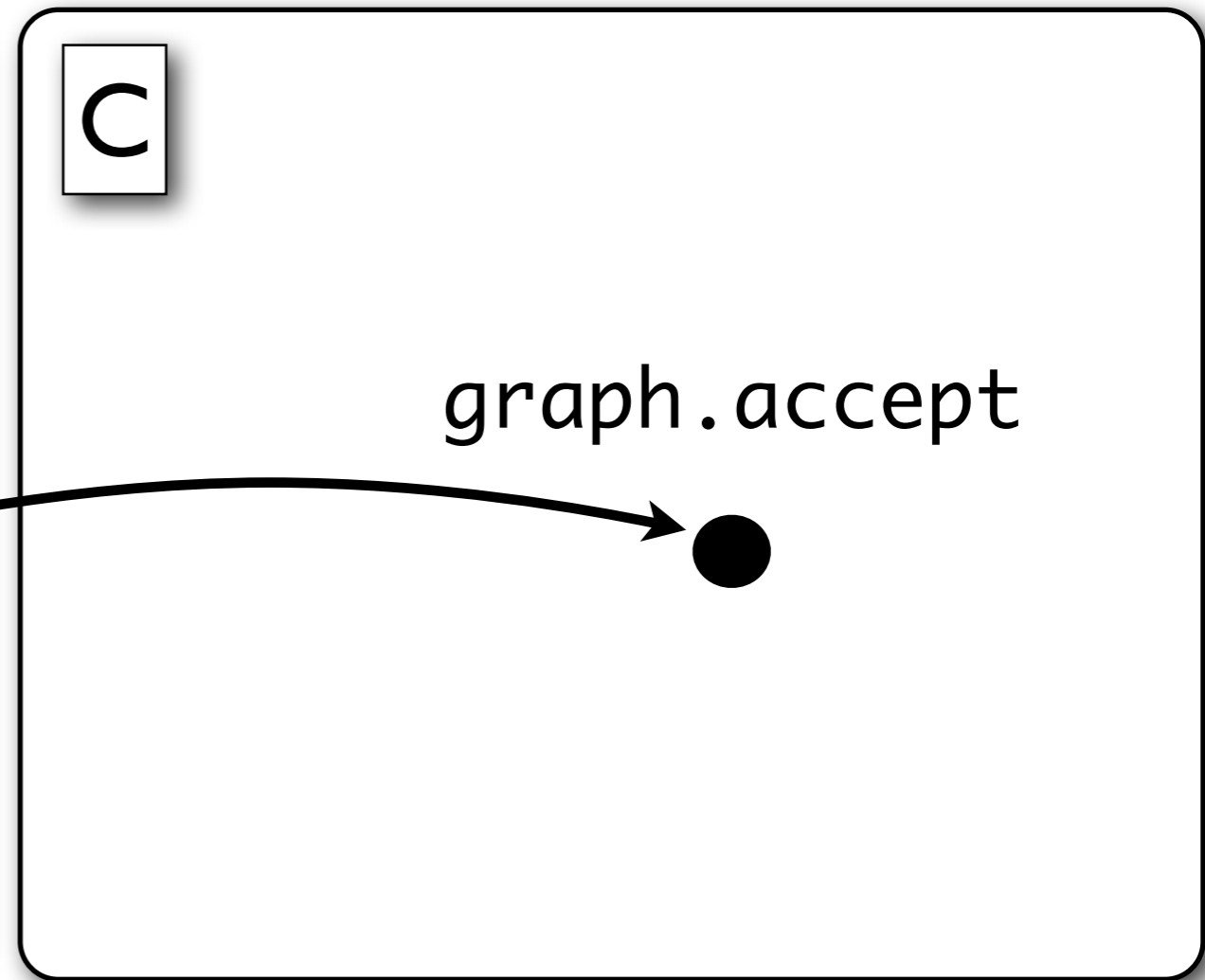
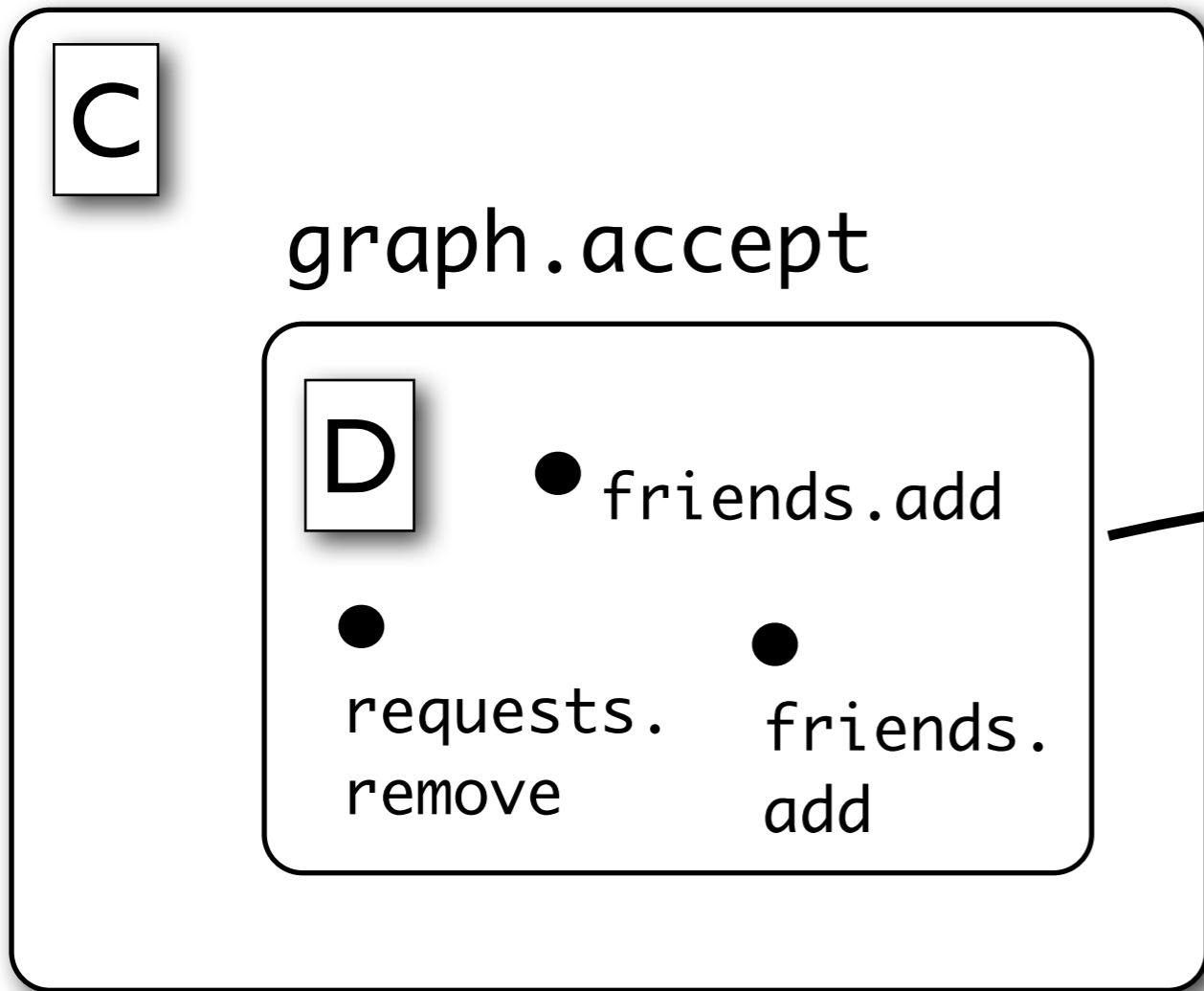
●

friends.  
add

**Fine-grain** semantics

[[C(D)]]

Method invocation →  
multiple events  
on internal objects



**Fine-grain** semantics

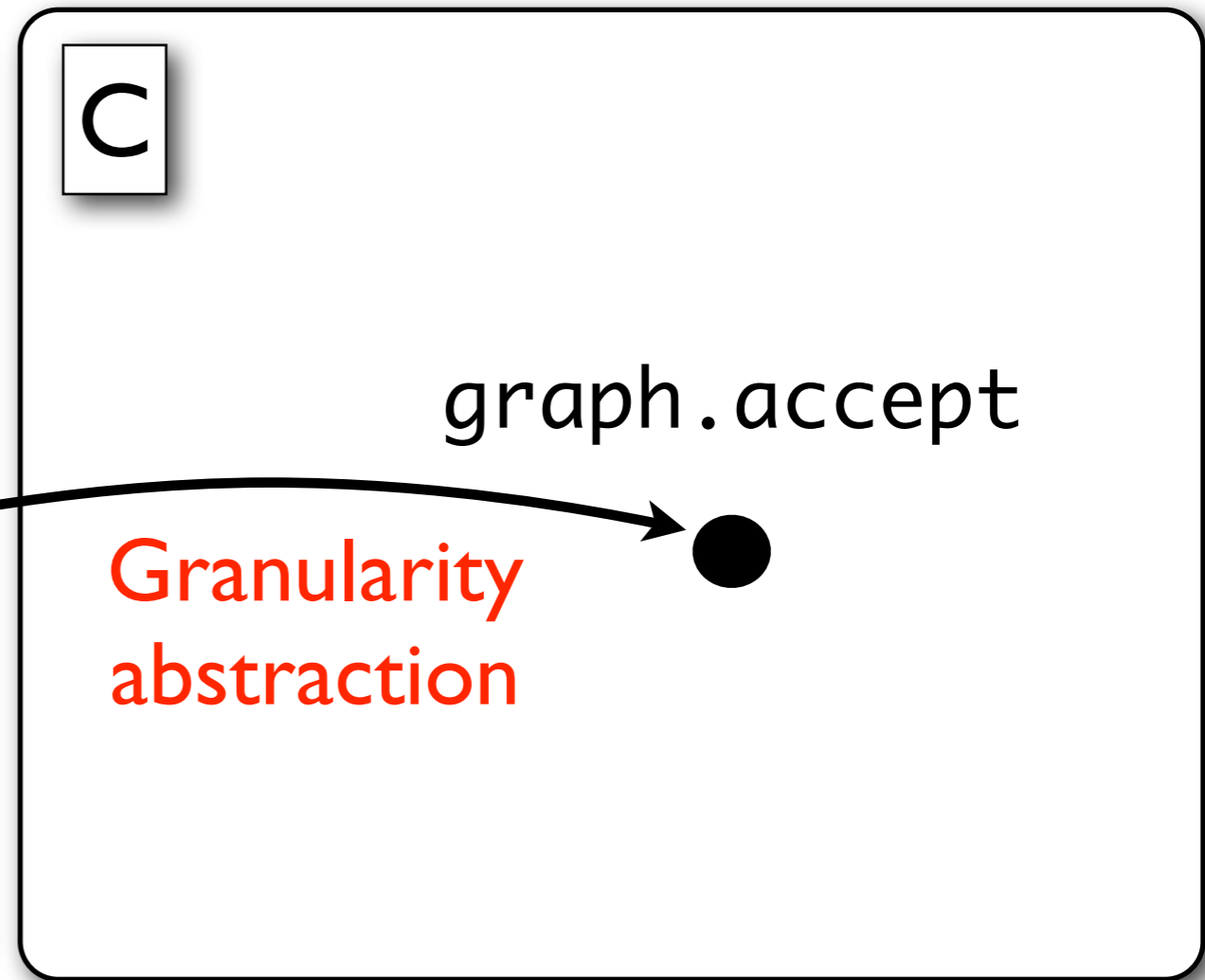
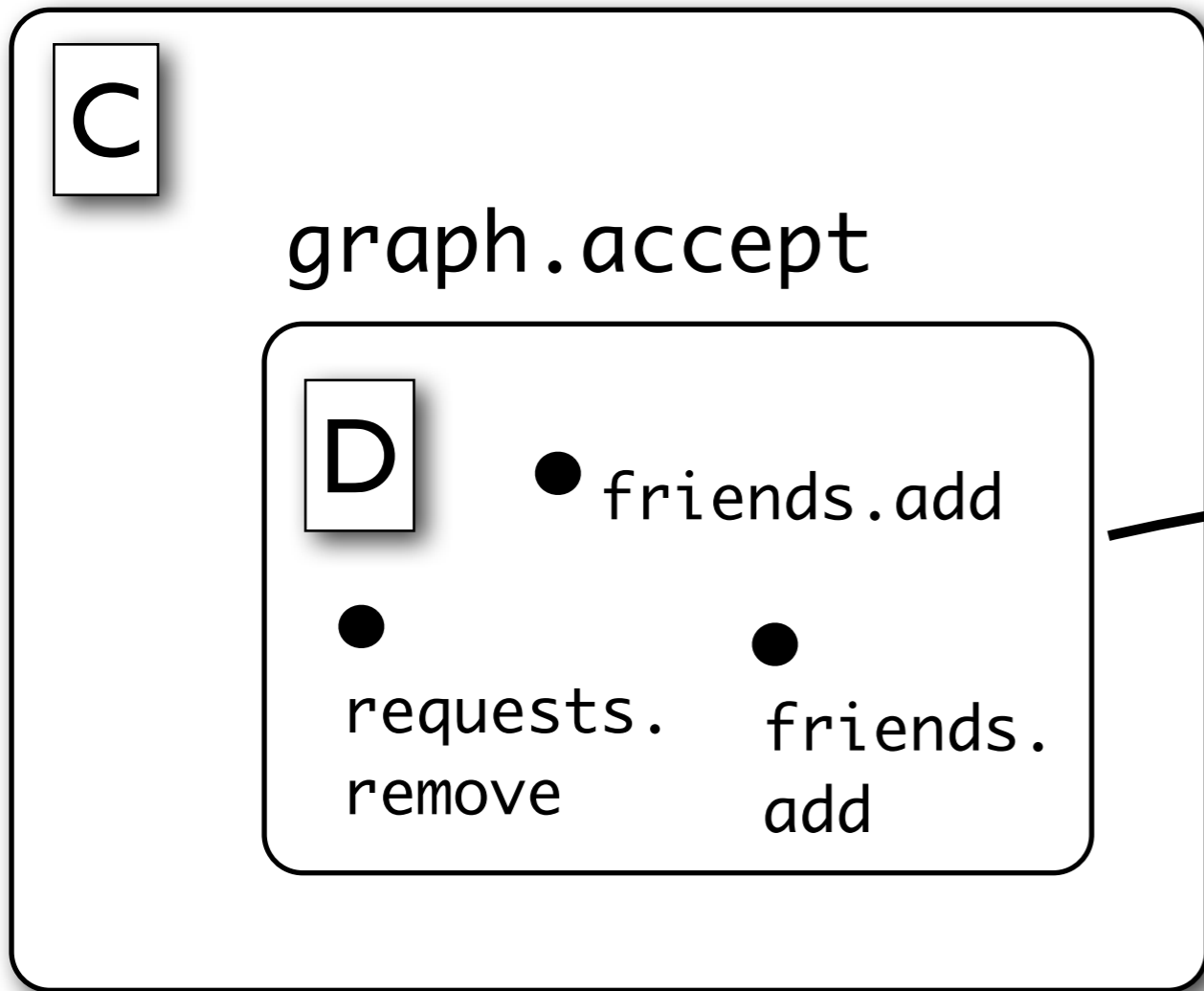
$\llbracket C(D) \rrbracket$

Method invocation  $\rightarrow$   
multiple events  
on internal objects

**Coarse-grain** denotational

semantics  $\llbracket C(\llbracket D \rrbracket) \rrbracket$

Method invocation  $\rightarrow$   
one event



**Fine-grain** semantics

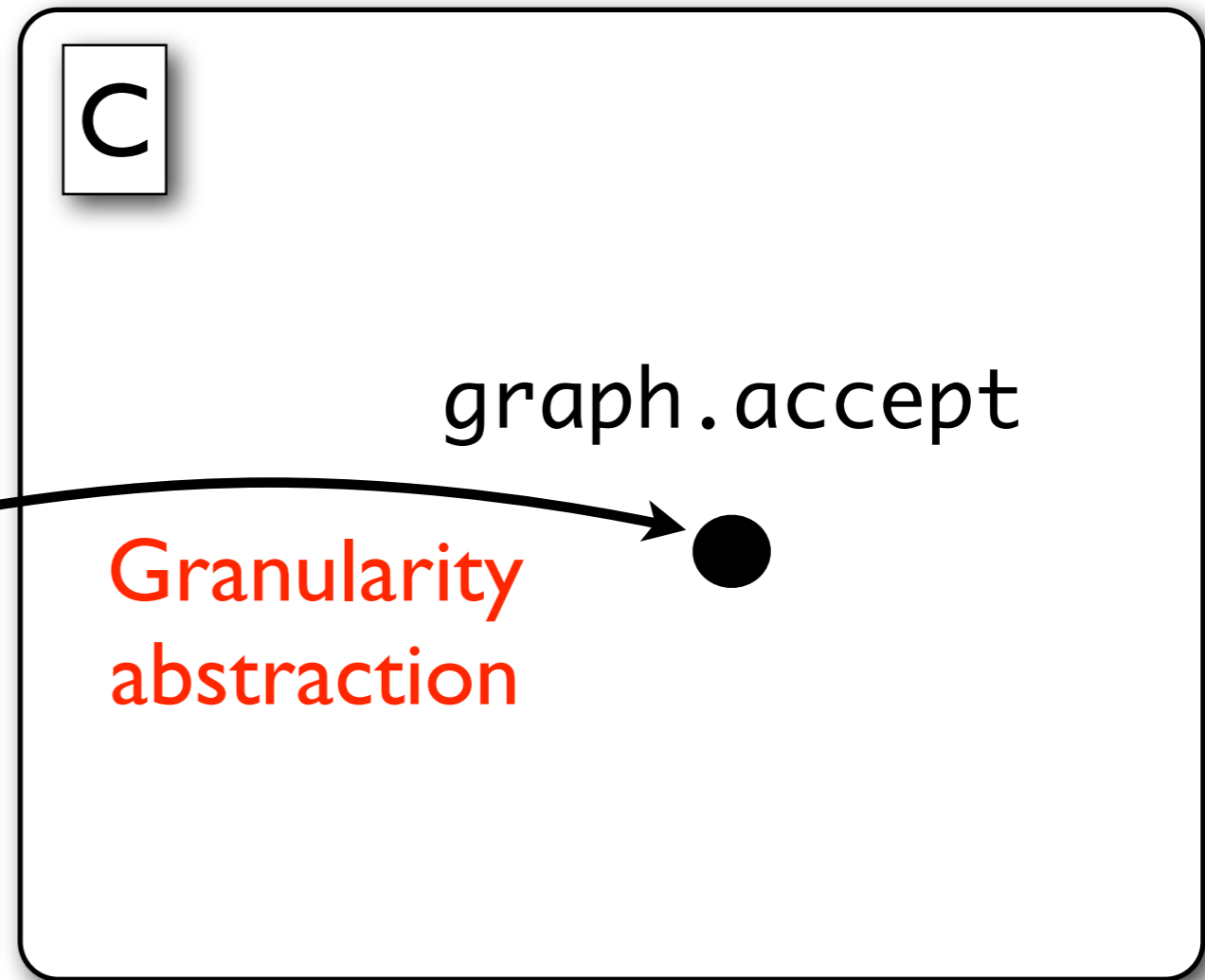
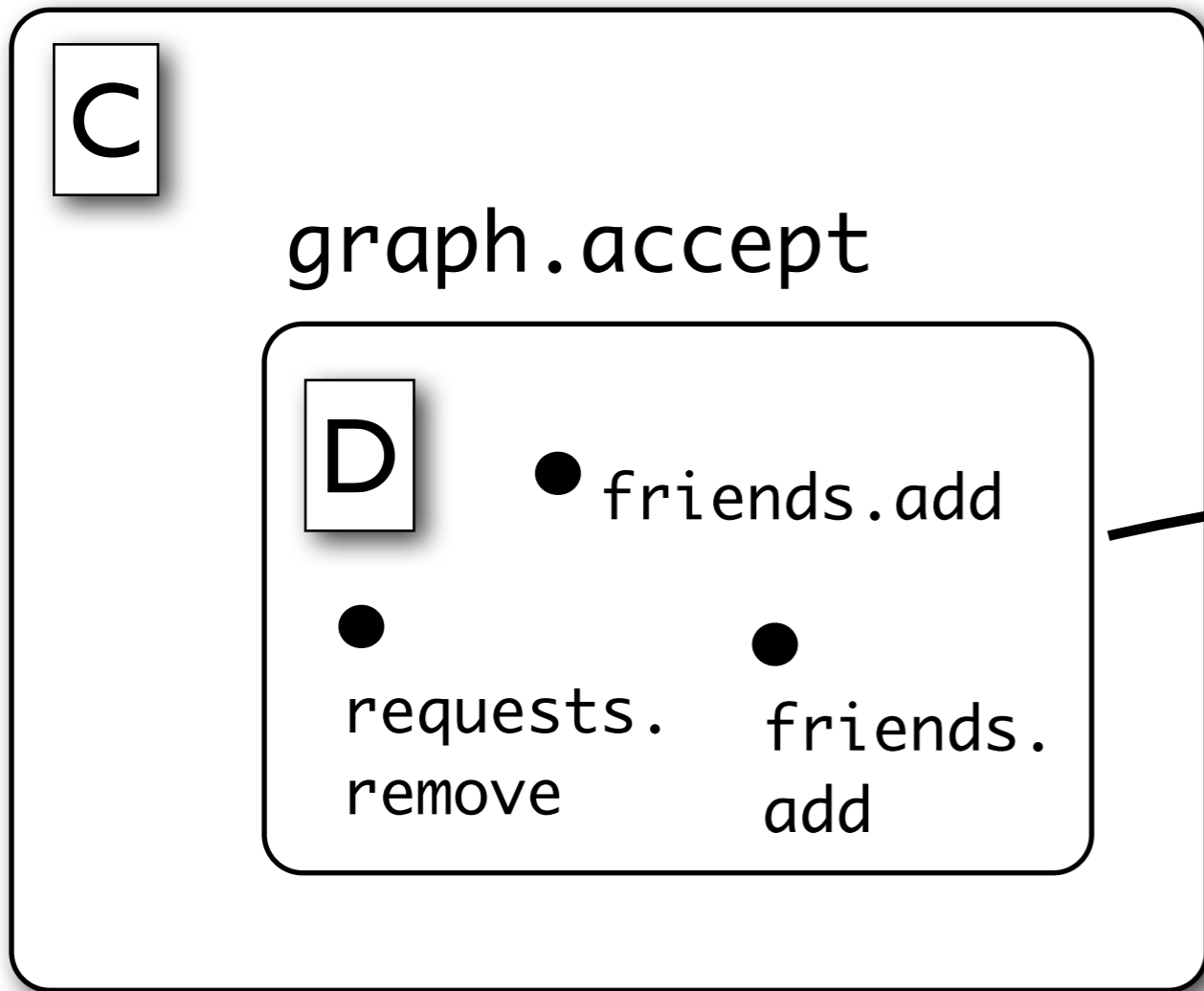
$\llbracket C(D) \rrbracket$

Method invocation  $\rightarrow$   
multiple events  
on internal objects

**Coarse-grain** denotational

semantics  $\llbracket C(\llbracket D \rrbracket) \rrbracket$

Method invocation  $\rightarrow$   
one event



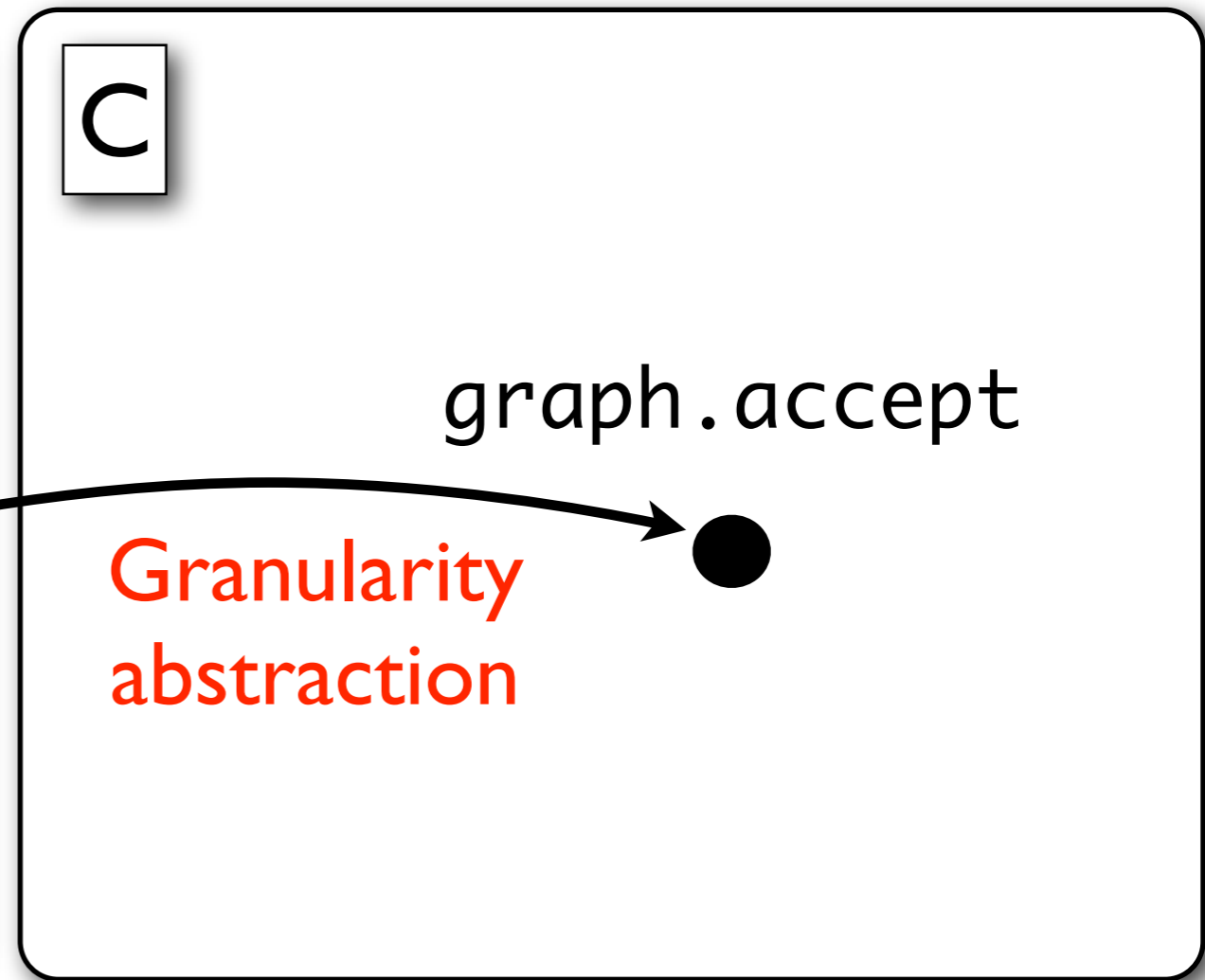
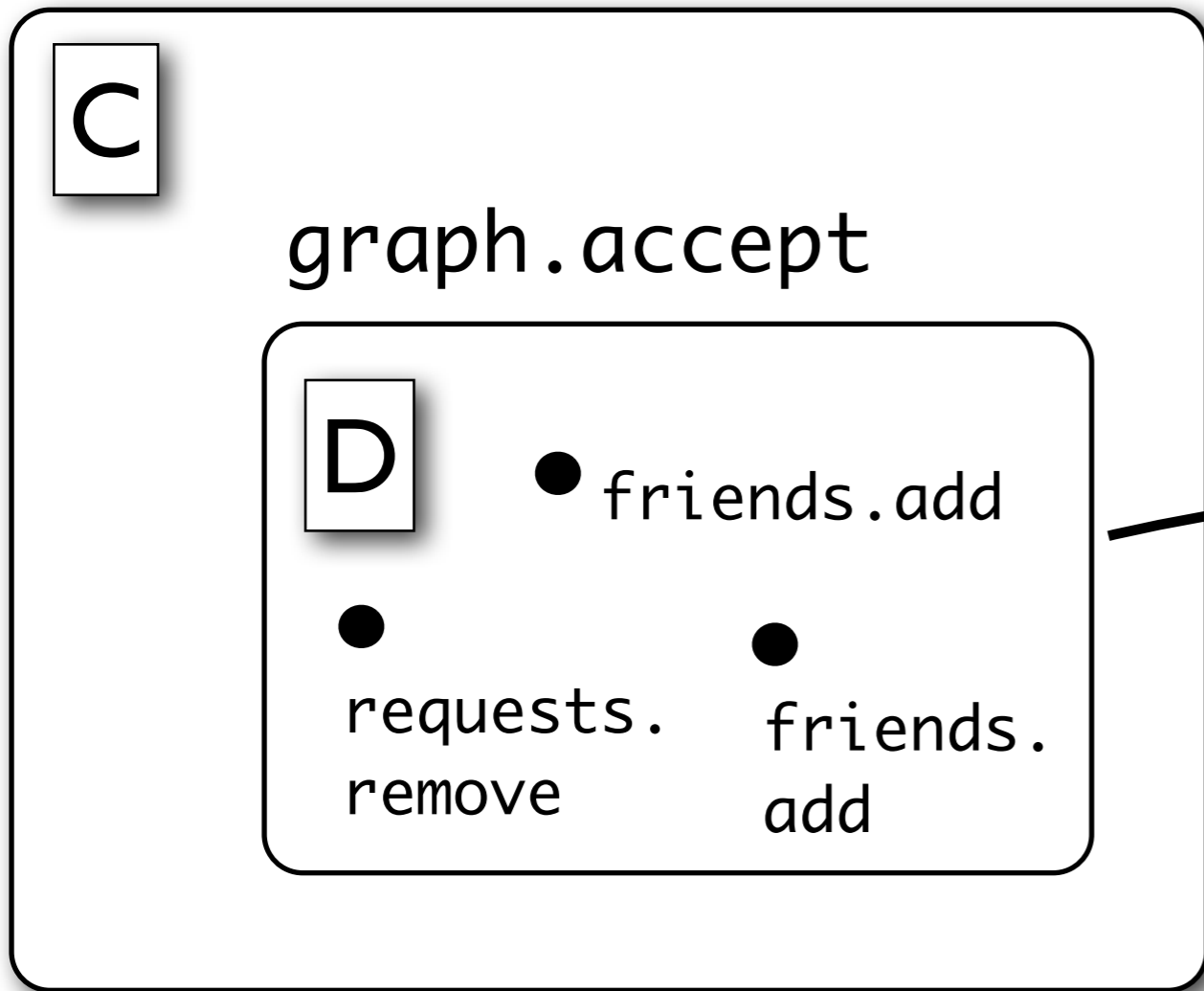
Fine-grain semantics

$\llbracket C(D) \rrbracket$

Coarse-grain denotational

semantics  $\llbracket C(\llbracket D \rrbracket) \rrbracket$

Method results determined using the spec  $\llbracket D \rrbracket$   
instead of the implementation  $D$



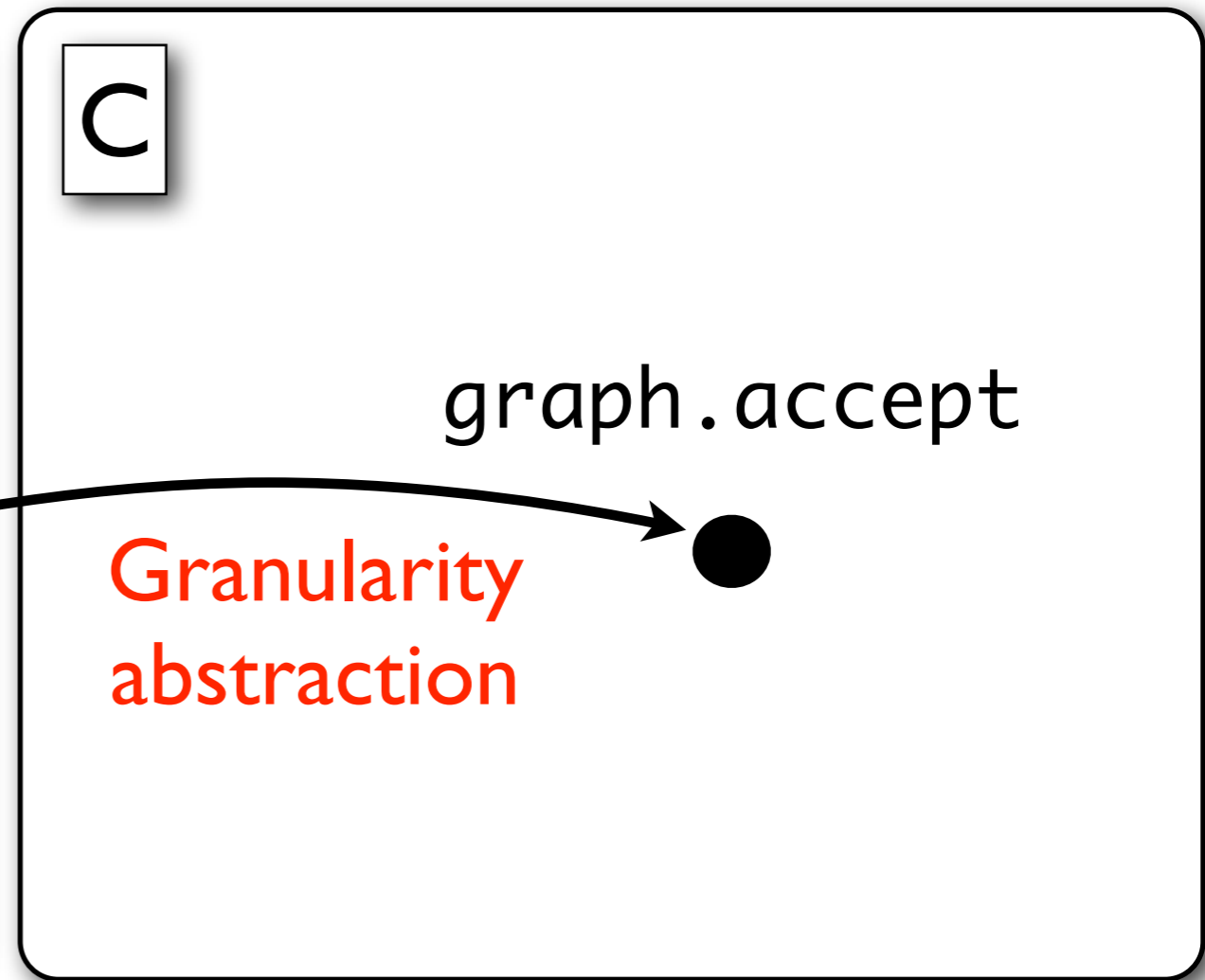
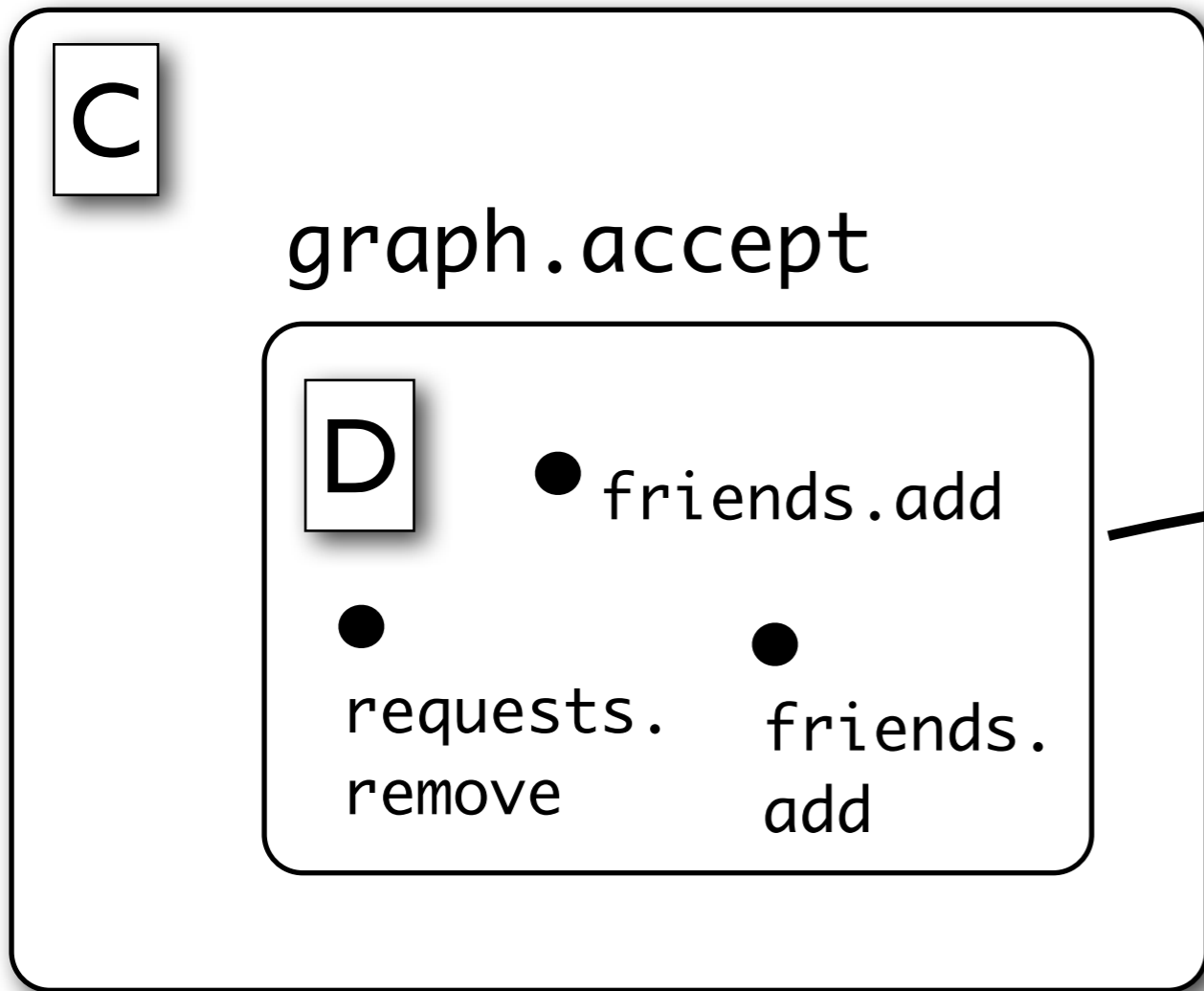
**Fine-grain semantics**  
 $\llbracket C(D) \rrbracket$

=

**Coarse-grain denotational semantics**  
 $\llbracket C(\llbracket D \rrbracket) \rrbracket$

Method results determined using the spec  $\llbracket D \rrbracket$   
 instead of the implementation  $D$





**Fine-grain** semantics

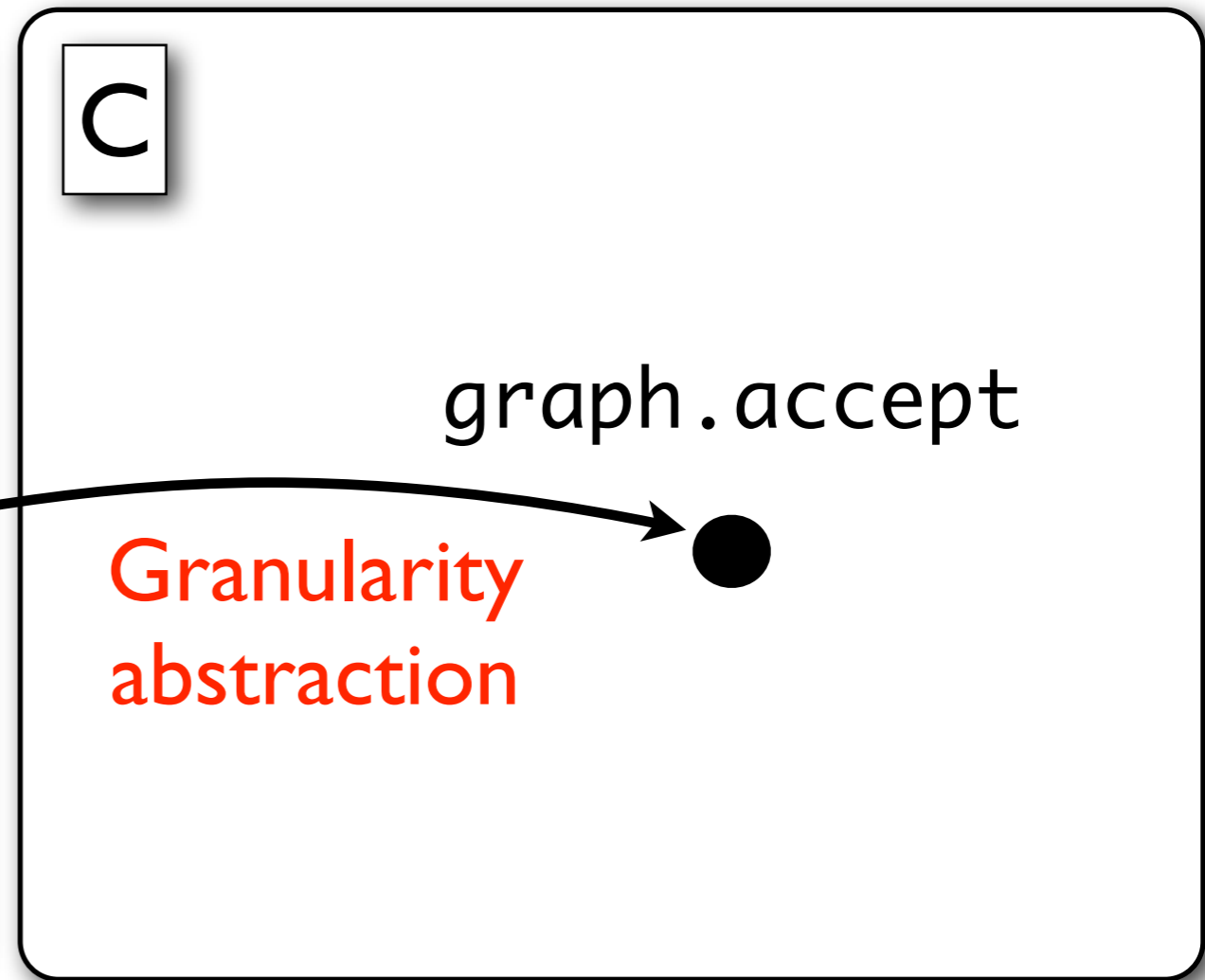
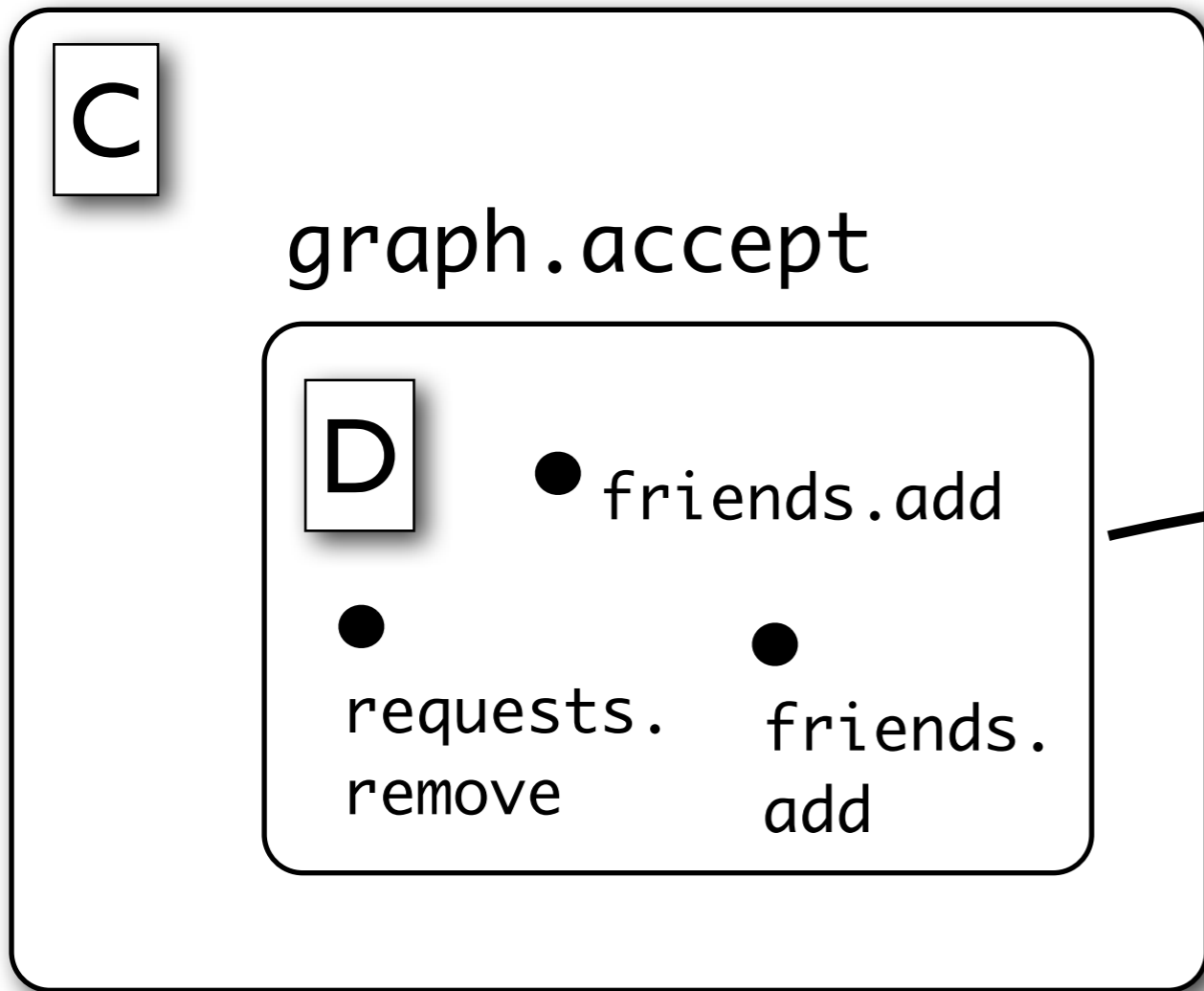
$\llbracket C(D) \rrbracket$

=

**Coarse-grain** denotational

semantics  $\langle\langle C(\llbracket D \rrbracket) \rangle\rangle$

I. Semantics of the bare database with primitive replicated data types



**Fine-grain** semantics

$\llbracket C(D) \rrbracket$

=

**Coarse-grain** denotational

semantics  $\langle\langle C(\llbracket D \rrbracket) \rangle\rangle$

1. Semantics of the bare database with primitive replicated data types

2. Definition of  $\llbracket D \rrbracket$

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

accept

requests[*Alice*].remove(*Bob*)  
friends[*Alice*].add(*Bob*)  
friends[*Bob*].add(*Alice*)

query *Alice*

friends[*Alice*].get: {*Bob*}  
requests[*Alice*].get: ∅

Database semantics = set of executions ~ graphs

Nodes: events describing operations performed

request *Bob* → *Alice*

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

requests[*Alice*].add(*Bob*)

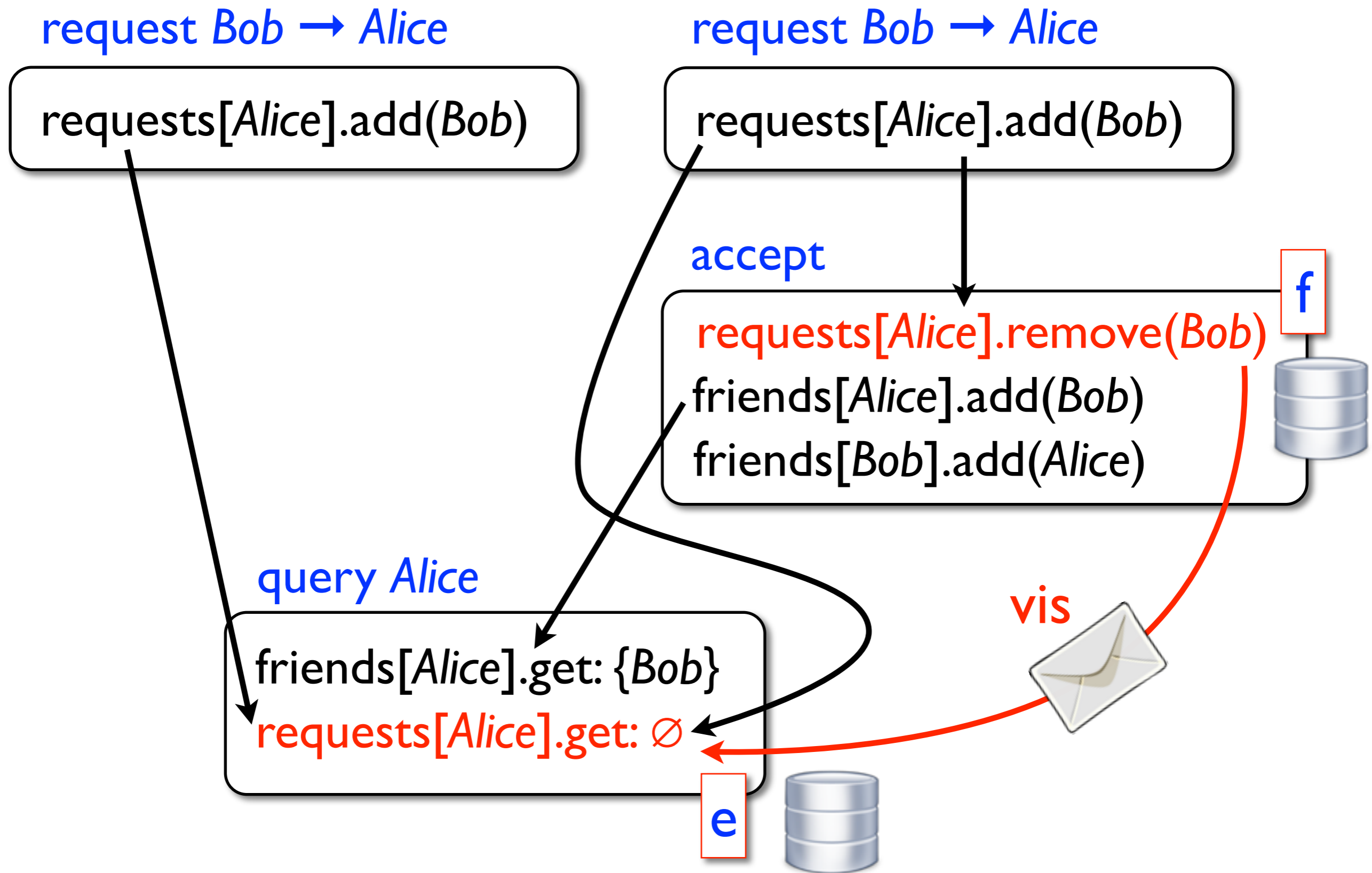
accept

requests[*Alice*].remove(*Bob*)  
friends[*Alice*].add(*Bob*)  
friends[*Bob*].add(*Alice*)

query *Alice*

friends[*Alice*].get: {*Bob*}  
requests[*Alice*].get: ∅

Events grouped into transactions



Visibility:  $(f, e) \in \text{vis} \rightarrow f$  delivered to the replica of  $e$

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

accept

requests[*Alice*].remove(*Bob*)  
friends[*Alice*].add(*Bob*)  
friends[*Bob*].add(*Alice*)

query *Alice*

friends[*Alice*].get: {*Bob*}  
requests[*Alice*].get: ∅

Visibility satisfies **consistency axioms**: e.g., transitivity

Restrict the anomalies allowed



request *Bob* → *Alice*

request *Bob* → *Alice*

requests[*Alice*].add(*Bob*)

requests[*Alice*].add(*Bob*)

accept

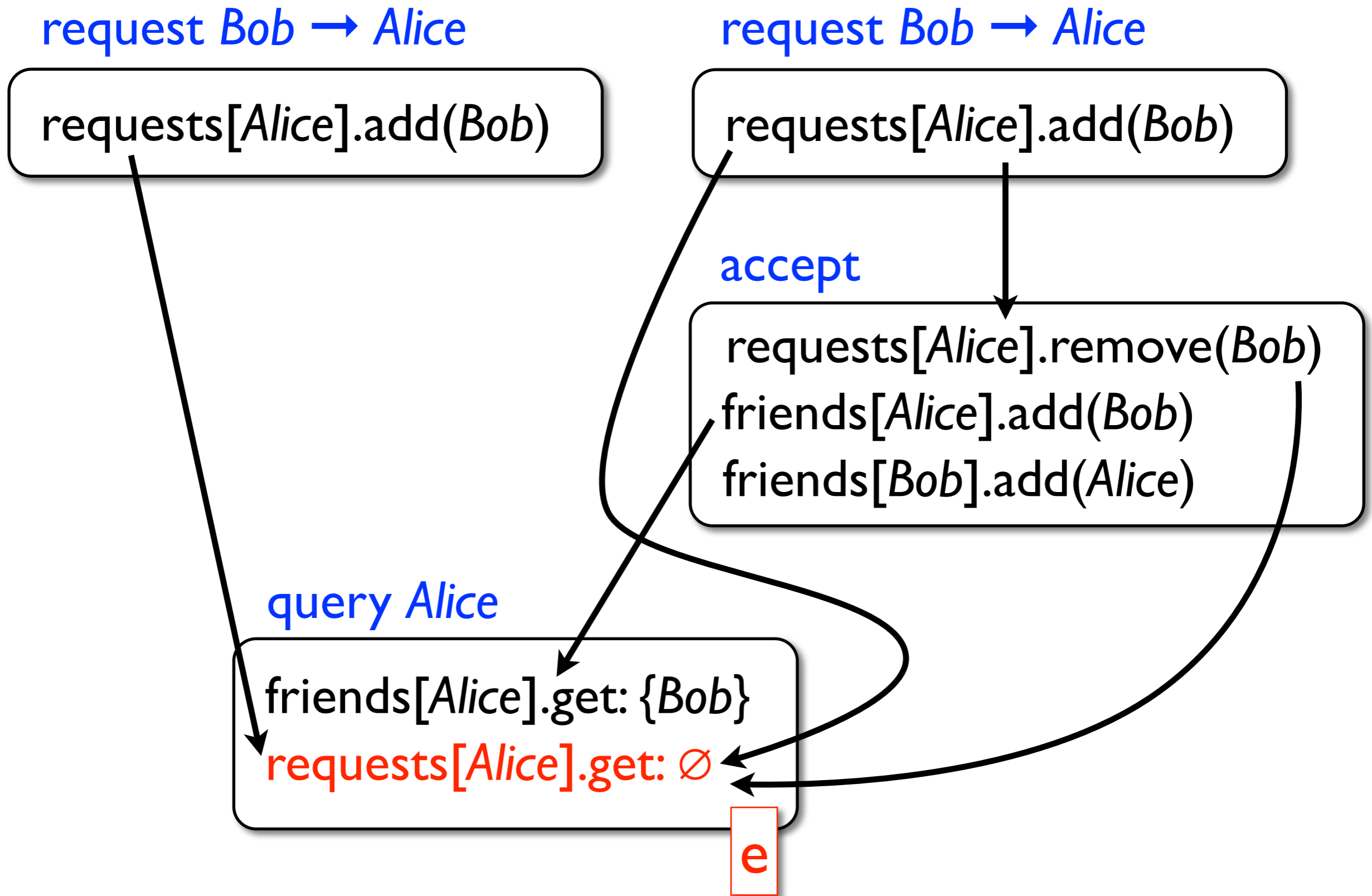
requests[*Alice*].remove(*Bob*)  
friends[*Alice*].add(*Bob*)  
friends[*Bob*].add(*Alice*)

query *Alice*

friends[*Alice*].get: {*Bob*}  
requests[*Alice*].get: ∅

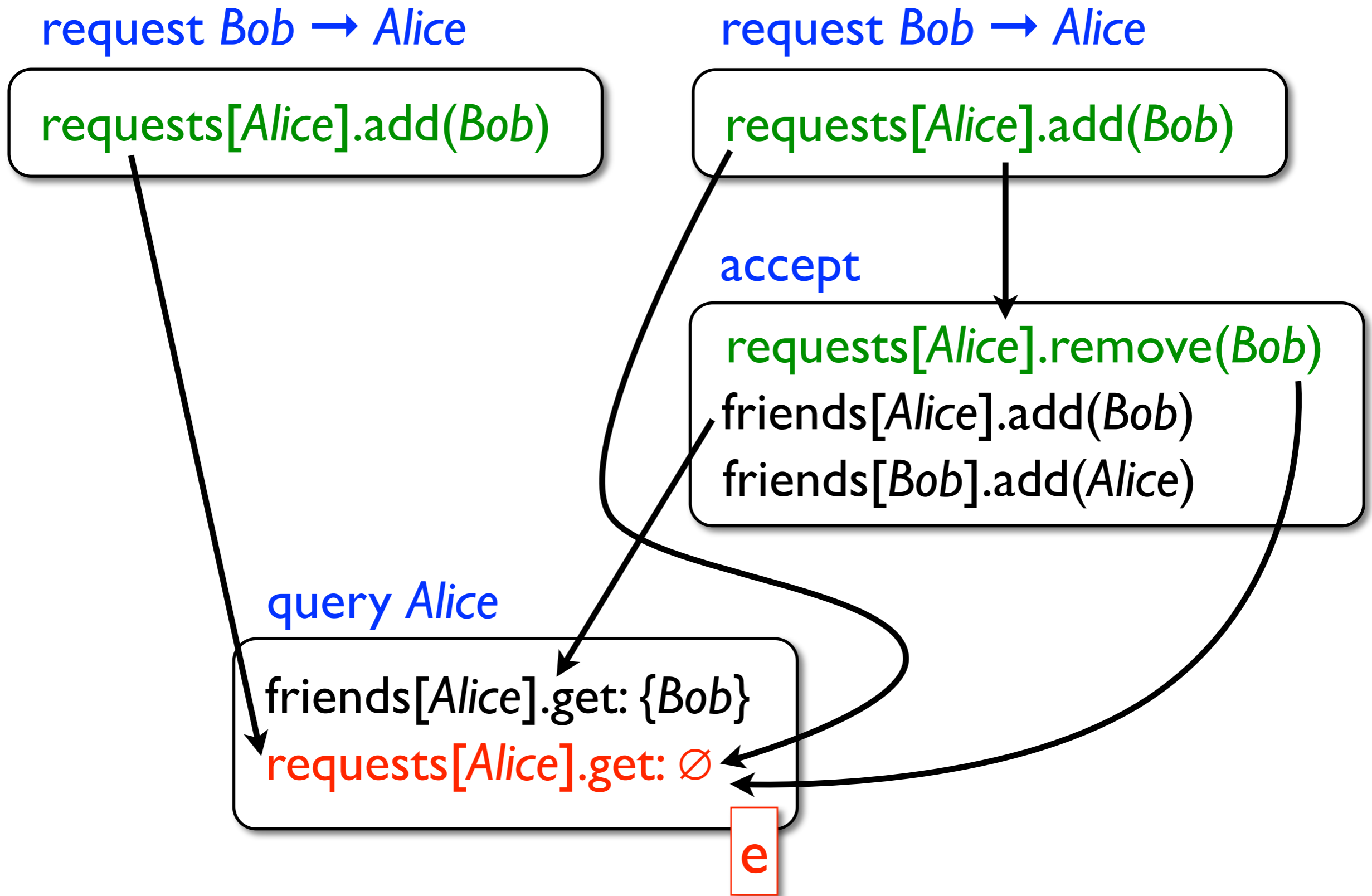
Database semantics = executions with correct return values

Determined by replicated data type specifications



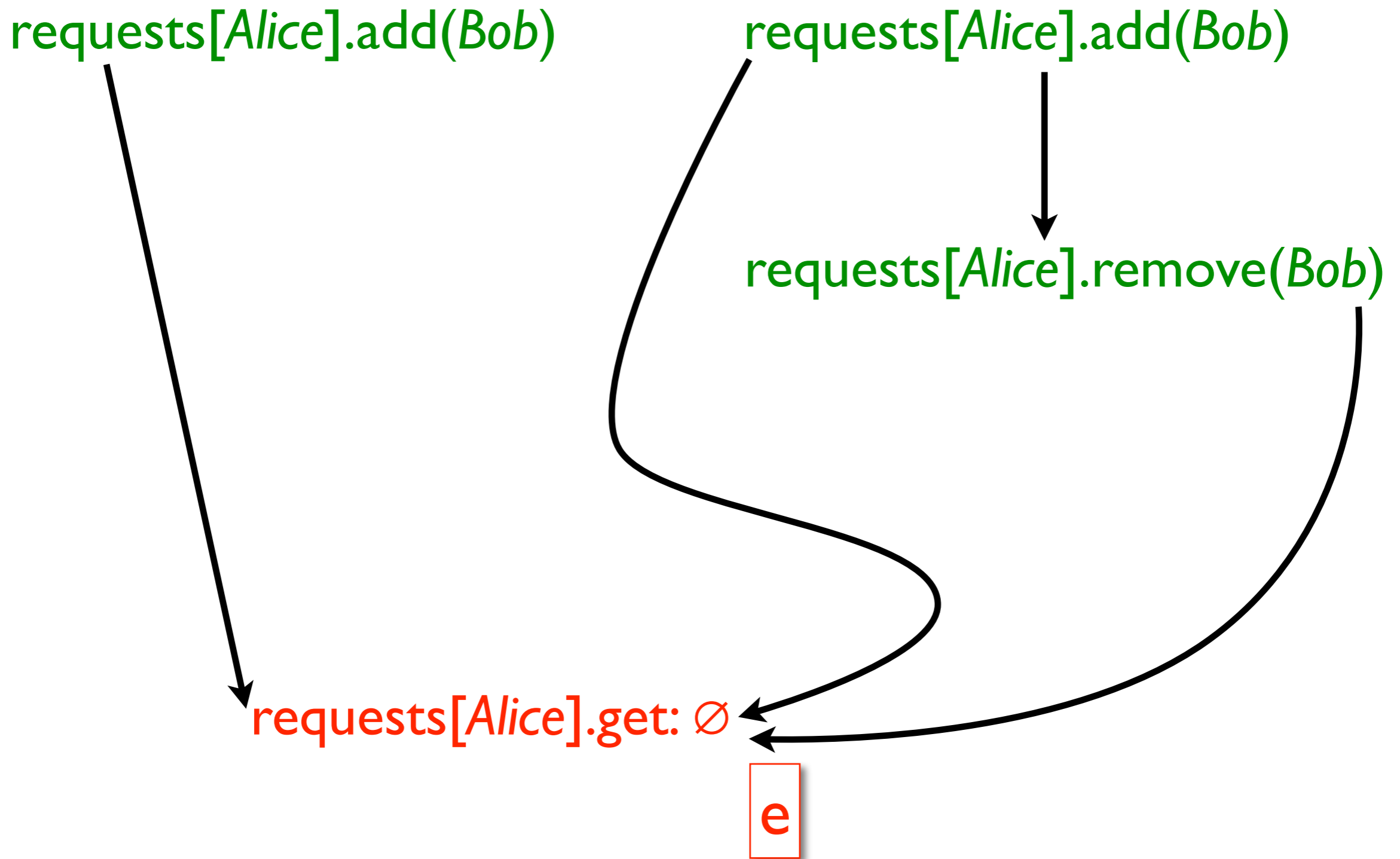
Database semantics = executions with correct return values

Determined by replicated data type specifications



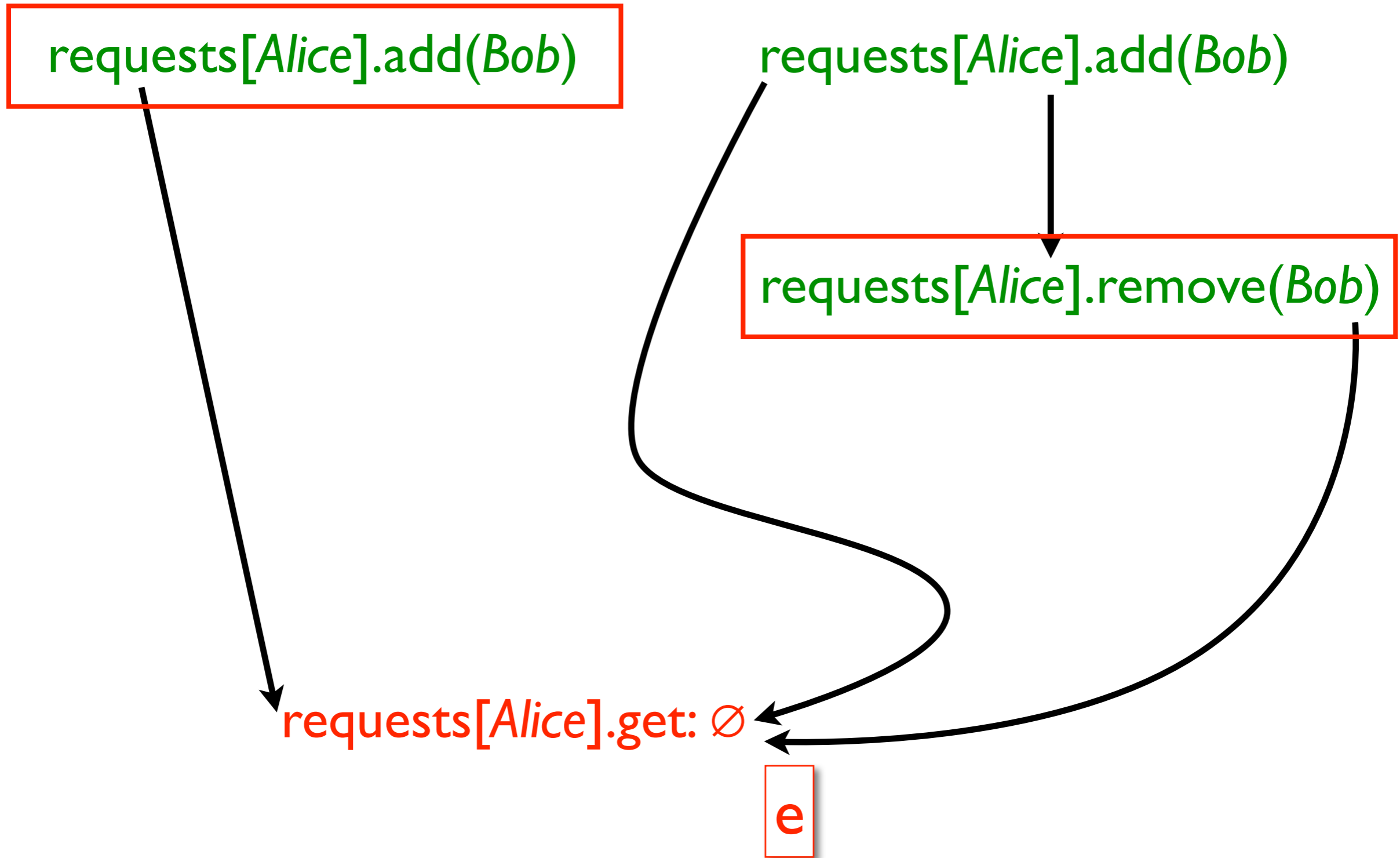
Database semantics = executions with correct return values

Determined by replicated data type specifications



Context(e): projection to events visible to e

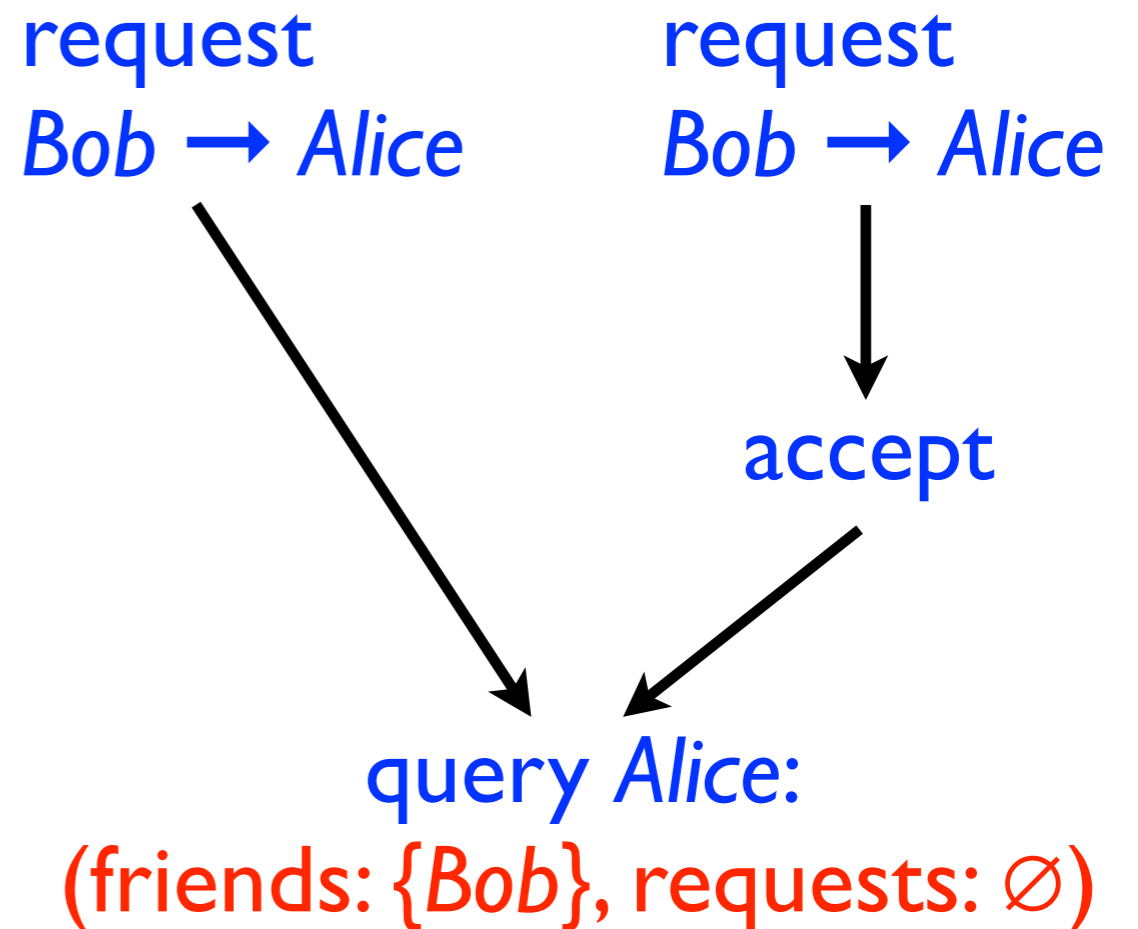
Spec S: Context(e)  $\rightarrow$  return value(e)



Remove-wins set: removes cancel adds  
unrelated to them in visibility

Contribution:  $D \rightarrow \llbracket D \rrbracket = S \in \text{Context} \rightarrow \text{Value}$

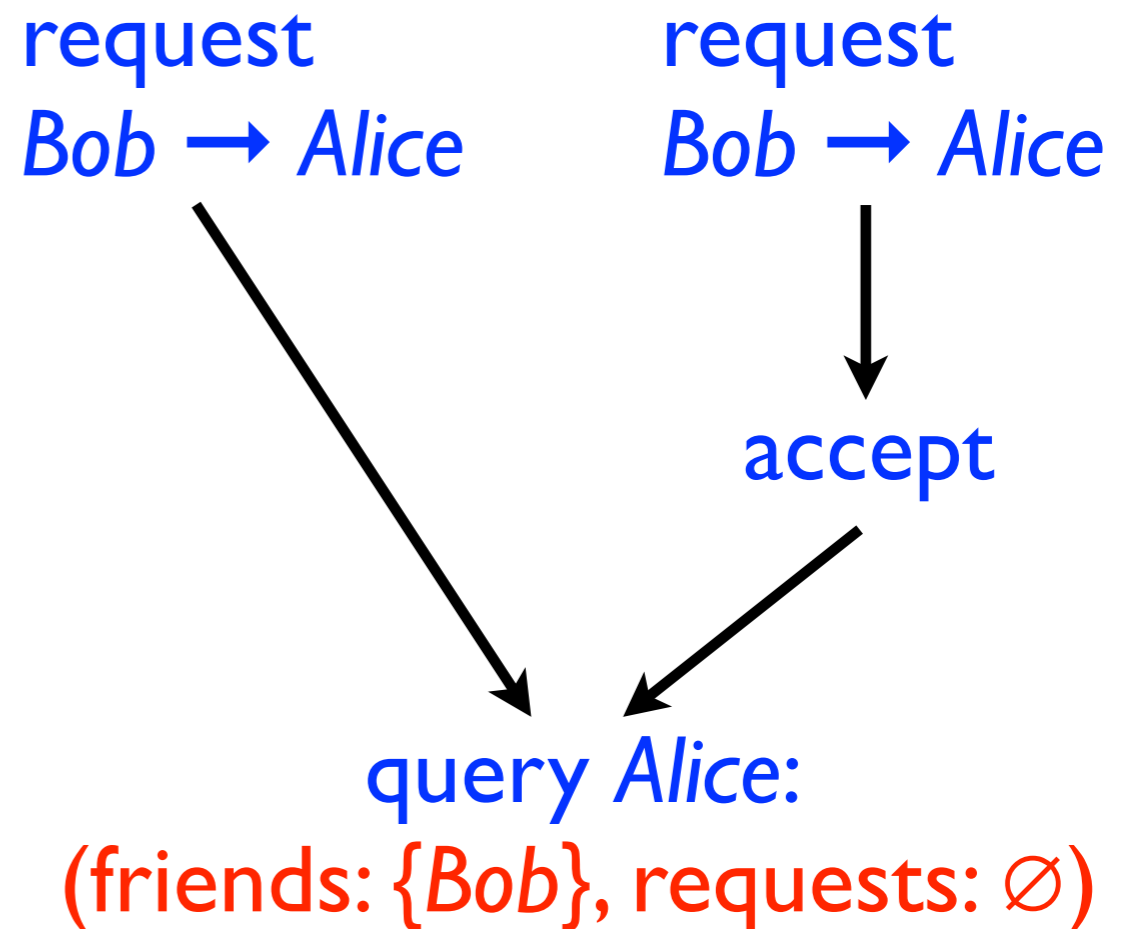
```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```



- Compute method return value based on **coarse-grain** events ~ executed methods
- No information about **fine-grain** internal events

Contribution:  $D \rightarrow \llbracket D \rrbracket = S \in \text{Context} \rightarrow \text{Value}$

```
datatype SocialGraph {  
  RemoveWinsSet friends[];  
  RemoveWinsSet requests[];  
  
  tx accept(Bob → Alice) {  
    ...  
    requests[Alice].remove(Bob);  
    friends[Alice].add(Bob);  
    friends[Bob].add(Alice);  
    ...  
  }  
}
```



- Compute method return value based on **coarse-grain** events ~ executed methods
- No information about **fine-grain** internal events

# Coarse-grain context N

request *Bob* → *Alice*



request *Bob* → *Alice*



accept



query *Alice*

?



Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$

request  $Bob \rightarrow Alice$



request  $Bob \rightarrow Alice$



accept



query  $Alice$

?

requests[Alice].add(Bob)



requests[Alice].add(Bob)



requests[Alice].remove(Bob)  
friends[Alice].add(Bob)  
friends[Bob].add(Alice)



friends[Alice].get: {Bob}  
requests[Alice].get:  $\emptyset$

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$

request  $Bob \rightarrow Alice$



request  $Bob \rightarrow Alice$



accept



query  $Alice$

(friends:  $\{Bob\}$ , requests:  $\emptyset$ )



requests[Alice].add(Bob)



requests[Alice].add(Bob)

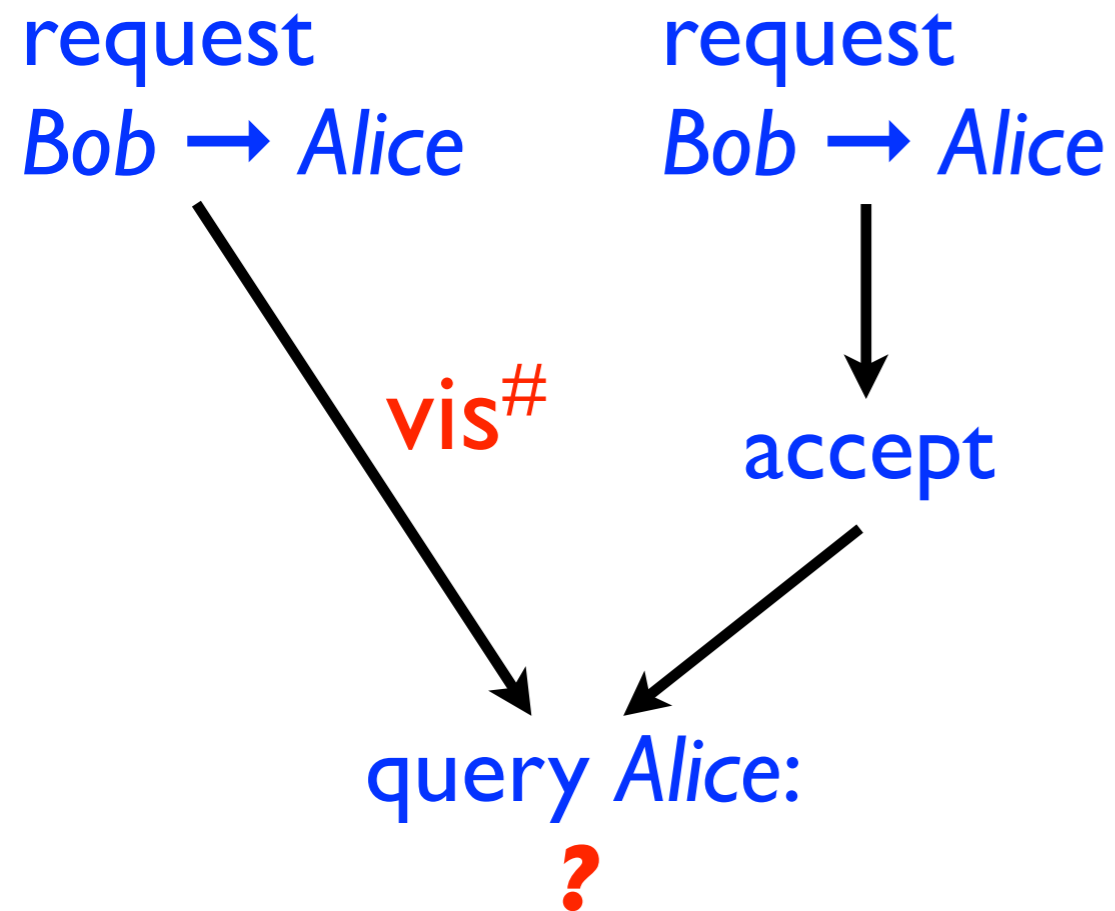


requests[Alice].remove(Bob)  
friends[Alice].add(Bob)  
friends[Bob].add(Alice)



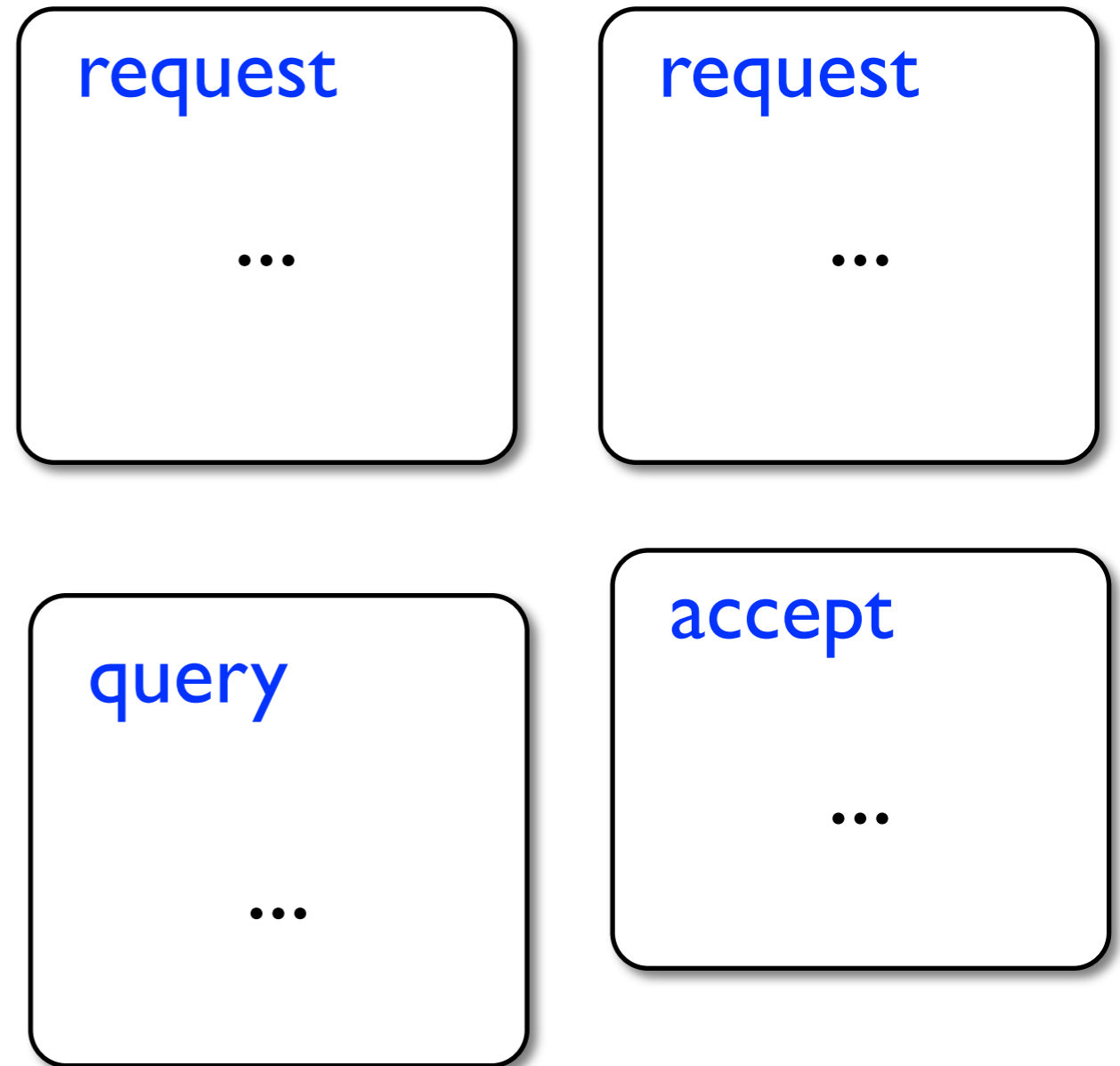
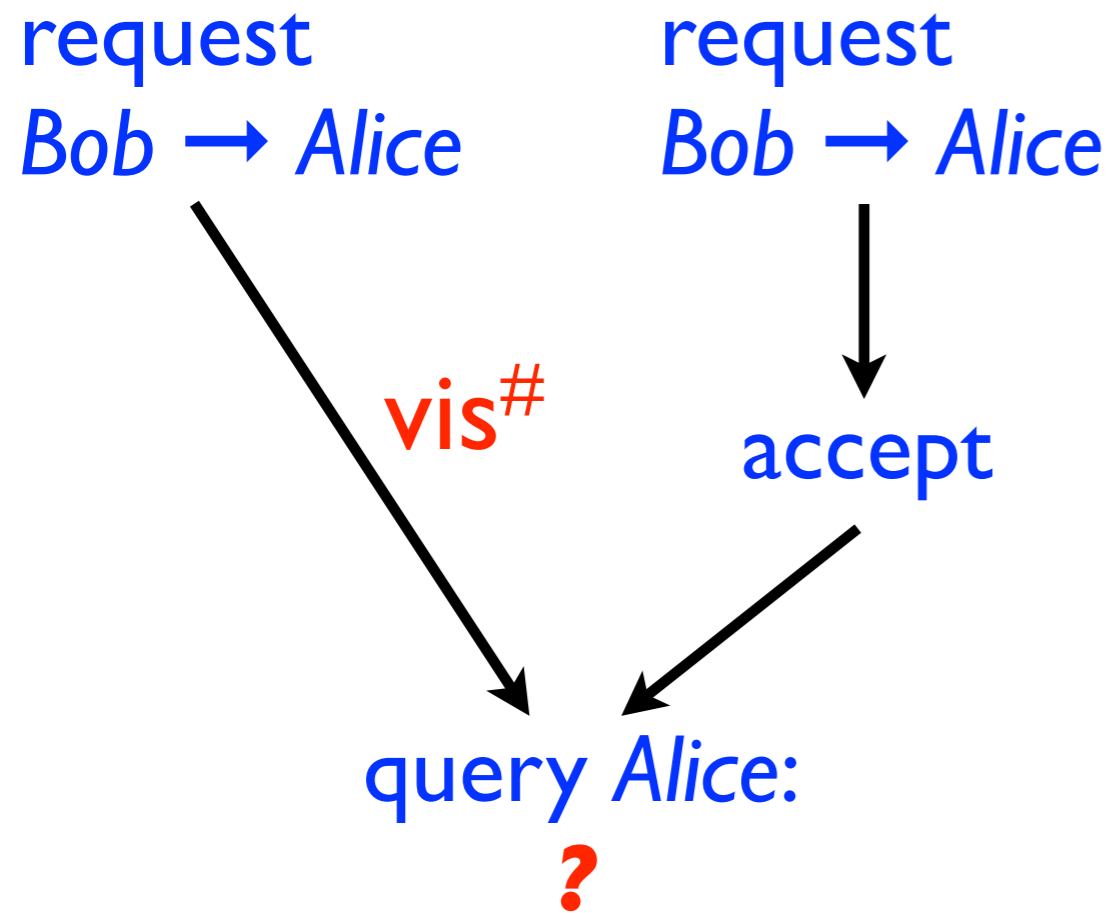
friends[Alice].get:  $\{Bob\}$   
requests[Alice].get:  $\emptyset$

# Coarse-grain context N



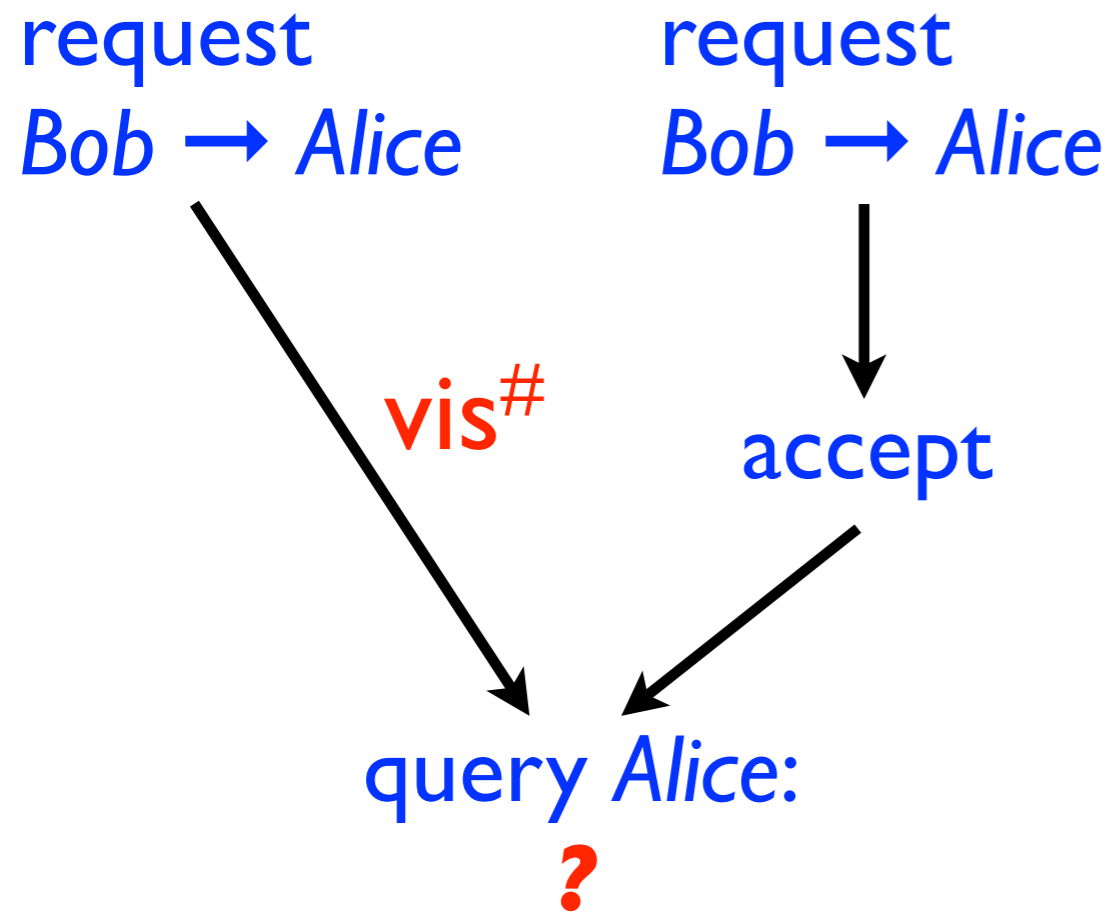
Given visibility  $vis^\#$  on whole methods = txns

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$

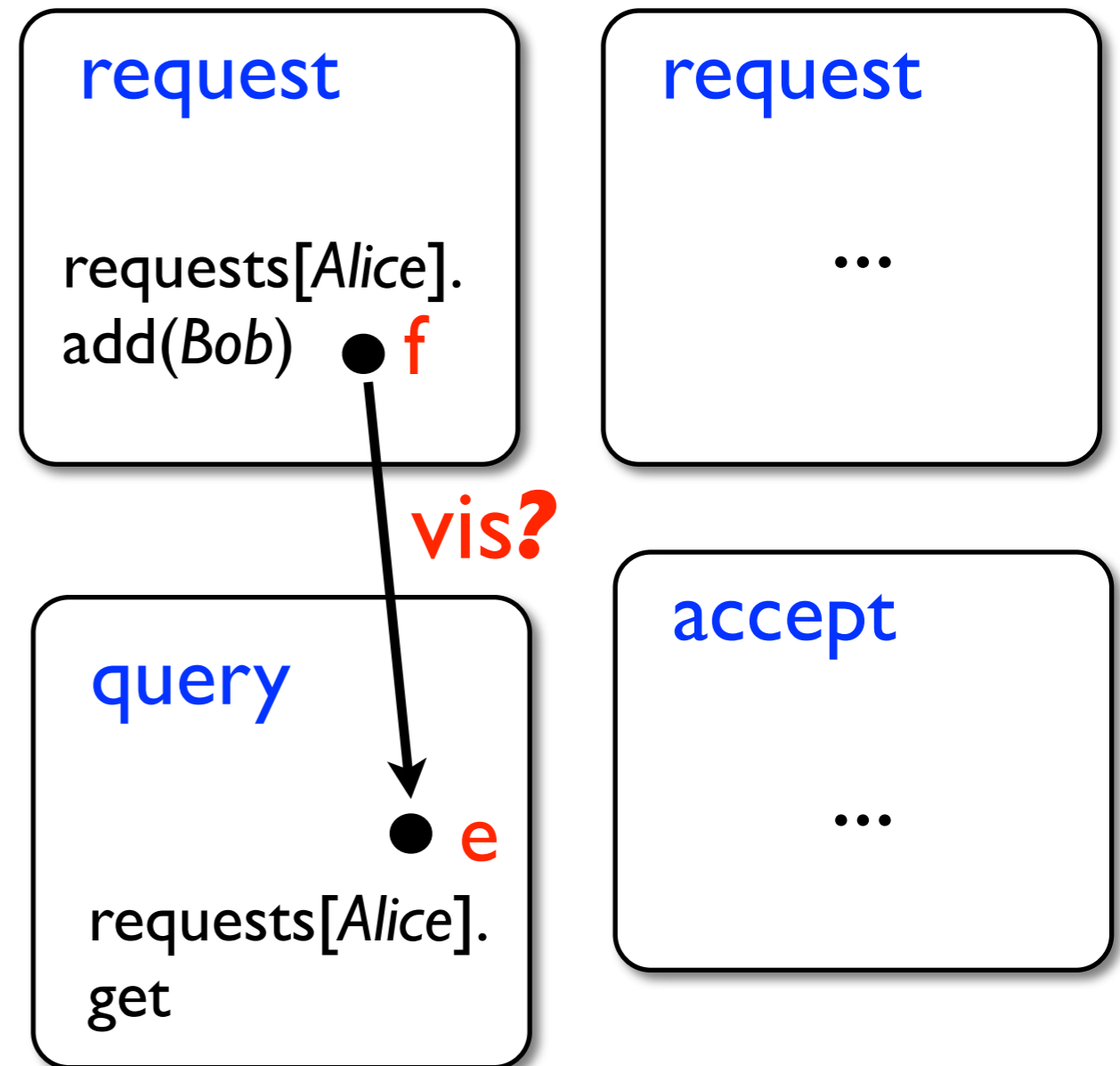


Given visibility  $vis^\#$  on  
whole methods = txns

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$

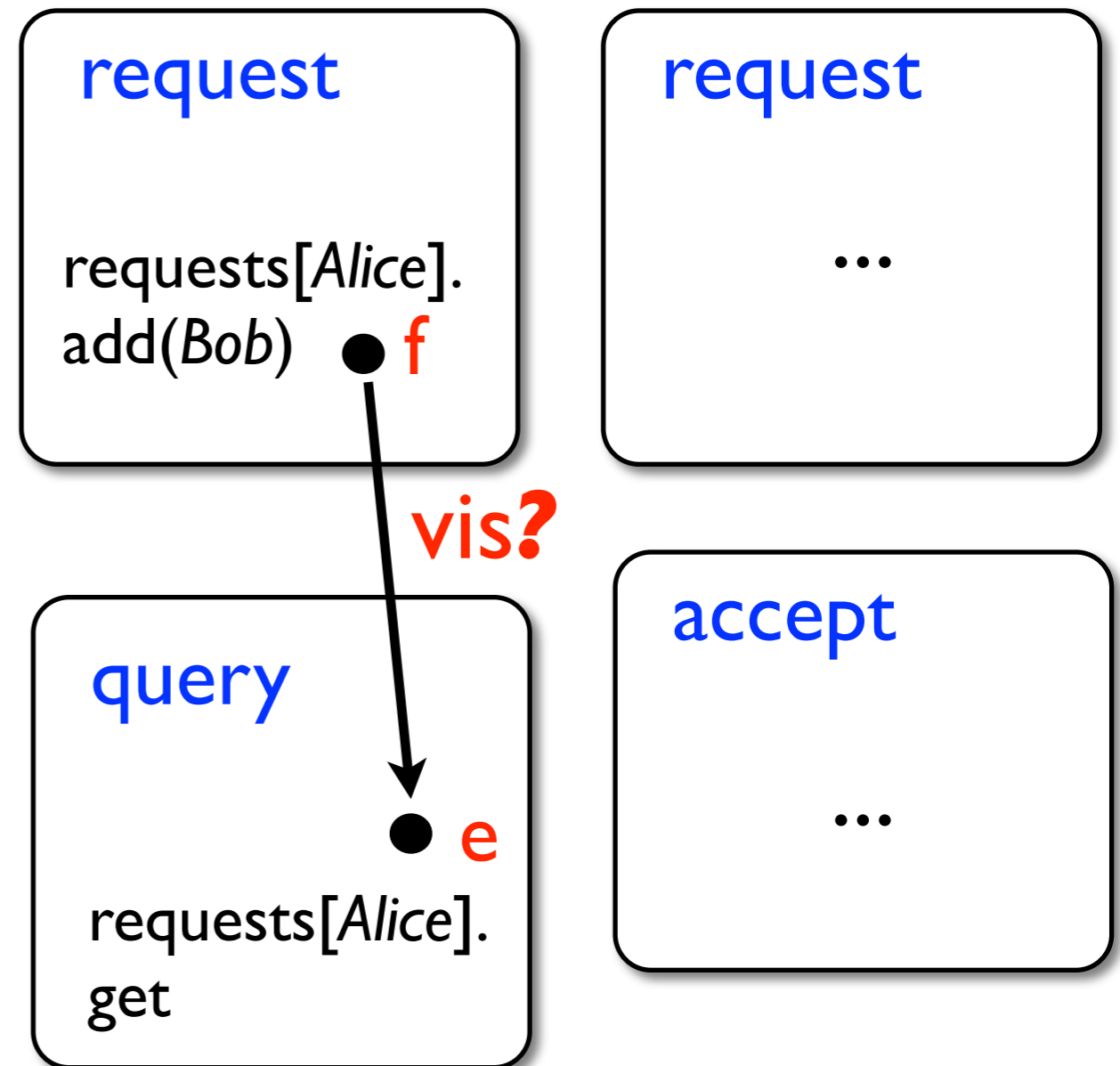
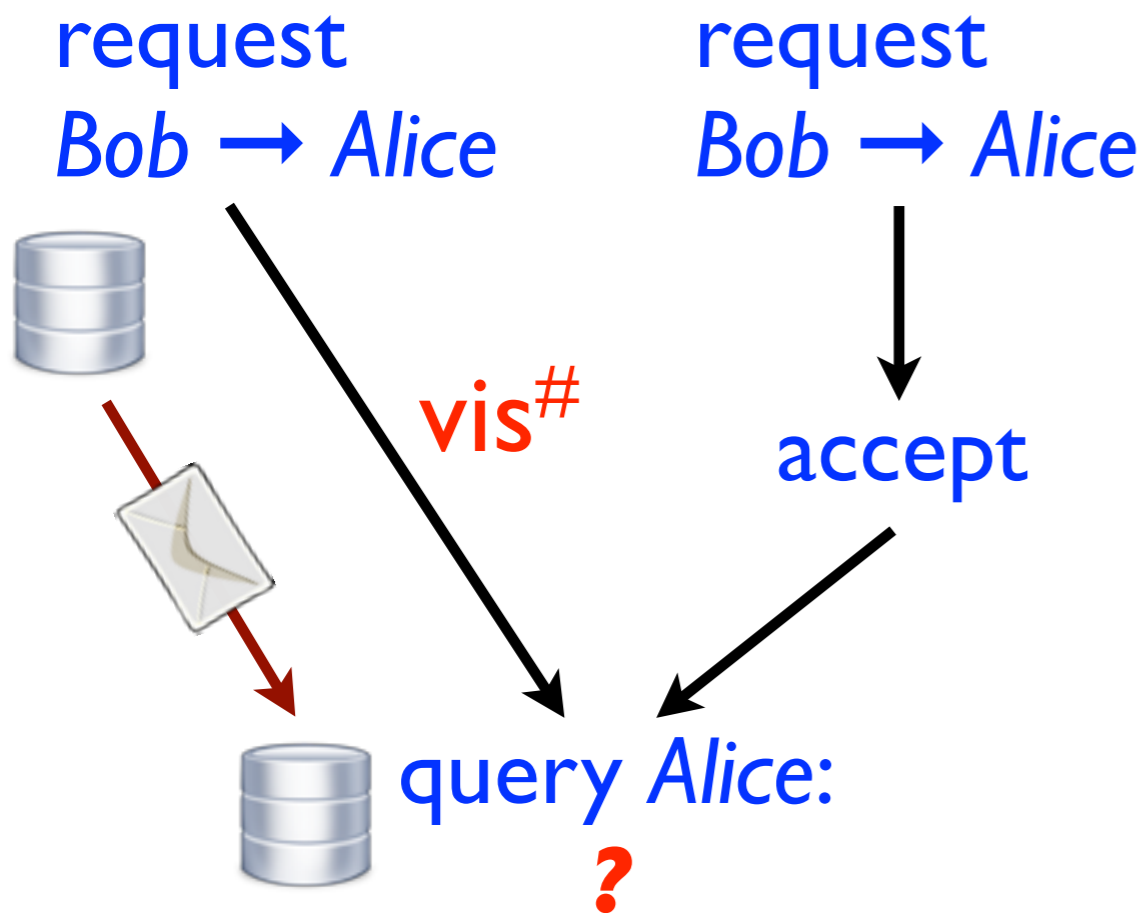


Given visibility  $\text{vis}^\#$  on whole methods = txns



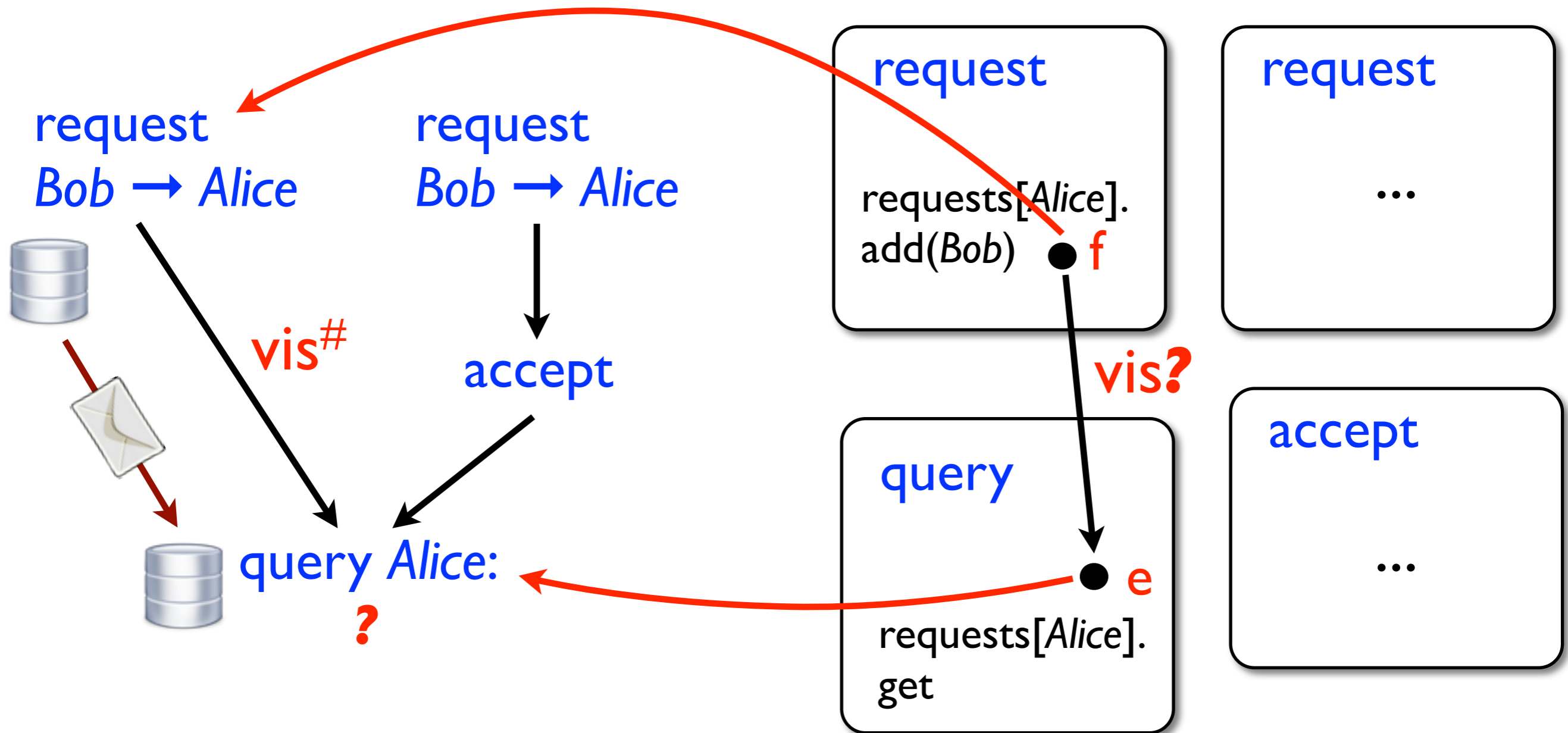
Find visibility  $\text{vis}$  on method-internal events

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$



Visibility = message delivery  
All txn updates delivered together

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$



Visibility = message delivery

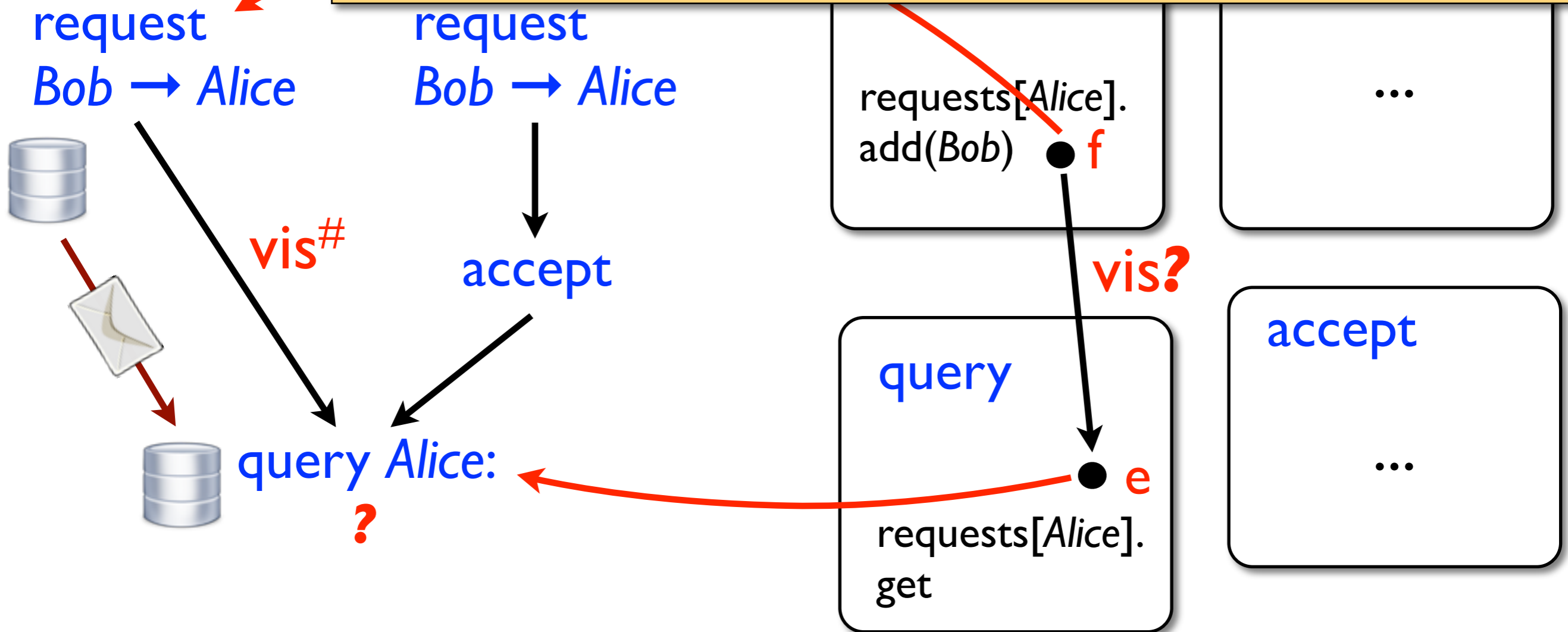
All txn updates delivered together

$$(f, e) \in vis \iff (\text{method}(f), \text{method}(e)) \in vis^\#$$

Coarse-grain

Only a constraint on fine-grain executions

**Theorem:**  $\gamma(N)$  is defined uniquely



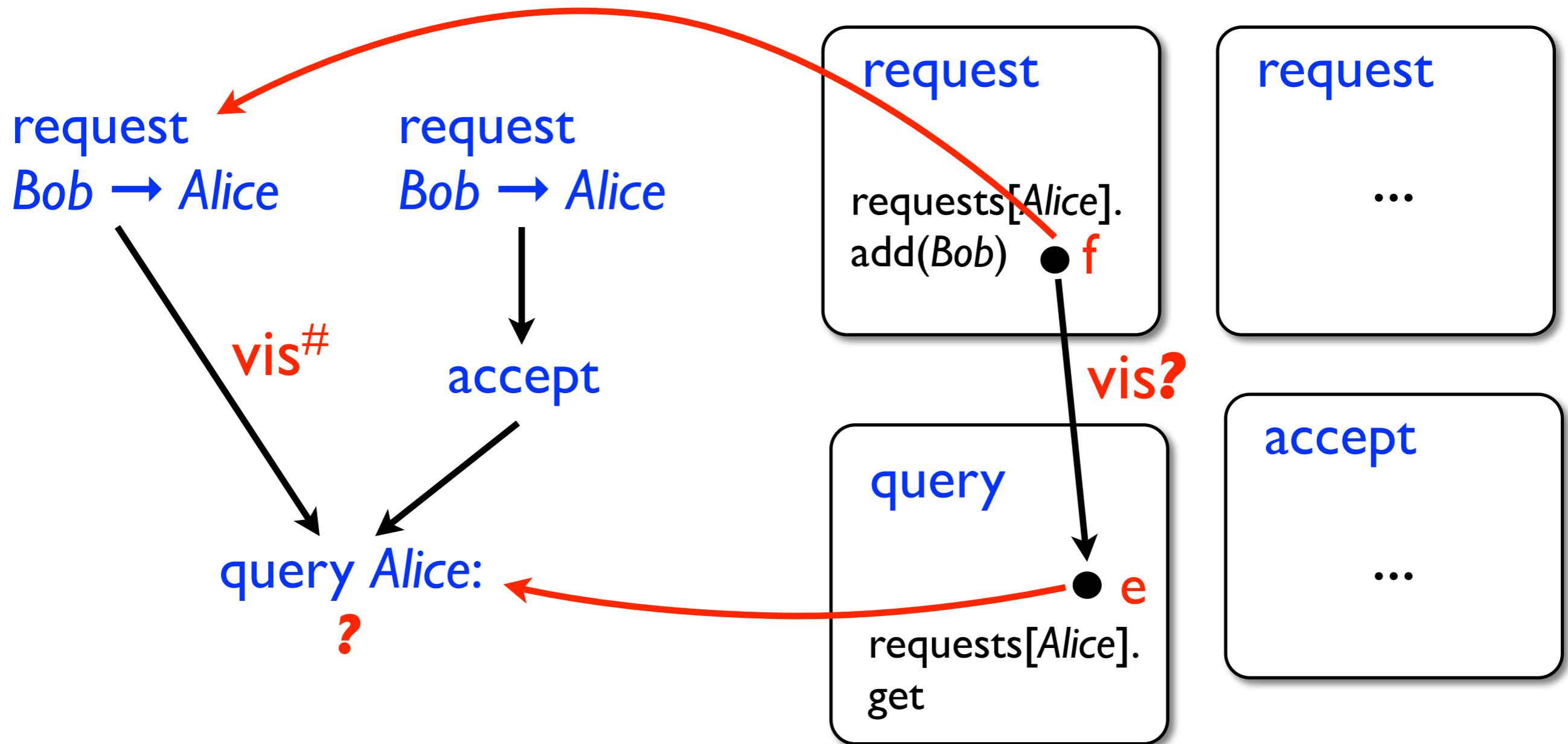
Visibility = message delivery

All txn updates delivered together

$$(f, e) \in vis \iff (\text{method}(f), \text{method}(e)) \in vis^\#$$



Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$

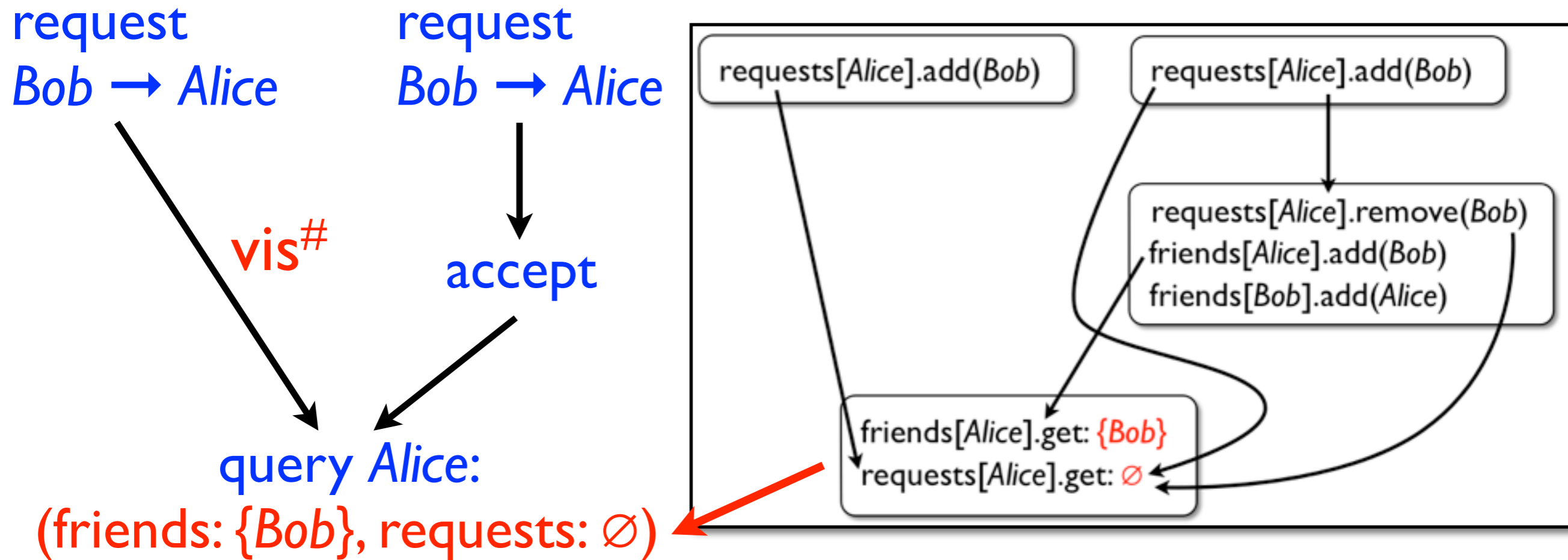


To determine  $\llbracket D \rrbracket$   
on context  $N$  of  
operation  $op$

Determine the  
concretisation  
 $\gamma(N)$  wrt  $D$

Take the return  
value of  $op$  in  
 $\gamma(N)$

Coarse-grain context  $N \rightarrow$  Fine-grain execution  $\gamma(N)$



To determine  $\llbracket D \rrbracket$   
on context  $N$  of  
operation  $op$

Determine the  
concretisation  
 $\gamma(N)$  wrt  $D$

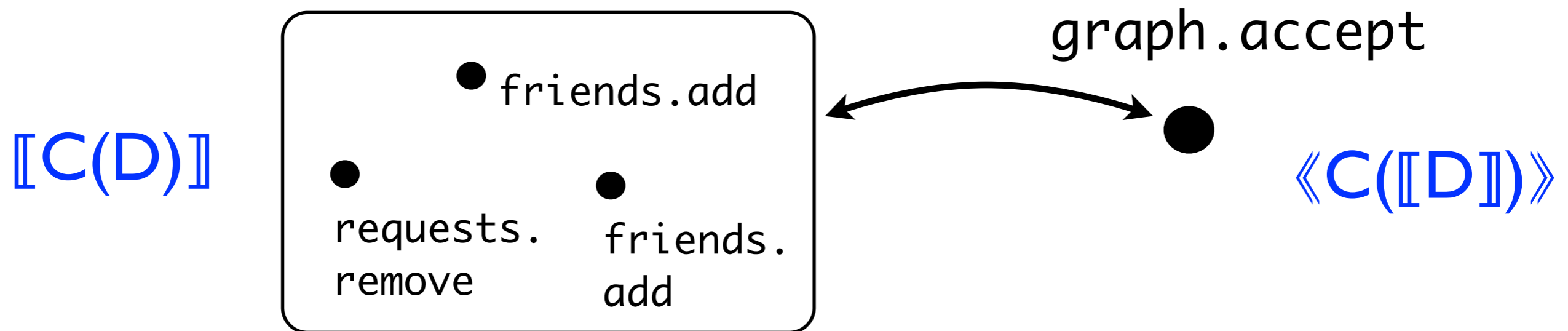
Take the return  
value of  $op$  in  
 $\gamma(N)$

**Theorem:**  $\llbracket C(D) \rrbracket = \langle\langle C(\llbracket D \rrbracket) \rangle\rangle$

- Coarse-grain and fine-grain semantics coincide
- Client behaviour the same whether using implementation  $D$  or spec  $\llbracket D \rrbracket$

**Theorem:**  $\llbracket C(D) \rrbracket = \langle\langle C(\llbracket D \rrbracket) \rangle\rangle$

- Consistency axioms are **global**: constrain the whole execution
- Preserve axioms when performing **local** execution transformations:



- 93-page extended version, enjoy!

# Using the denotation

- $\llbracket D \rrbracket$  not effectively computable:
  1. Write a spec  $S$  describing desirable behaviour of  $D$  declaratively
  2. Prove  $S = \llbracket D \rrbracket$
- Definition of  $\llbracket D \rrbracket$  gives a **proof method**
- Examples: social graph, shopping cart

# Conclusion

- Modular reasoning method for applications using eventually consistent databases: reflects application modularity
- Only causal consistency model: general theory?
- Opportunity: extend rich programming language theories to modern distributed systems