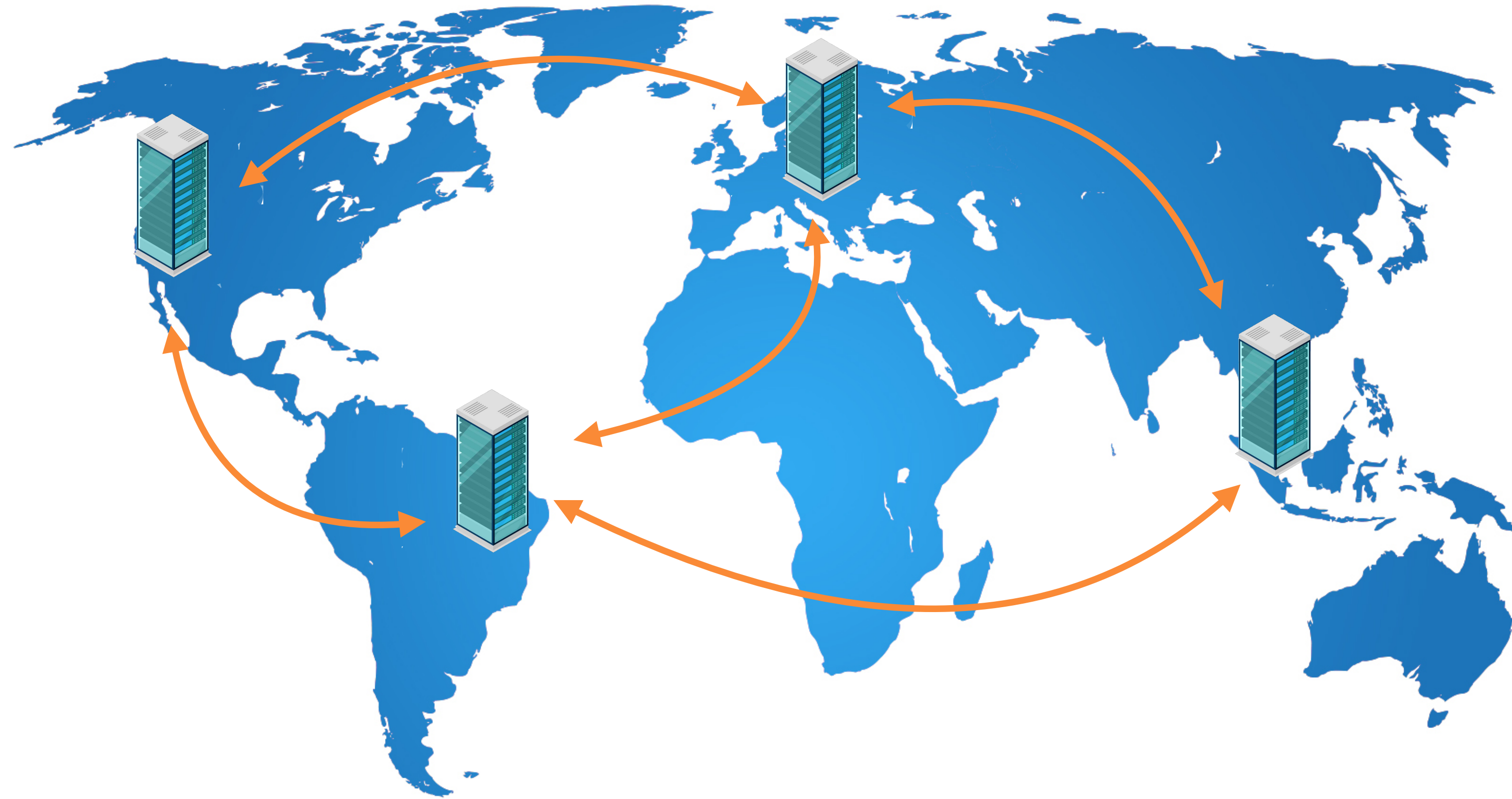# CRDTs in Production

Dmitry Martyanov, Software Engineer @ PayPal

QCon, 2018

# Geo-Distributed Datastore

# Context

- More than 200 countries

- Regulatory requirements

- State Machine of Compliance Status

- Modified by multiple Actors
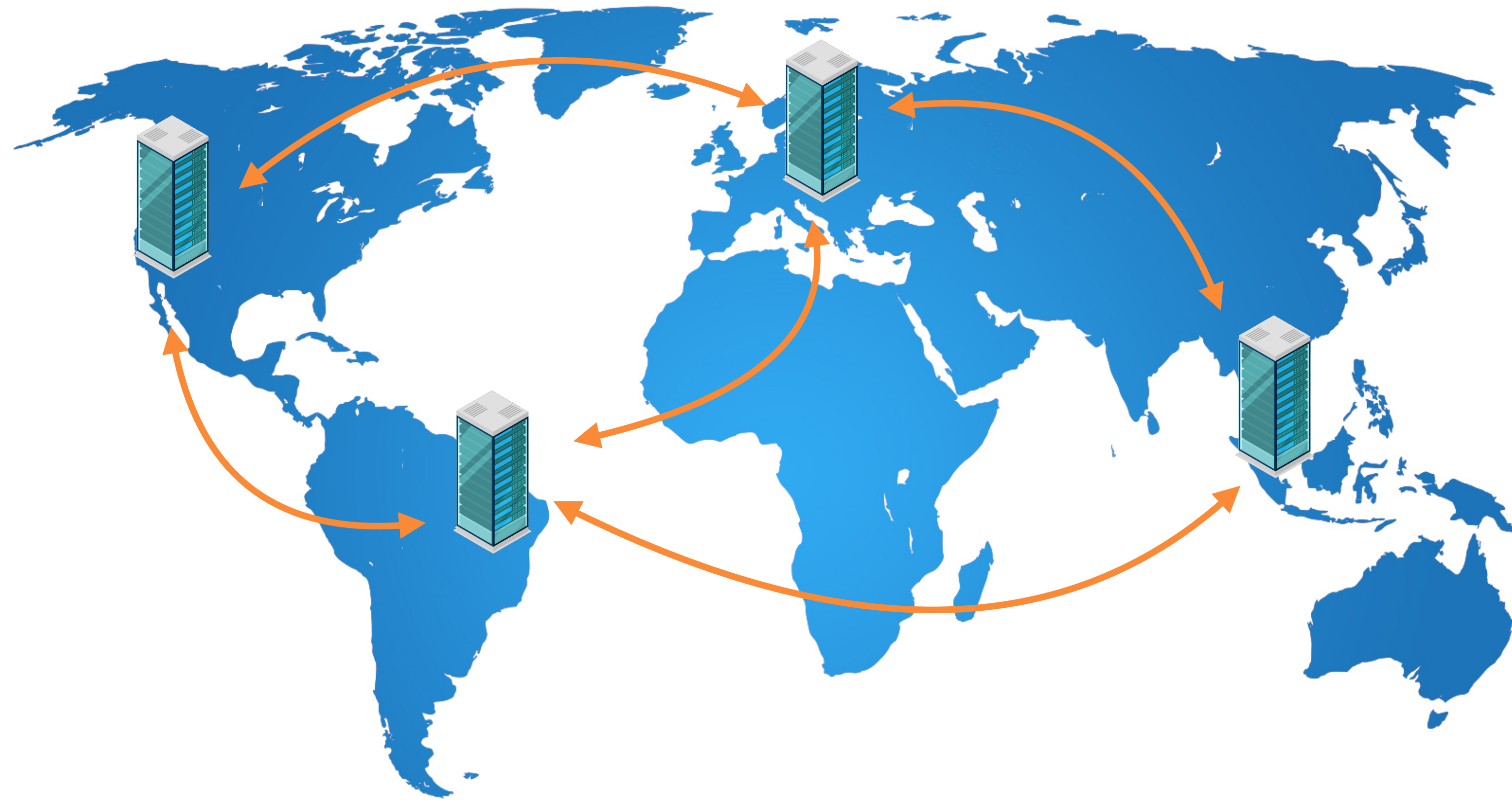
# Shared Mutable State

# Shared Mutable State

# Mutex
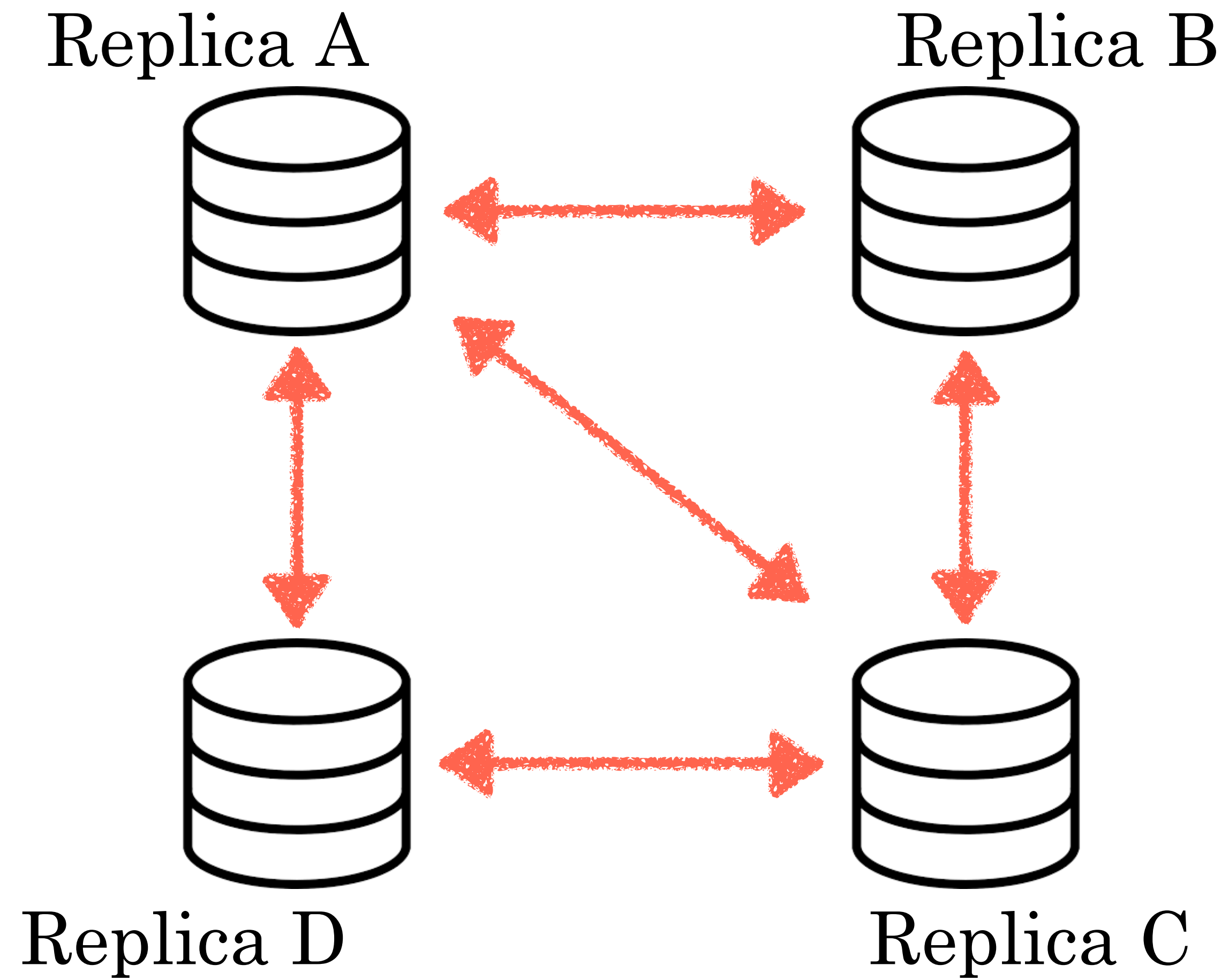
# Shared Mutable State
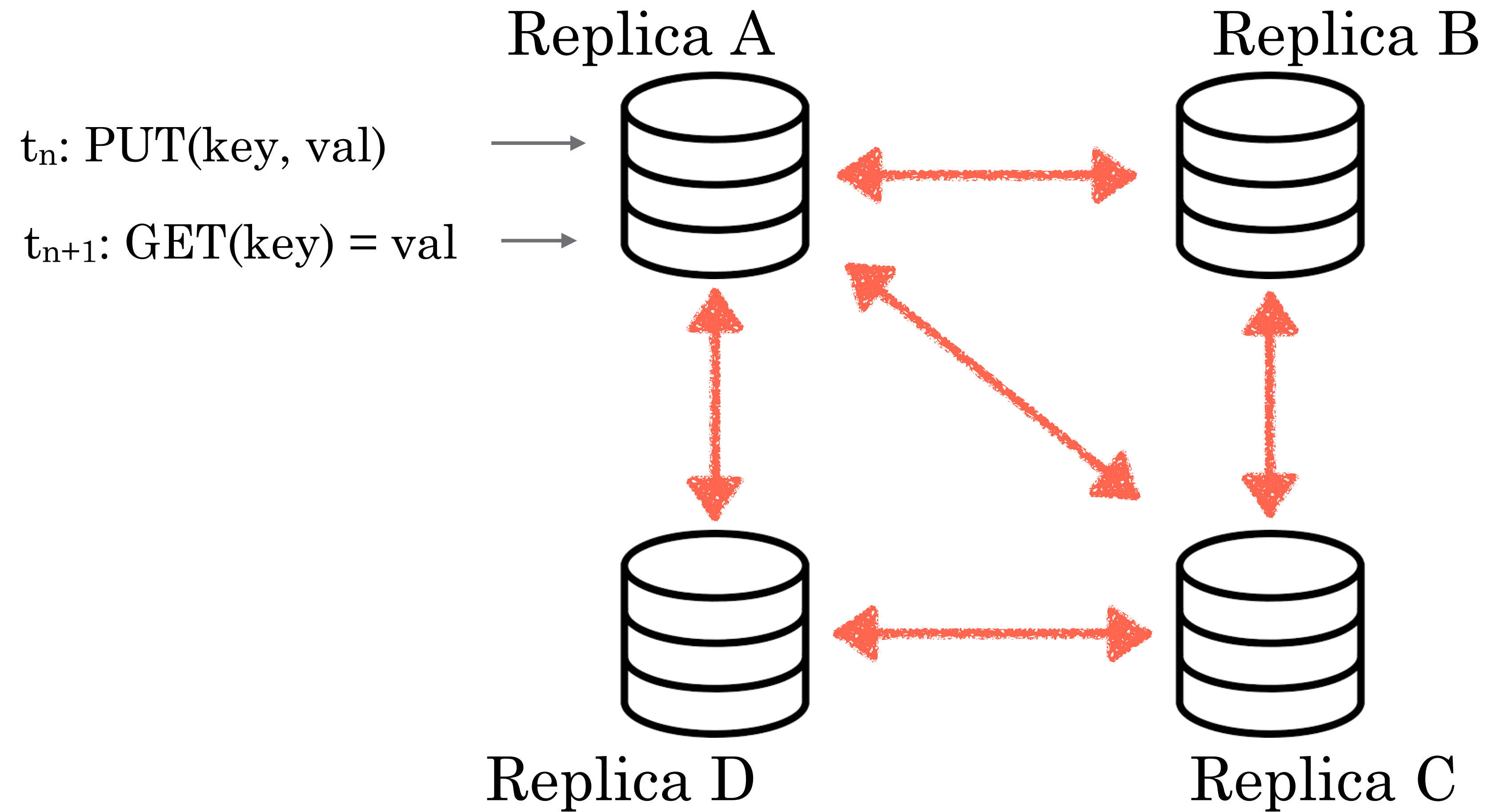
# Mutex

# Transactions

# Geo-Distributed Datastore

# Eventual Consistency

# Distributed System

# Distributed System

$t_n$: PUT(key, val)

$t_{n+1}$: GET(key) = val

Replica A

Replica B

Replica D

Replica C

# Distributed System



$t_n$: PUT(key, val)

$t_{n+1}$: GET(key) = val

Replica A

Replica B

$t_{n+1}$: GET(key) = **?**

Replica D

Replica C

# Affinity Based Approaches

# Affinity Based Approaches

Replica A

Replica B

Replica D

Replica C

# Coordinator Based Approaches



Replica A

Replica B

Replica D

Replica C

# Consensus Based Approaches

**Paxos, Raft, etc.**

# Service Stack



**What type of documents**
**Business Rules**
**Filtering Logic, etc.**

**Entity objects**
**DAO Layer**
**Flow Control**

**Service discovery**
**Routing & Balancing**
**Failover strategy**

**Service**
- Business Logic
- Domain Platform
- Service Infrastructure

**Datastore**
- Domain Data
- Data Infrastructure

**Data Model, Records**

**Deployment configuration**
**Namespaces & Schemas**
**Replication params**

# Service Stack
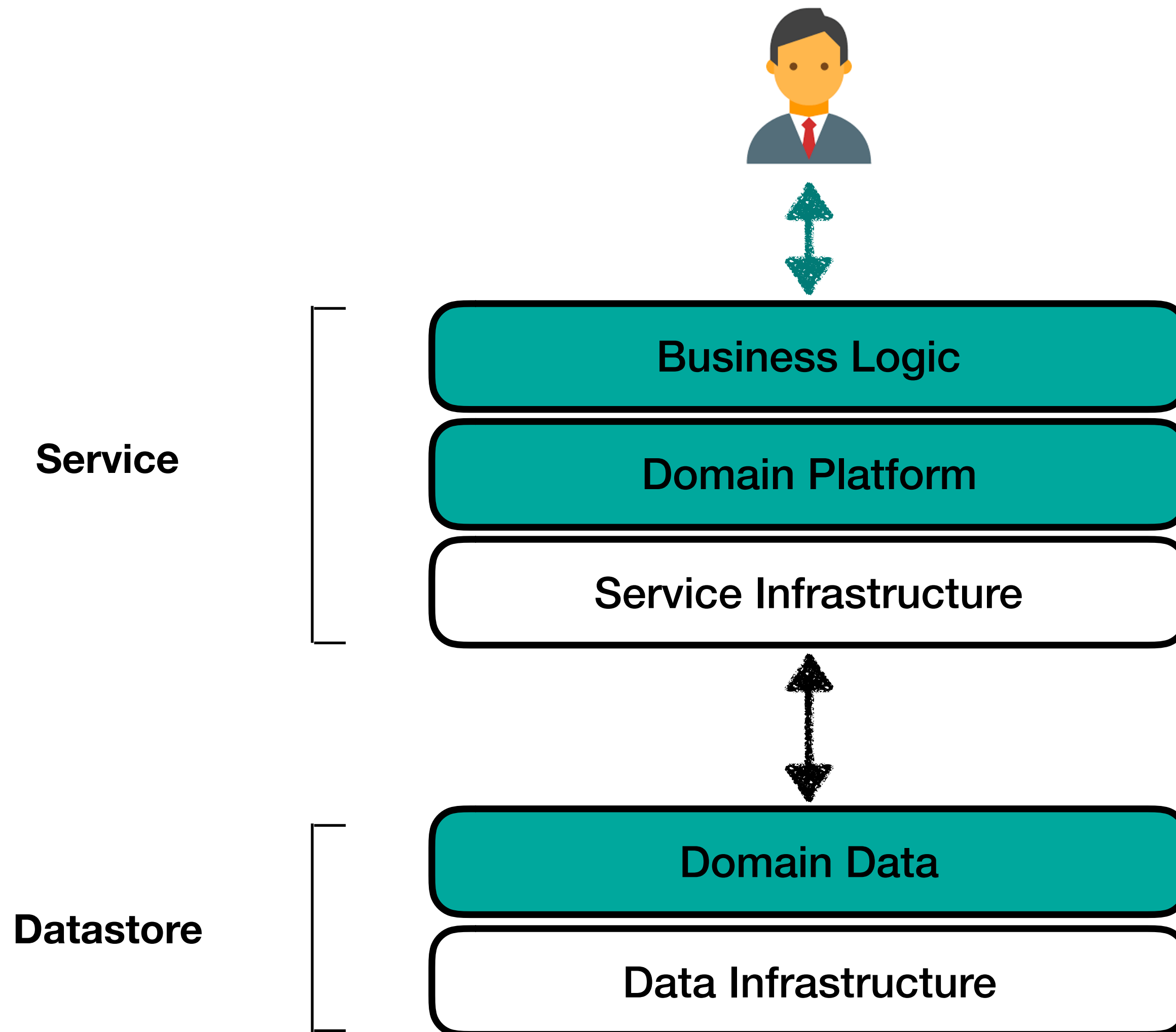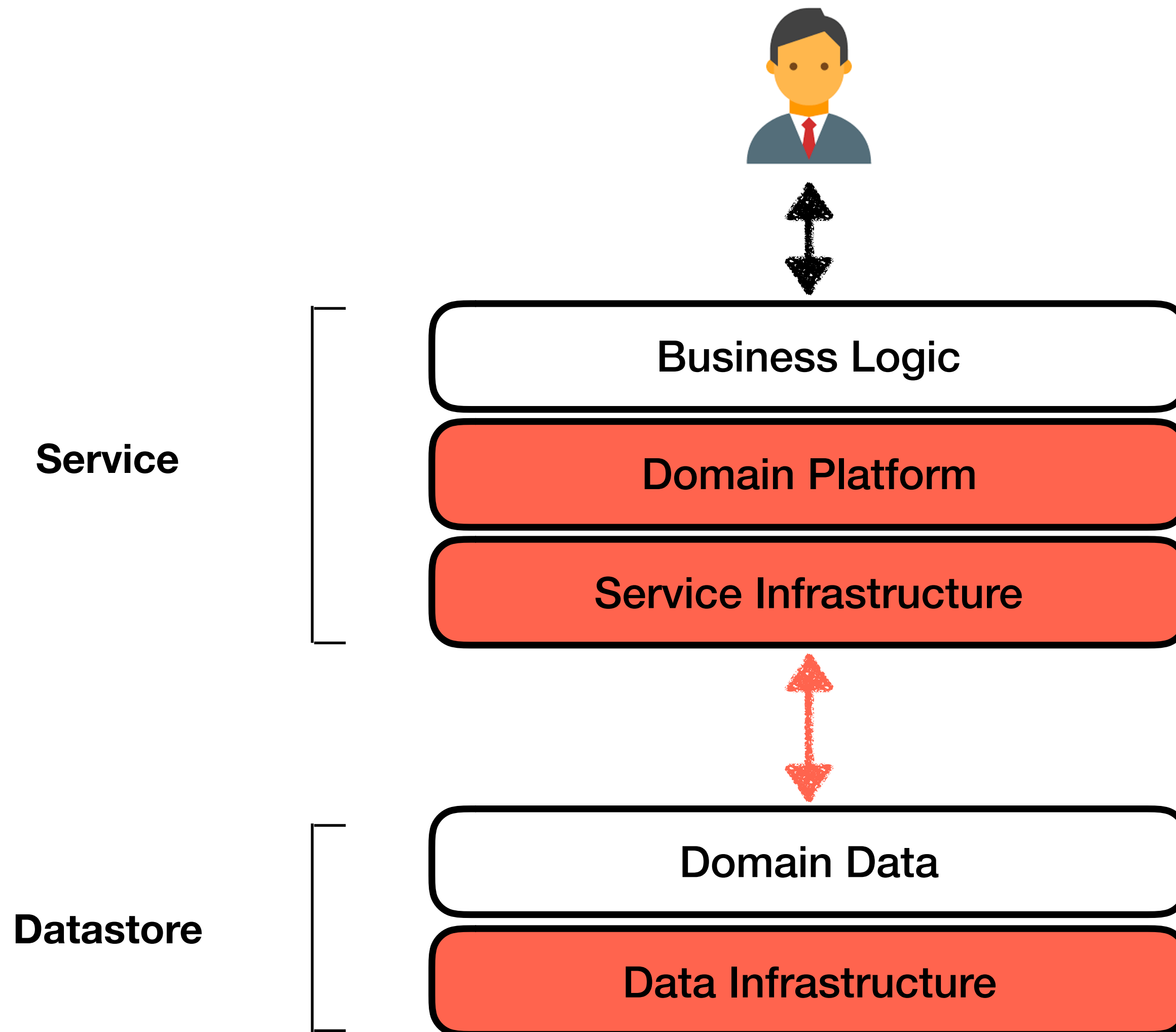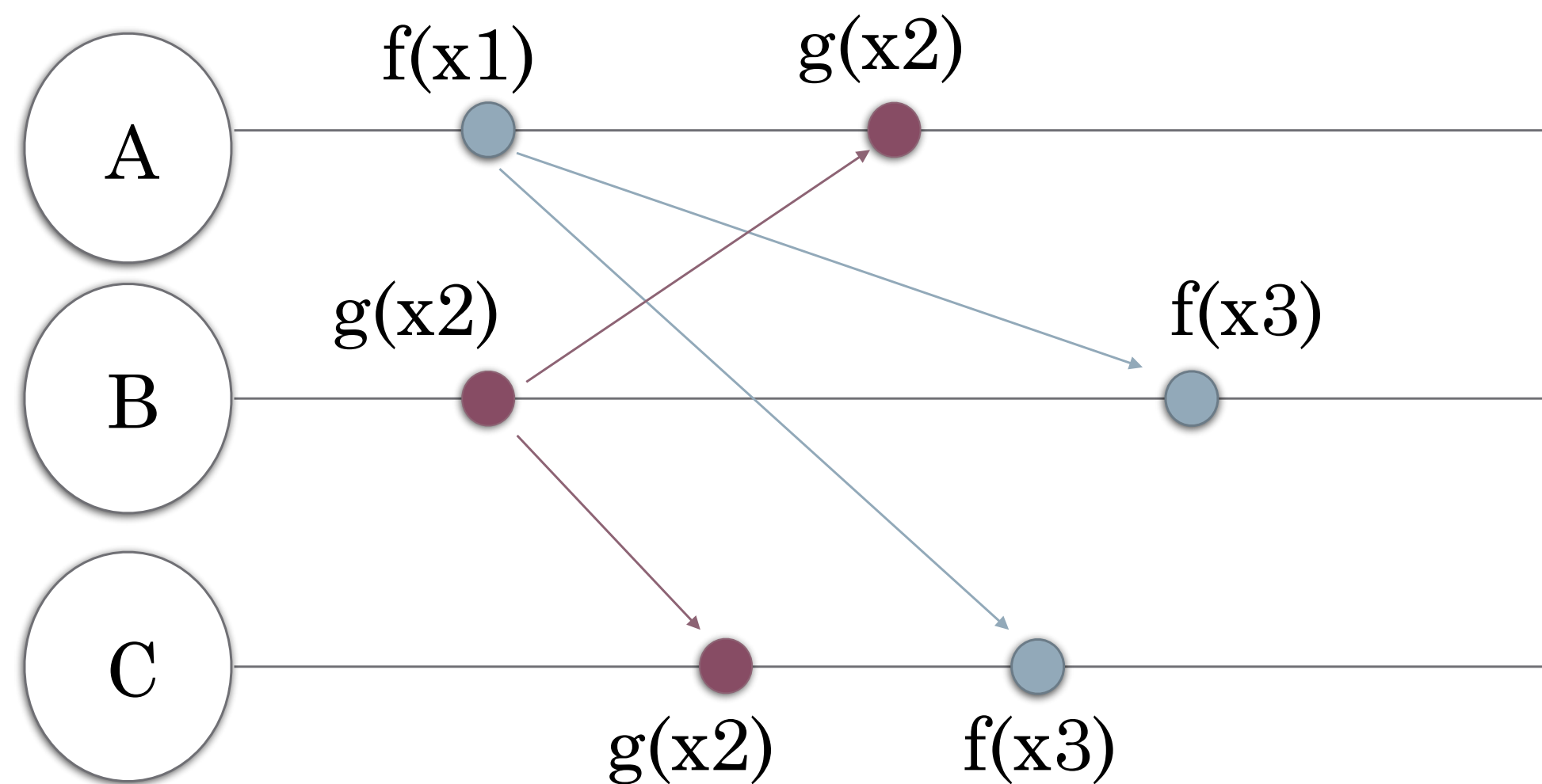
# Service Stack

# Conflict-free Replicated Data Types

# CRDTs

## commutative
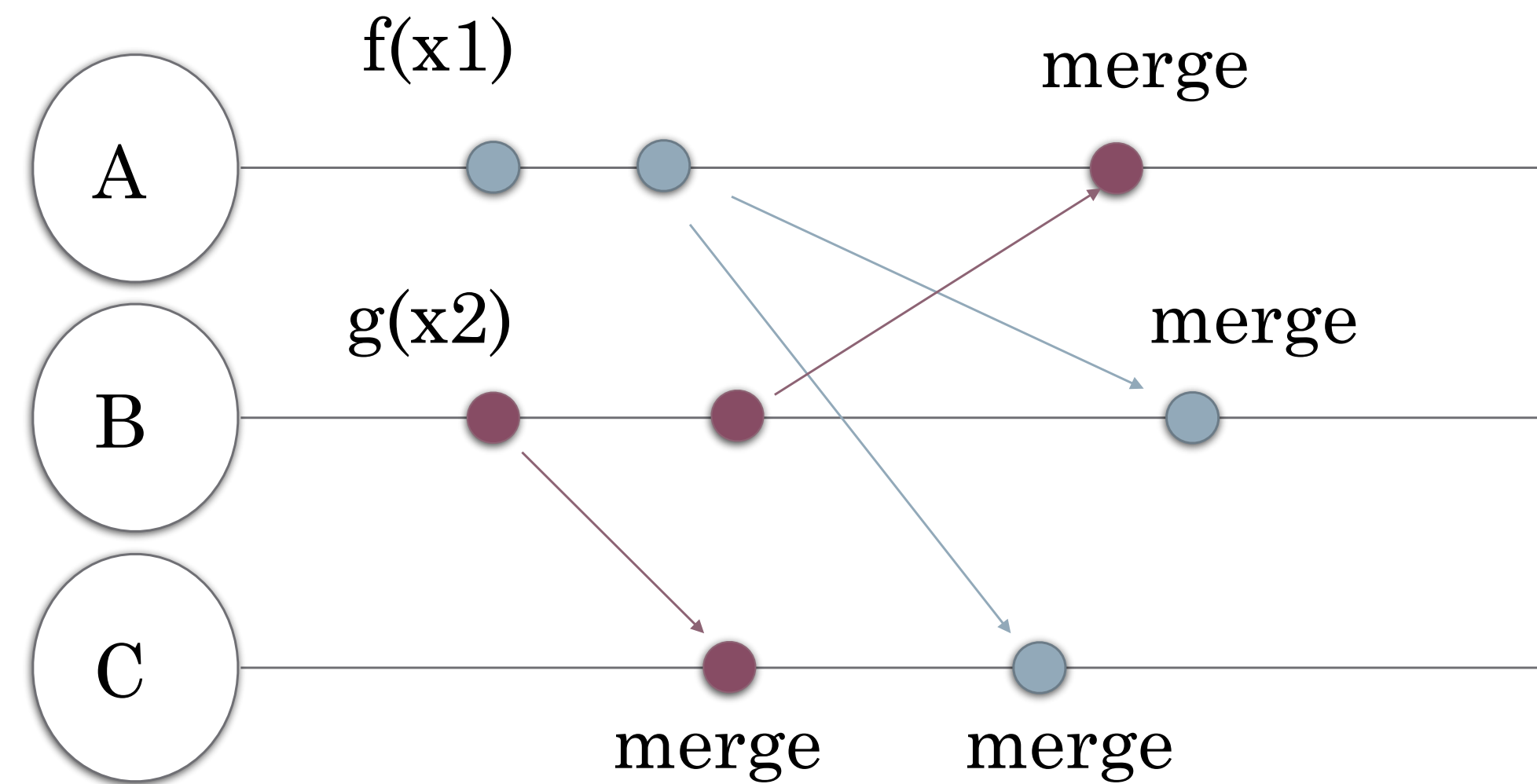
Requirements:
- + Commutativity
- + Associativity
- + Exactly once delivery
- - Idempotence

## convergent

Requirements:
- + Commutativity
- + Associativity
- + Idempotence
- - Exactly once delivery

# Convergent CRDTs

- $M(a, b) = M(b, a)$

- $M(M(a, b), c) = M(a, M(b, c))$

- $M(a, b) = M(M(a, b), b) = M(M(M(a, b), b), b)$

# Impacted Components for CRDTs

# Online Flight Check-in System

# Online Flight Check-in System



t1 — seat: 12F (220)     XDR     seat: 16D (150)

t2 — seat: 12F (220)     LWW     seat: 12F (220)

a

b

TIME

**M(a, b) for LWW = MAX(a, b)**

# Online Flight Check-in System



**t1** seat: $\{a_1:12F\}$

**XDR**

**t1** seat: $\{b_1: 16D\}$

**t2** seat: {
   $b_1$: 16D,
   $a_1$: 12F
}

**t2** seat: {
   $b_1$: 16D,
   $a_1$: 12F
}

**a**

**b**

**TIME**

# Online Flight Check-in System



**t1**

seat: $\{a_1:12F\}$

**t2**

seat: {
  $b_1$: 16D,
  $a_1$: 12F
}

**a**

**XDR**

seat: $\{b_1: 16D\}$

seat: {
  $b_1$: 16D,
  $a_1$: 12F
}

**Add-O Map**

**b**

**TIME**

# Causality

```
seat: {
  b₁: 16D,
  a₁: 12F
}
```

$$\downarrow$$

**a1;b1**

**Causality
Vector (cv)**

# Causality

seat: {
　$b_1$: 16D,
　$a_1$: 12F
}

↓

**a1;b1**

**Causality
Vector (cv)**



12F is causal to 10A - we can drop 12F

10A is causal to 5C - we can drop 10A

# Causality

seat: {
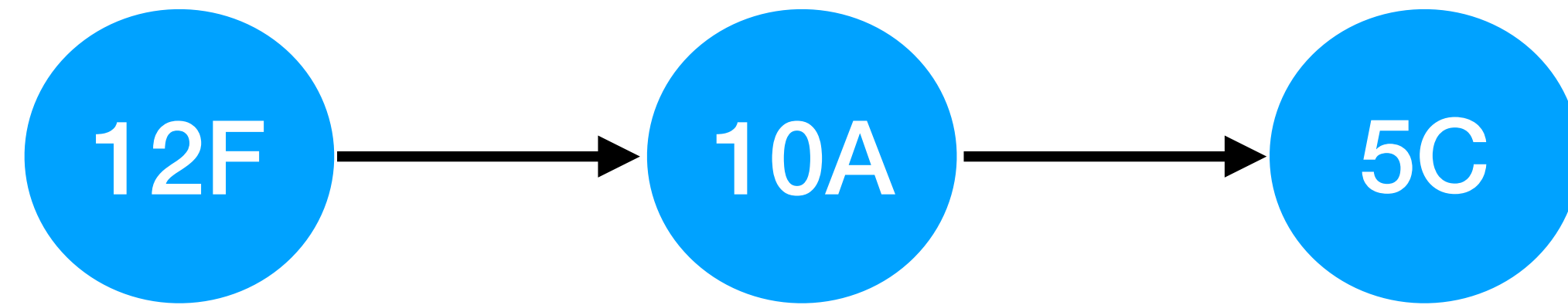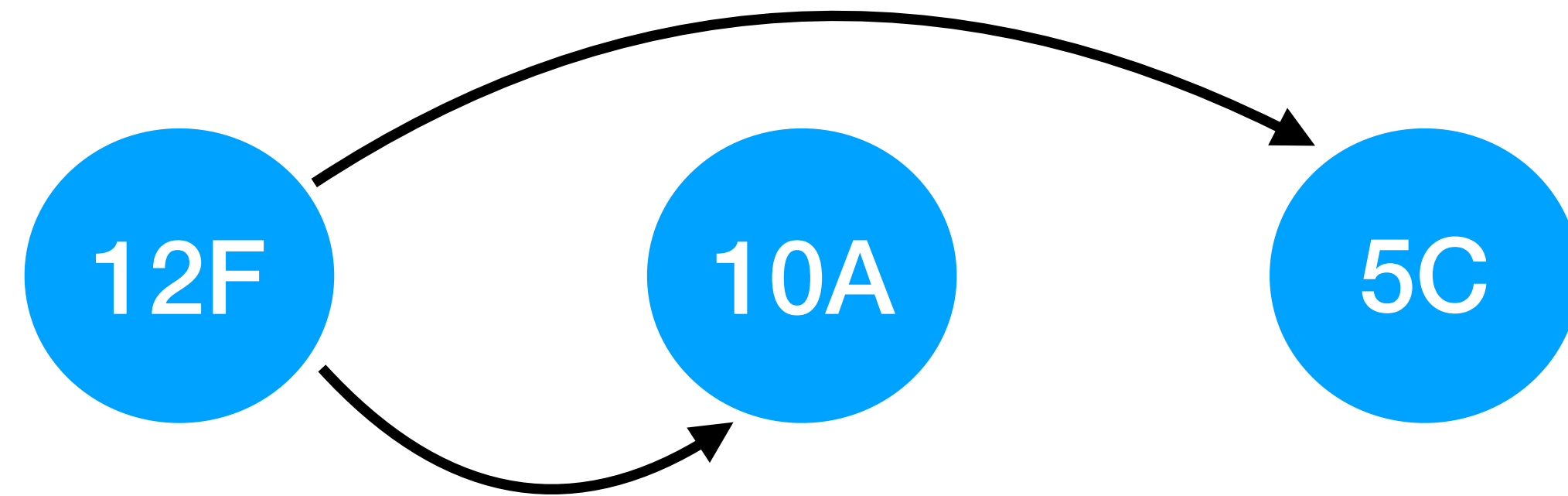  b₁: 16D,
  a₁: 12F
}

$a_1; b_1$

**Causality Vector (cv)**



12F is causal to 10A - we can drop 12F

10A is **NOT** causal to 5C - we can **NOT** drop 10A

# Causality

seat: {
<br>
  $b_1$: 16D,
<br>
  $a_1$: 12F
<br>
}

$\downarrow$

**a1;b1**

**Causality**
**Vector (cv)**

## Client Operations:

GET(key): value  =>  GET(key): (value**, cv**)

PUT(key, value)  =>  PUT(key, value**, cv**)

# Causality

seat: {
   $b_1$: (16D, **cv**),
   $a_1$: (12F, **cv**)
}

$$\downarrow$$

**a1;b1**

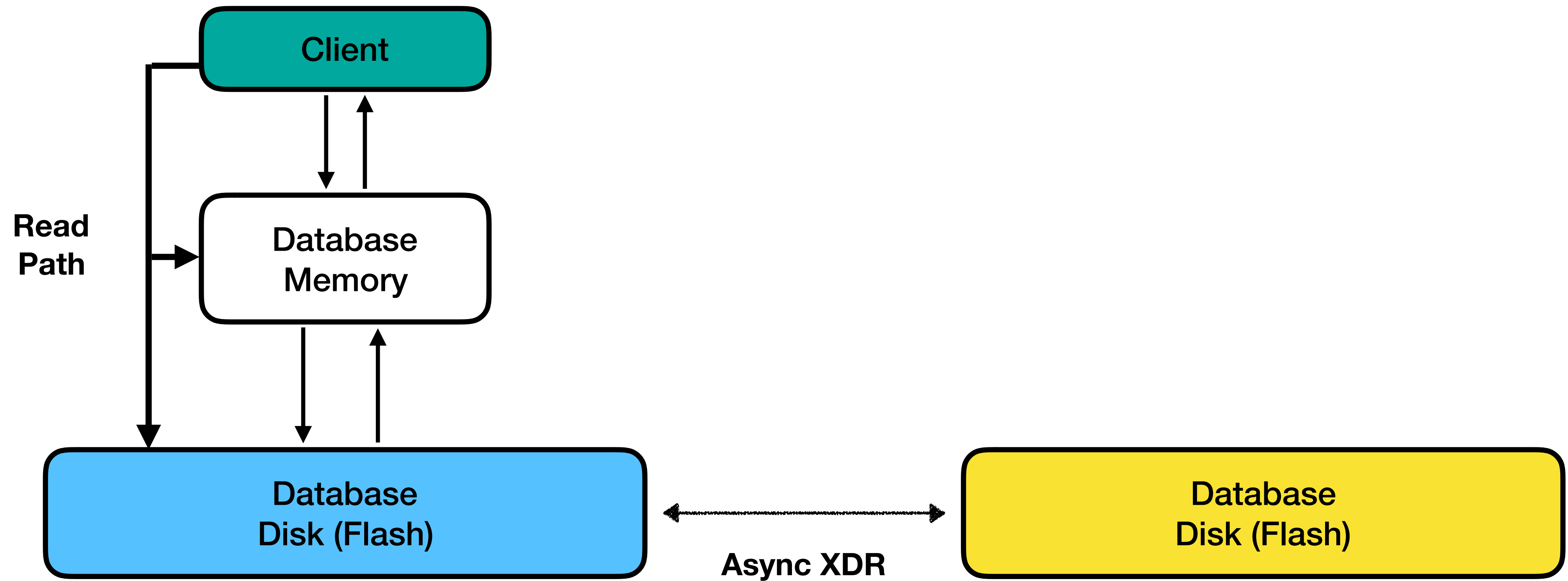**Causality
Vector (cv)**

**Client Operations:**

GET(key): value => GET(key): (value, **cv**)

PUT(key, value) => PUT(key, value, **cv**)

# Aerospike Datastore

# Aerospike Datastore



**Record**

| Key |
|---|
| Metadata |

**Bins**

| Bin1 |
|---|
| Bin2 |
| Bin3 |

**1** ## User-Defined Functions

**2**

| A | B | | Result |
|---|---|---|---|
| [  bin1: 16B,  bin2: 12A ] | [  bin2: 14C,  bin3: 10A ] | **XDR** ← → | [  bin1: 16B,  bin2: (12A or 14C),  bin3: 10A ] |

12F(_)

a1:(12F,_)

| Bins | 1 |
|------|-----------|
| a1 | (12F, _) |

**a**

| Bins | 1 |
|------|---|
| | |

**b**

10D(_)

12F(_)

a     a1:(12F,_)

b     b1:(10D,_)

| Bins | 1 | 2 |
|------|------|------|
| a1 | (12F, _) | (12F, _) |

a

| Bins | 1 | 2 |
|------|------|------|
| b1 | | (10D, _) |

b

| Bins | 1 | 2 | 3 |
|------|-----------|-----------|-----------|
| a1 | (12F, _) | (12F, _) | (12F, _) |
| b1 | | | (10D, _) |

a

| Bins | 1 | 2 | 3 |
|------|---|-----------|-----------|
| a1 | | | (12F, _) |
| b1 | | (10D, _) | (10D, _) |

b

10D(_)

[12F] -> 10F

12F(_)

[12F](a1)

10F(a1)

**a**

a1:(12F,_)

a1: (12F, _)

a1: (12F, _)
b1:(10D,_)
a2: (10F, a1)

a1: (12F, _)
b1: (10D, _)

**b**

b1:(10D,_)

| Bins | 1 | 2 | 3 | 4 |
|------|-----|-----|-----|-----|
| a1 | (12F, _) | (12F, _) | (12F, _) | (12F, _) |
| b1 | | | (10D, _) | (10D, _) |
| a2 | | | | (10F, a1) |

**a**

| Bins | 1 | 2 | 3 | 4 |
|------|-----|-----|-----|-----|
| a1 | | | (12F, _) | (12F, _) |
| b1 | | (10D, _) | (10D, _) | (10D, _) |

**b**

| Bins | 1 | 2 | 3 | 4 | 5 | 6 |
|------|---|---|---|---|---|---|
| a1 | (12F, _) | (12F, _) | (12F, _) | (12F, _) | (12F, _) | (12F, _) |
| b1 | | | (10D, _) | (10D, _) | (10D, _) | (10D, _) |
| a2 | | | | (10F, a1) | (10F, a1) | (10F, a1) |
| a3 | | | | | (5C, a2b1) | (5C, a2b1) |

**a**

| Bins | 1 | 2 | 3 | 4 | 5 | 6 |
|------|---|---|---|---|---|---|
| a1 | | | (12F, _) | (12F, _) | (12F, _) | (12F, _) |
| b1 | | (10D, _) | (10D, _) | (10D, _) | (10D, _) | (10D, _) |
| a3 | | | | | | (5C, a2b1) |

**b**

# Learnings

- CRDTs allowed us to achieve convergent **predictable** state of our data

# Learnings

- CRDTs allowed us to achieve convergent **predictable** state of our data

- Education about right trade-off between **Consistency** and **Correctness**
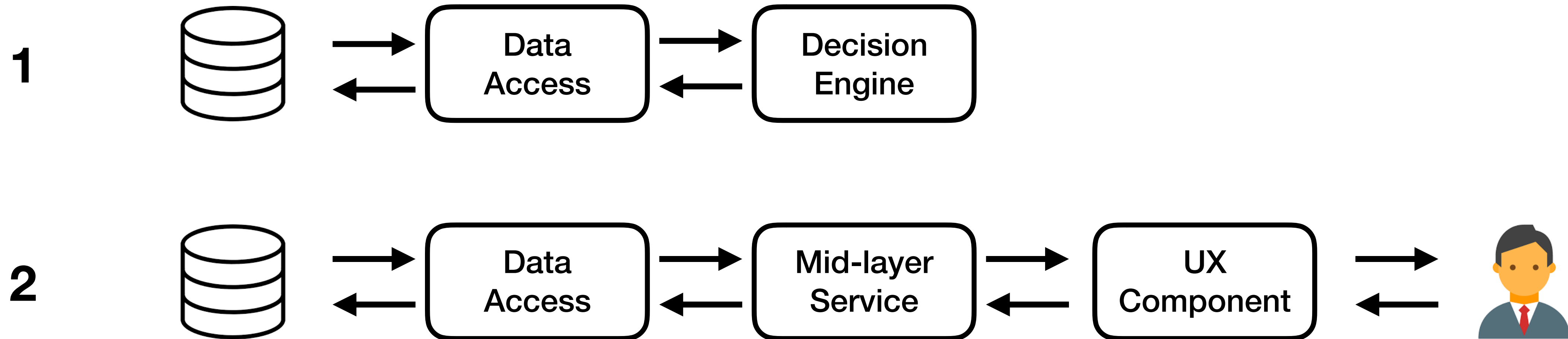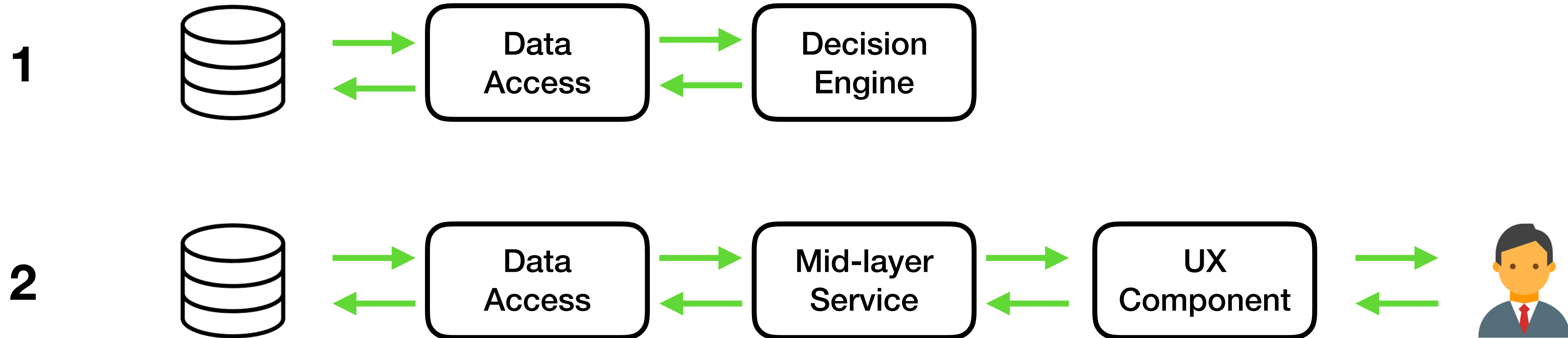
# Learnings

- CRDTs allowed us to achieve convergent **predictable** state of our data

- Education about right trade-off between **Consistency** and **Correctness**

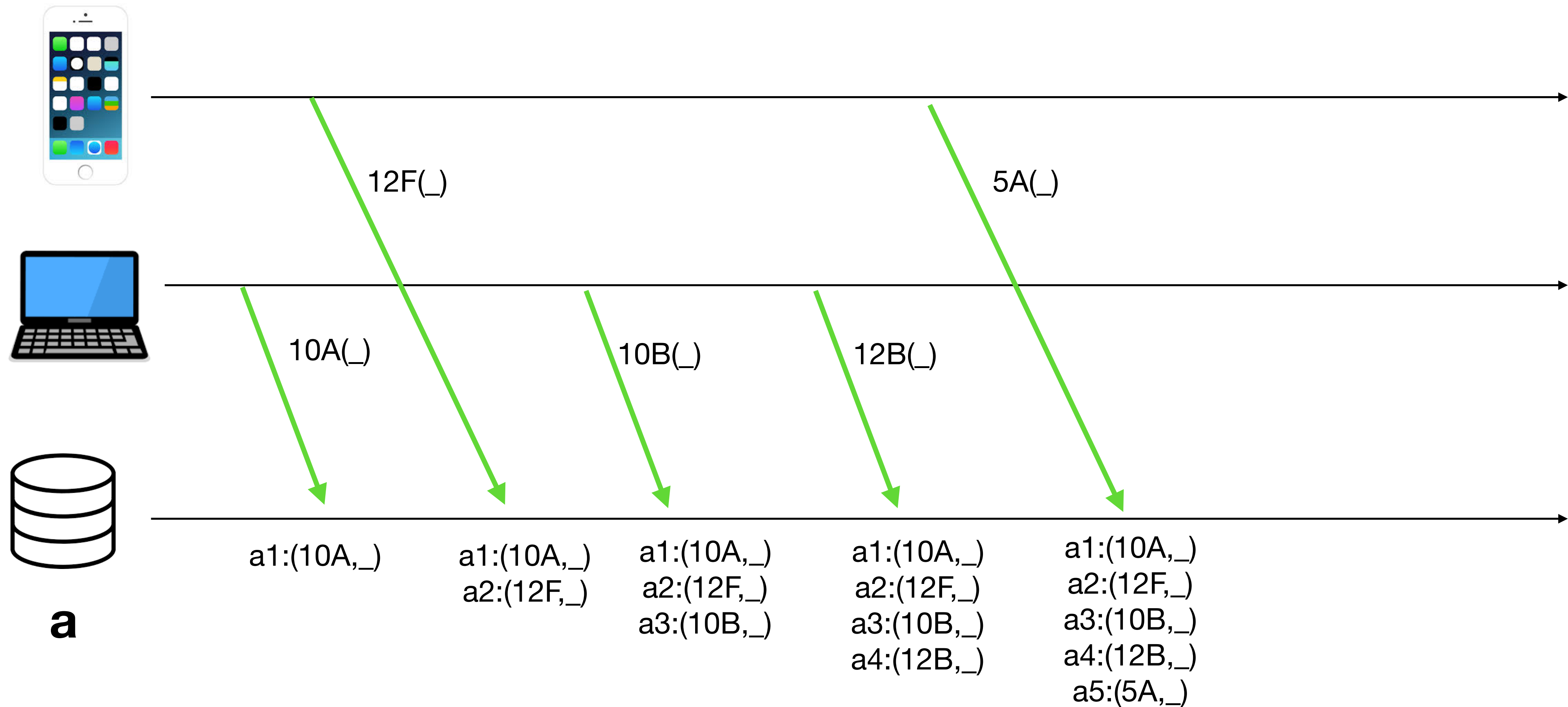- Do not **underestimate** concurrent data access

# Caveat #1: CV Propagation

# Caveat #1: CV Propagation
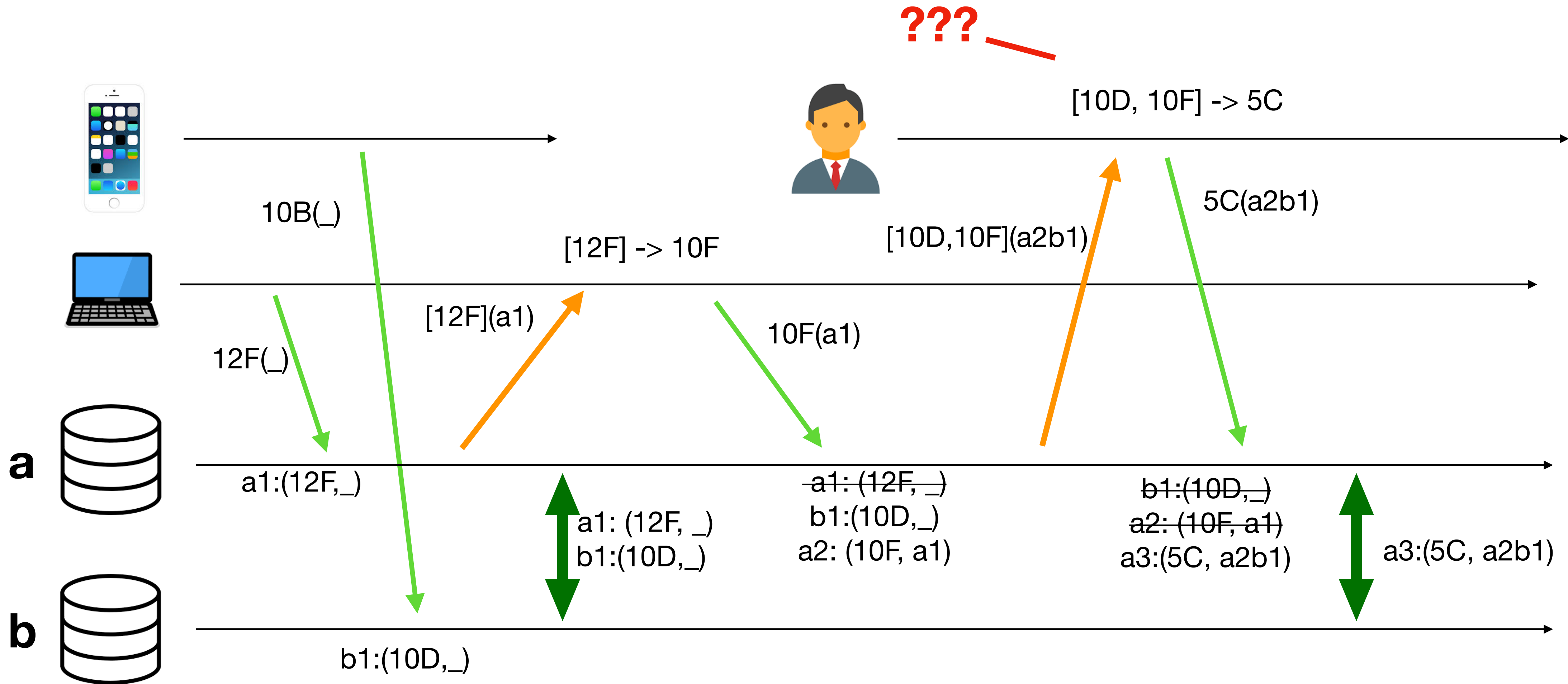
# Caveat #2: Siblings Explosion

12F(_)

5A(_)

10A(_)

10B(_)

12B(_)

**a**

a1:(10A,_)

a1:(10A,_)
a2:(12F,_)

a1:(10A,_)
a2:(12F,_)
a3:(10B,_)

a1:(10A,_)
a2:(12F,_)
a3:(10B,_)
a4:(12B,_)

a1:(10A,_)
a2:(12F,_)
a3:(10B,_)
a4:(12B,_)
a5:(5A,_)

# Caveat #3: Wait, Siblings ?

**???**

[10D, 10F] -> 5C

10B(_)

[12F] -> 10F

[10D,10F](a2b1)

5C(a2b1)

[12F](a1)

10F(a1)

12F(_)

**a**

a1:(12F,_)

a1: (12F, _)
b1:(10D,_)

a1: (12F, _)
b1:(10D,_)
a2: (10F, a1)

b1:(10D,_)
a2: (10F, a1)
a3:(5C, a2b1)

a3:(5C, a2b1)

**b**

b1:(10D,_)

# Thanks!