



C8DB

Geo-Replicated, Conflict-Free Document Database

Christopher S. Meiklejohn

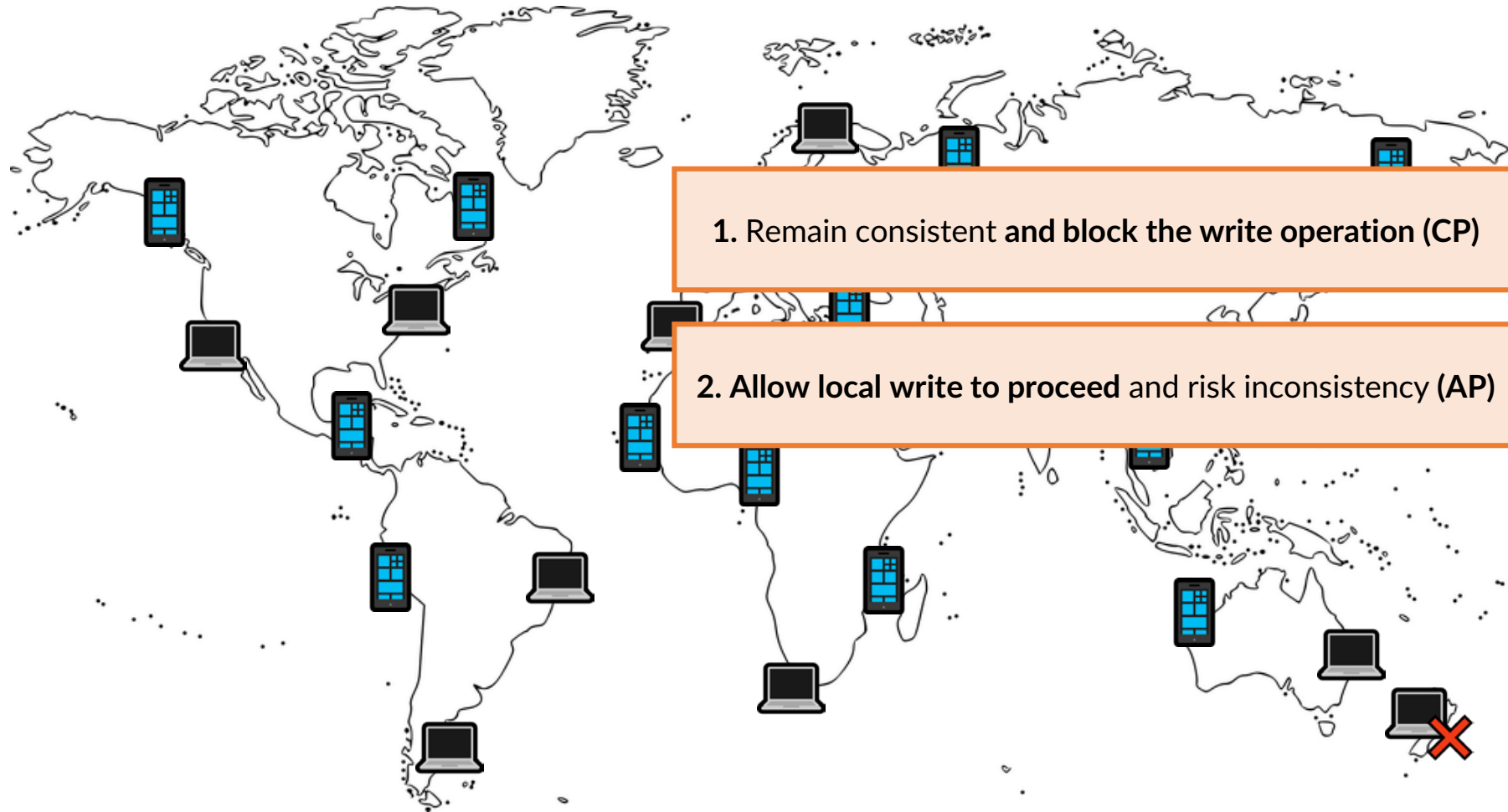
Macrometa

HPTS 2019

November 5th, 2019



Geo-Distribution Challenges



Conflict-Free Replicated Data Types

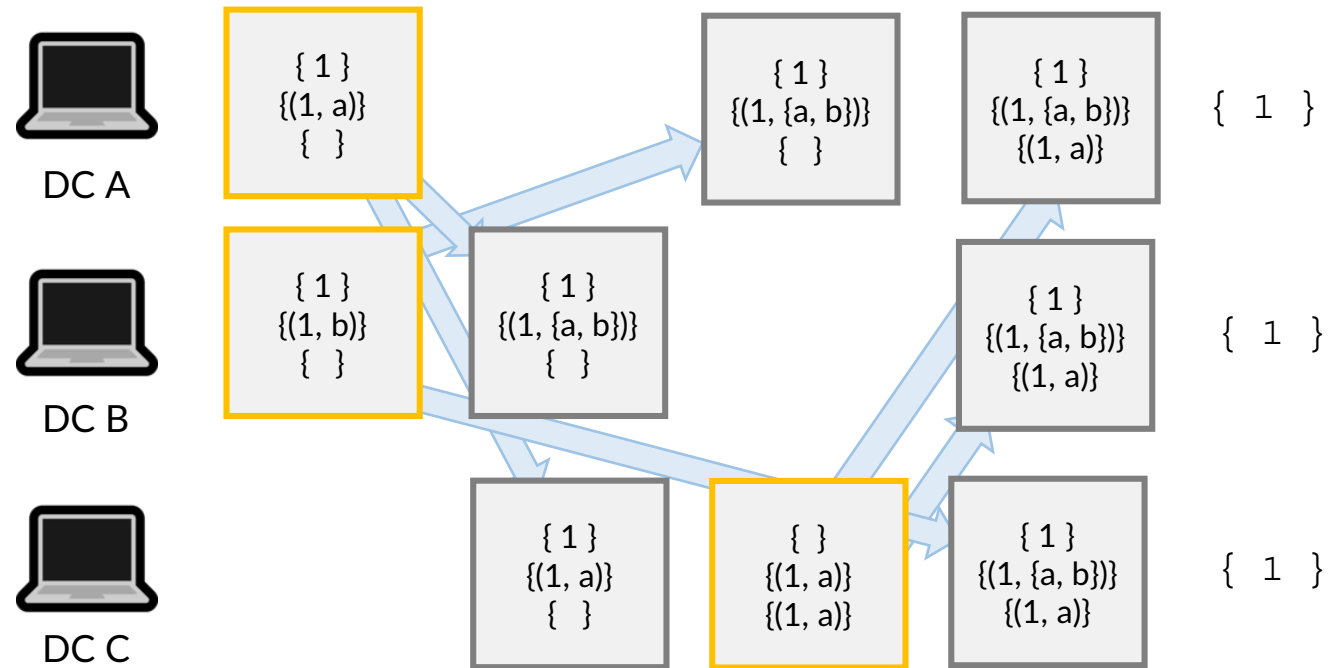
Distributed data structures, mimic sequential counterparts

Designed for convergence under concurrency

Multiple flavors (e.g., state-based, delta-based)

What do we do when things don't commute?

Counters G-Counter PN-Counter	Registers LWW-Register MV-Register
Sets G-Set 2P-Set OR-Set AW-Set RW-Set	Flags Sequence Treedoc RGA RGA-Split



CRDTs in Production



Riak distributed database

First major production implementation of state-based CRDTs
Used by NHS, Riot Games, LO/JACK, bet365



AntidoteDB (HPTS '17)

Developer on operation-based CRDT database with Transactional Causal Consistency
Part of the LightKone, SyncFree research projects

Others

Mesosphere's DC/OS (Lashup, Minuteman), Adobe, Comcast, Macrometa
TomTom, Roshi, Automerge, Ditto, Cosmos DB, etc.



Macrometa: Overview



Document database for the edge

Documents are **JSON** and grouped into **collections**

(think: namespaces for each document, clients presented with dictionaries)



105 physical regions, 175 logical regions

36 **reserved core** regions, 12 available through public beta today

34 **customer requested** regions (on-prem, specific colo, etc.)

35 **burstable** regions available



Configurable consistency

Strong consistency in a single DC with propagation to remote DCs in order

Eventual consistency in an multi-master mode with all DCs receiving updates



Macrometa: Architecture



Logging-based Architecture

Updates within a data-center are **durably written to a log**
Logs are maintained with **Zookeeper** and **BookKeeper**



Causal broadcast

There is a single log per destination DC delivered in **FIFO order**
Updates are buffered until log entries **dependencies are met**

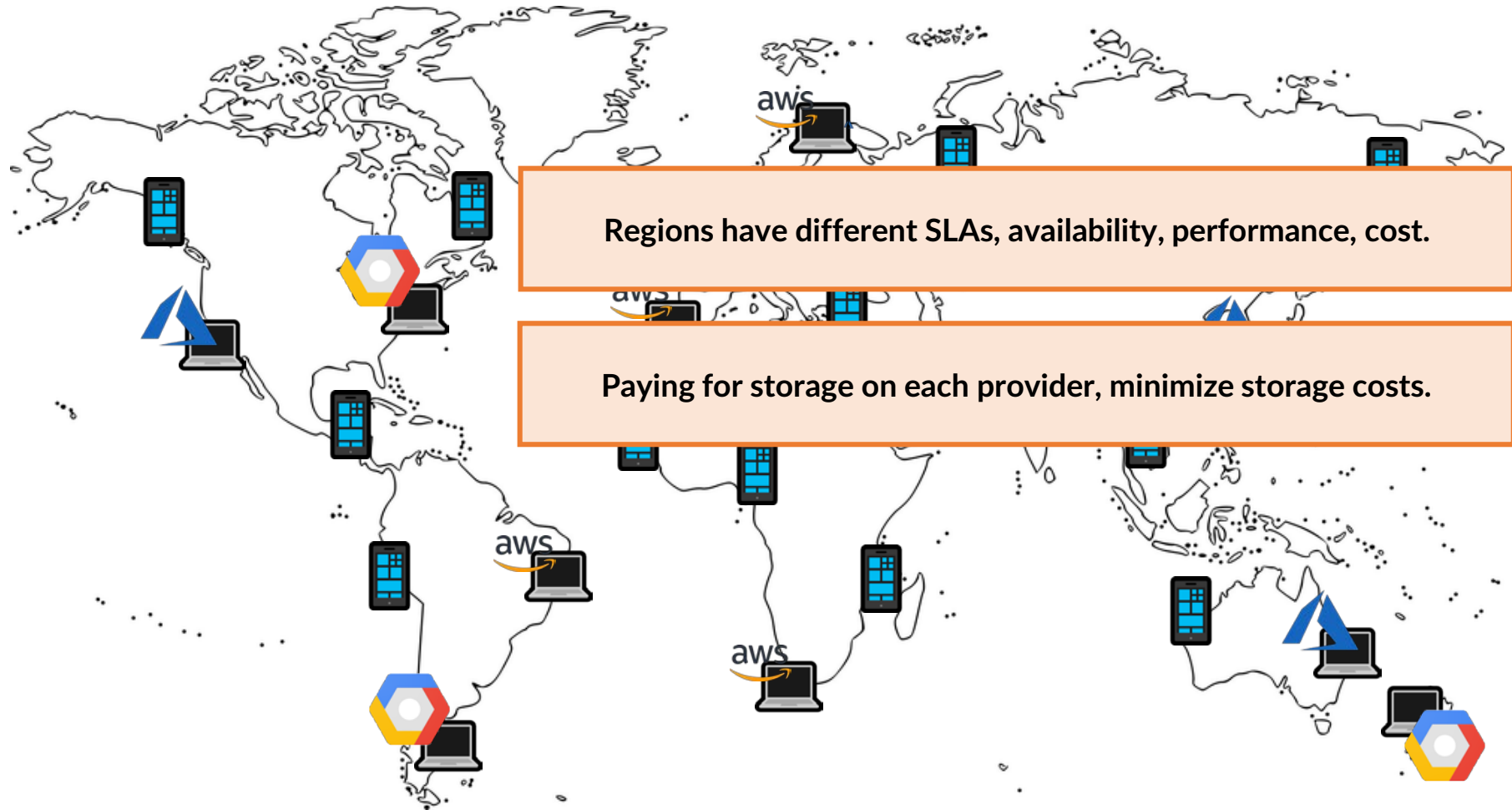


Client operate on the **document** abstraction level

Clients submit documents that are **converted to updates** by Macrometa nodes
Macrometa nodes **materialize values** for keys based on incoming log updates



Geo-Distribution Concerns



Registers



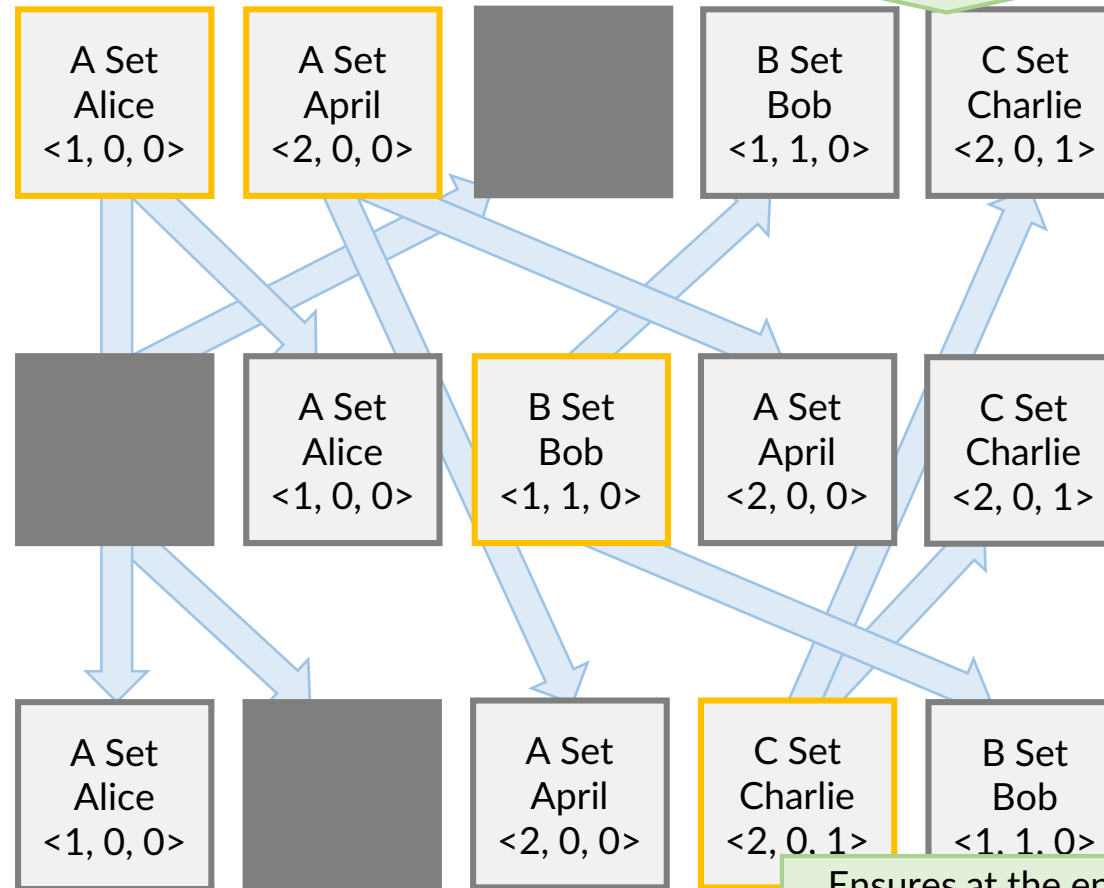
DC A



DC B



DC C



Concurrent updates merge arbitrate based on the maximum DC identifier.

Charlie

Charlie

Charlie

Ensures at the end of convergence we need to only keep a single update in the log for each field in the document.



Registers: Read Anomalies



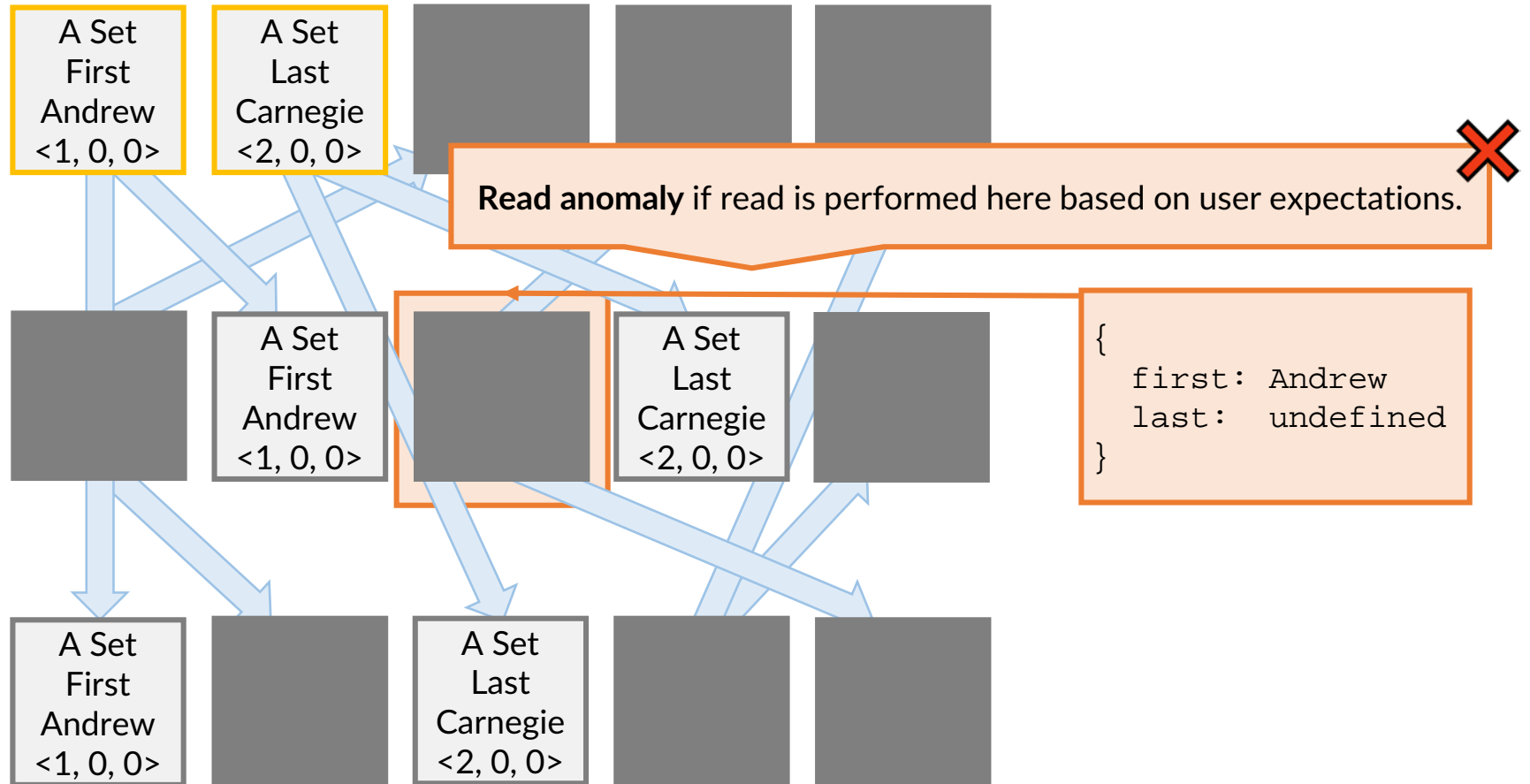
DC A



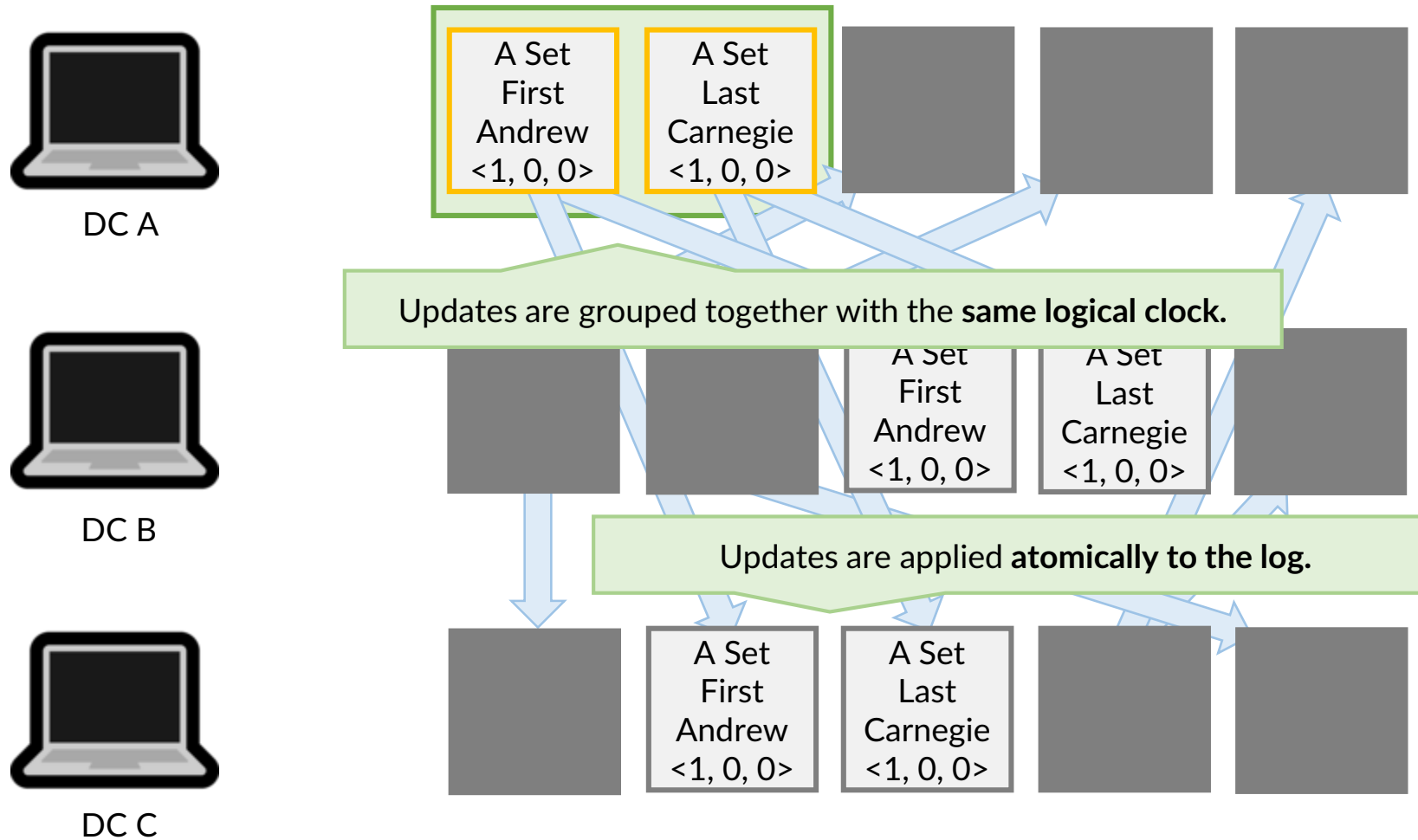
DC B



DC C



Transactions v1



Transactions v1: Concurrency

Concurrent operations are merged together using the existing convergence strategy by **maximum DC identifier**.



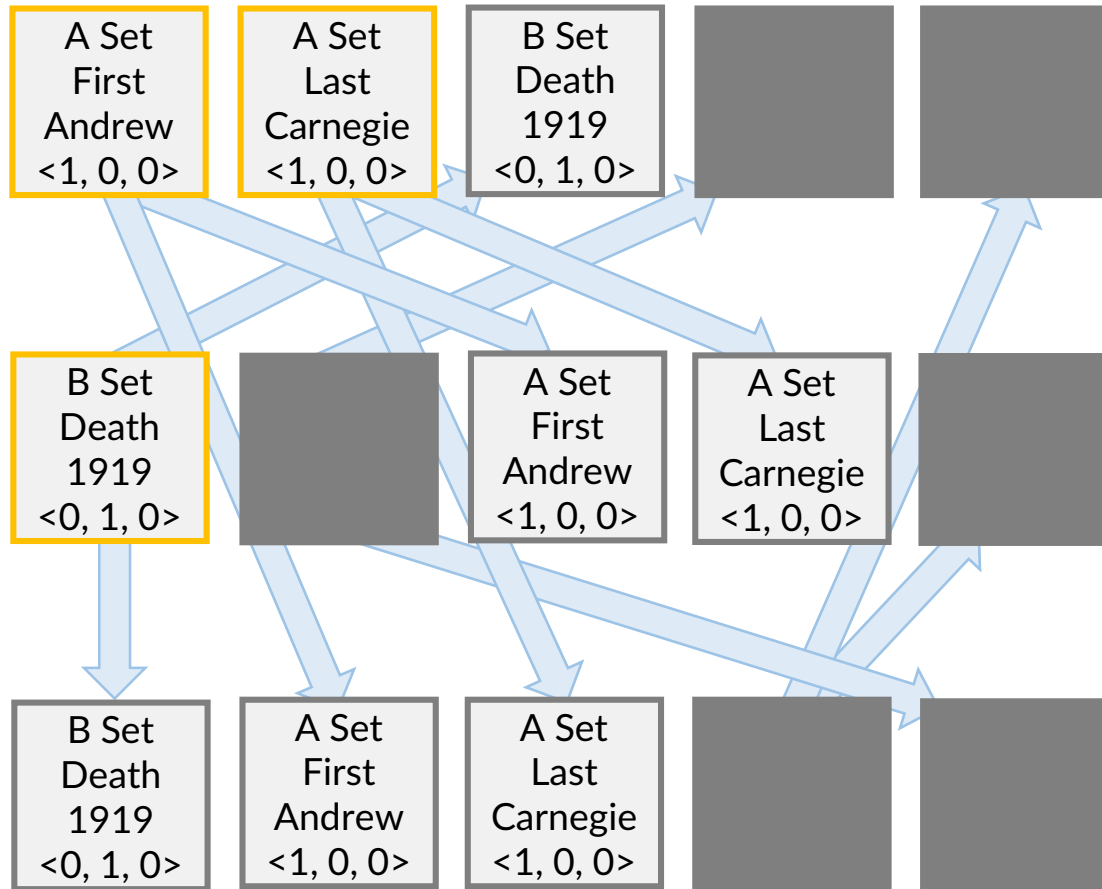
DC A



DC B



DC C



```
{
  first: Andrew
  last: Carnegie
  death: 1919
}
```

```
{
  first: Andrew
  last: Carnegie
  death: 1919
}
```

```
{
  first: Andrew
  last: Carnegie
  death: 1919
}
```



Transactions v1: Concurrency Anomalies



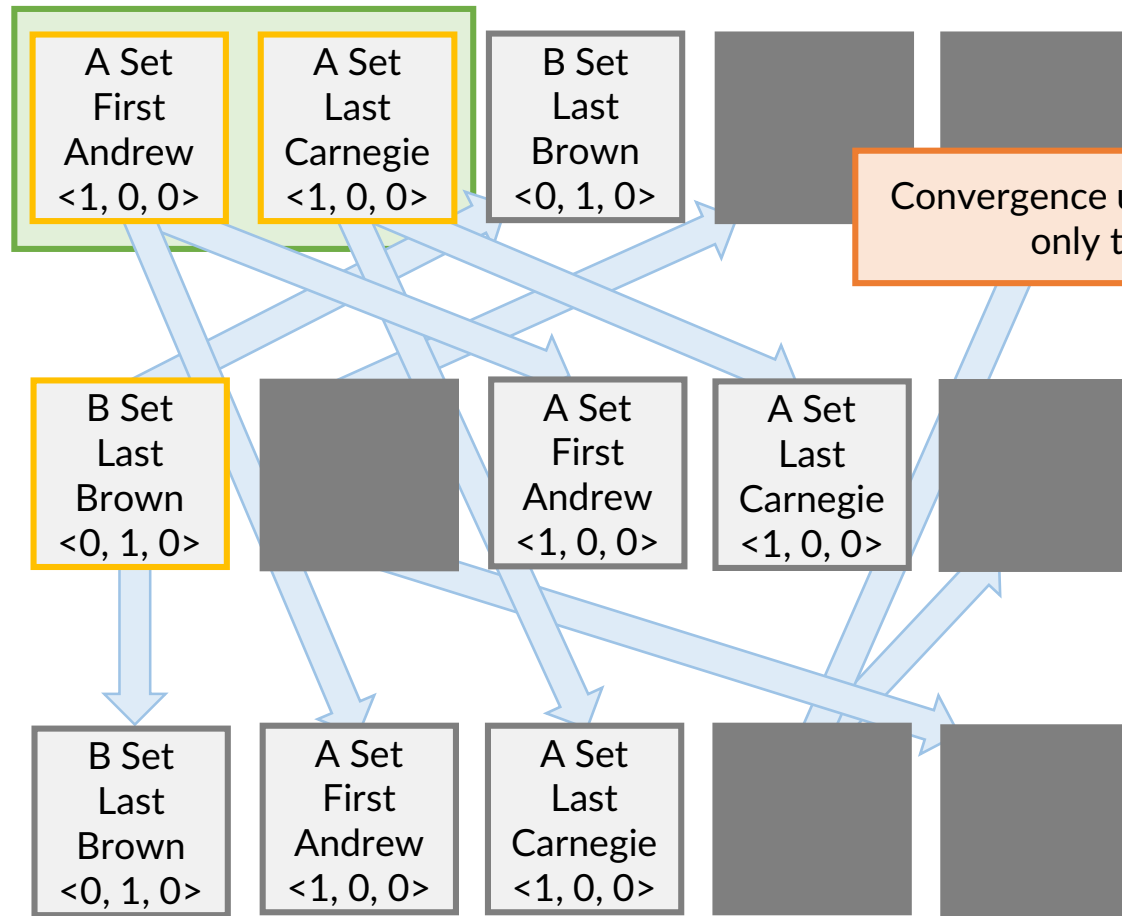
DC A



DC B



DC C

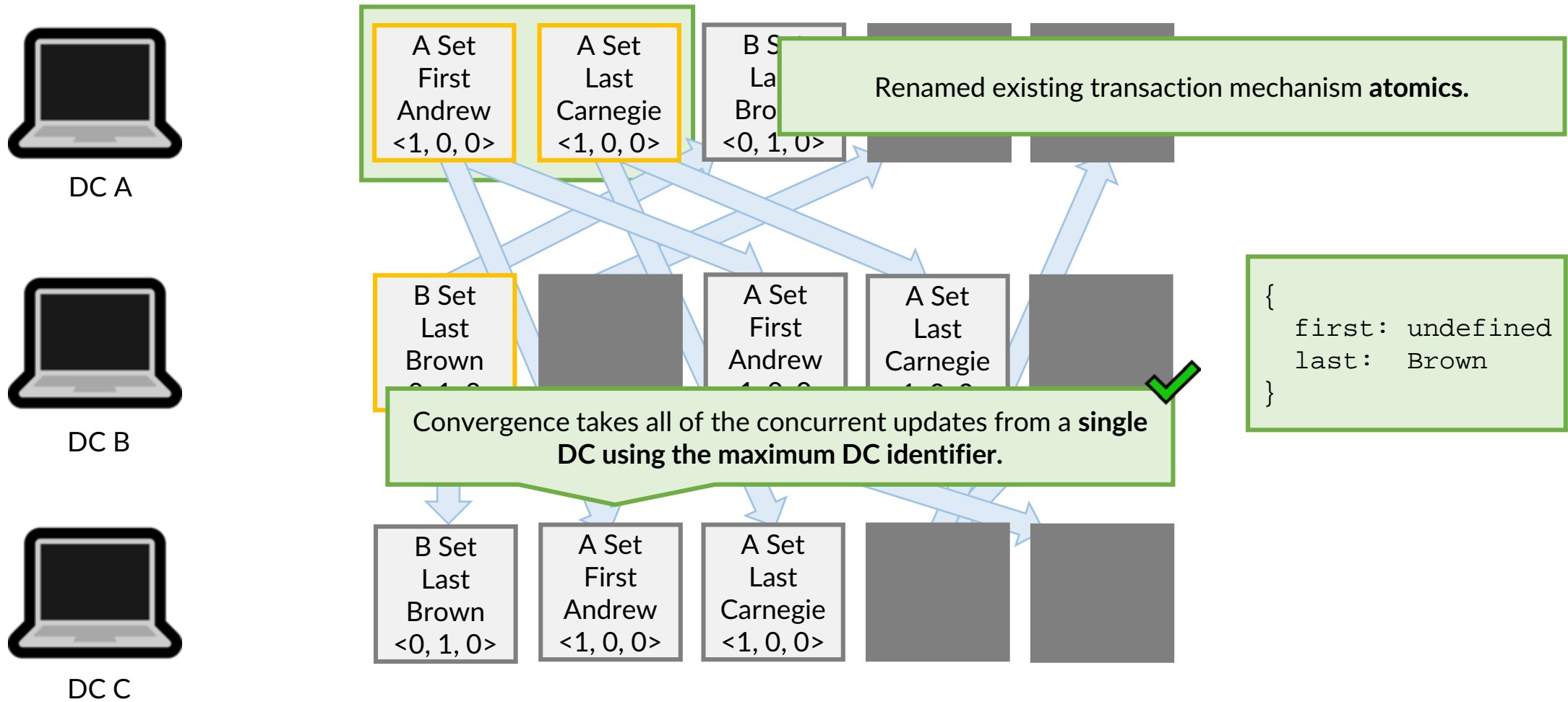


Convergence uses **maximum DC identifier** and only takes part of the update.


```
{  
  first: Andrew  
  last:  Brown  
}
```



Transactions: v2







Counters?

 **Lindsey Kuper**
@lindsey

"Oh, you wanted to *increment a counter*?! Good luck with that!" -- the distributed systems literature

11:55 AM · Mar 9, 2015 · [Twitter Web Client](#)

367 Retweets **557** Likes

Can't use a register, **assignment doesn't commute.**

Sum double counts, **max** under counts.

Store sum per node or...

...store individual updates.



Types of CRDT Counters



State-based CRDTs

Lattice-based, lower for storage overhead and require that developers ensure that all operations are always mergeable by proper design.

e.g., Riak

How do we find a solution that is **compatible** with the **transactions** that meets both the **garbage collection** and **storage** requirements?



Operation-based CRDTs

Optimized for individual operations where all operations are designed to be commutative, but do not need for form a lattice. Also requires causal broadcast.

e.g. Antidote



State-based



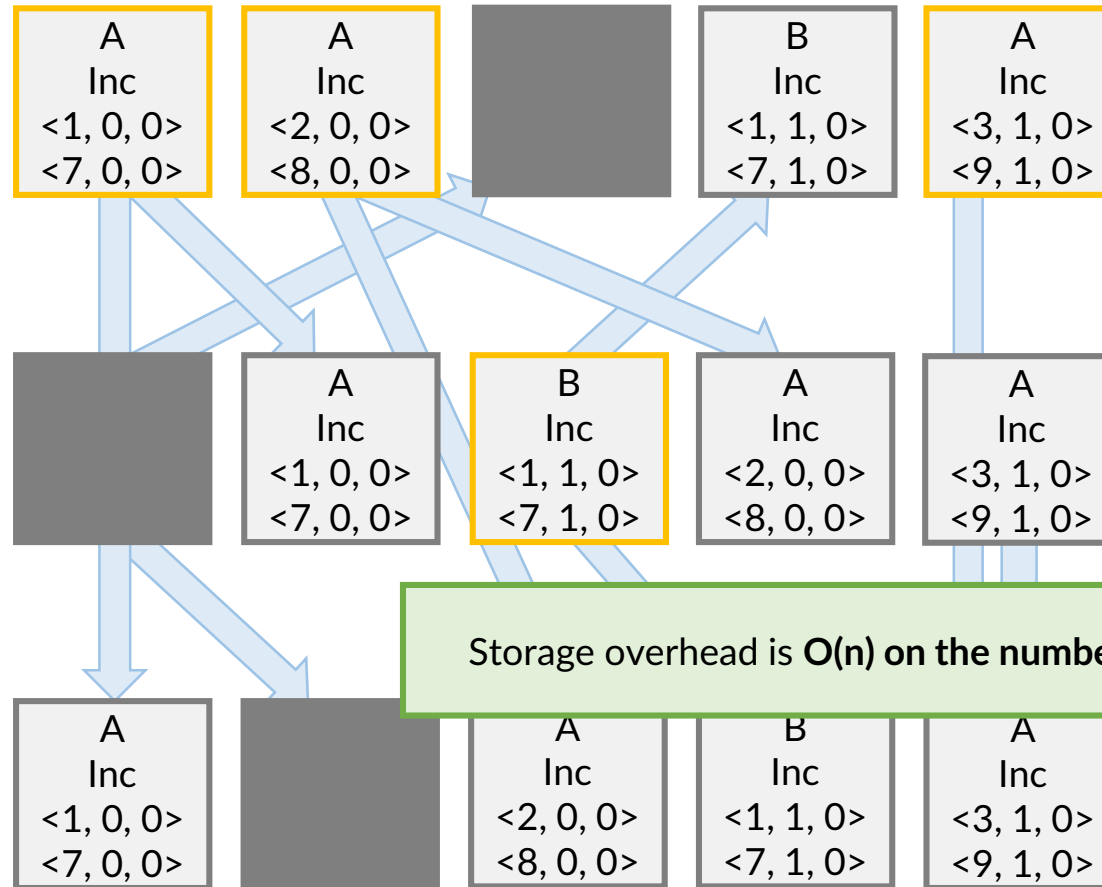
DC A



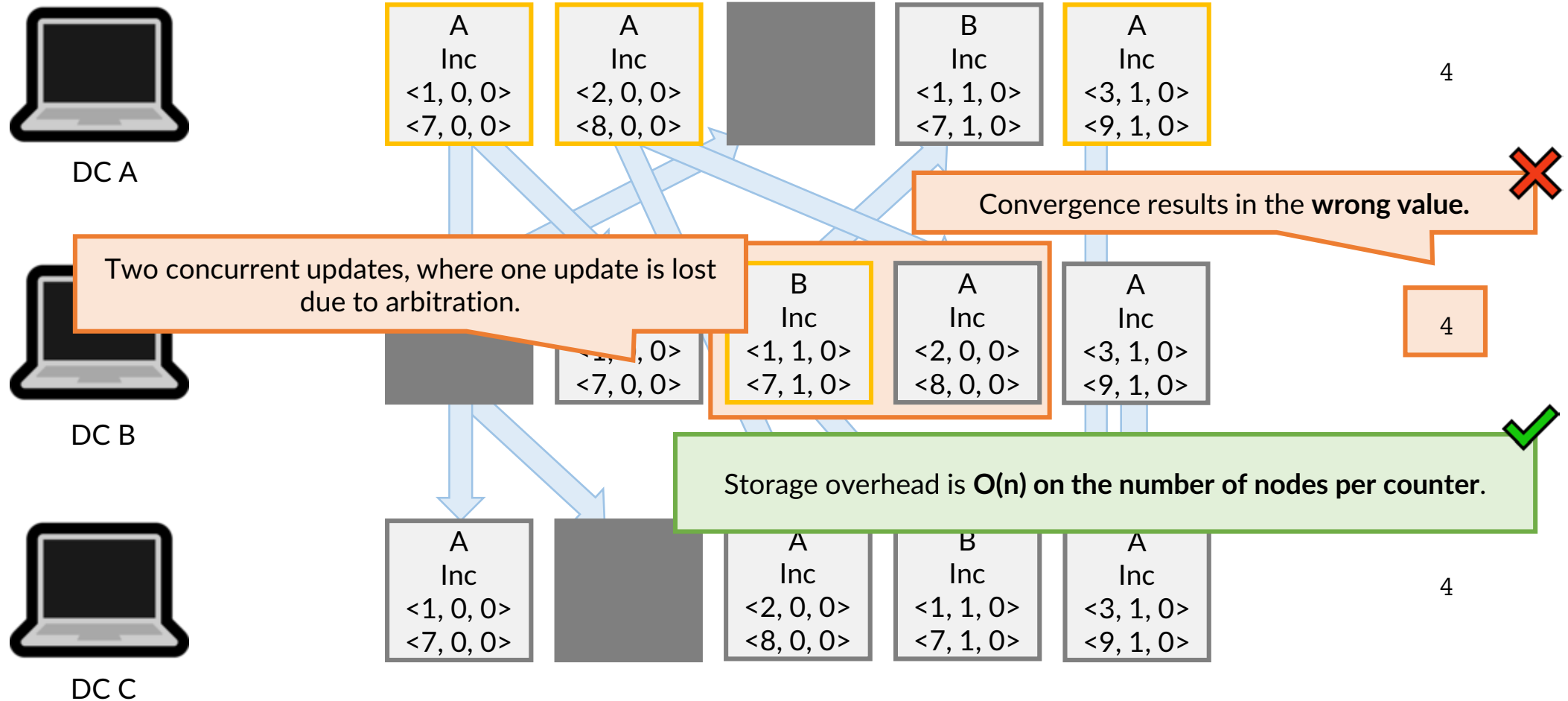
DC B



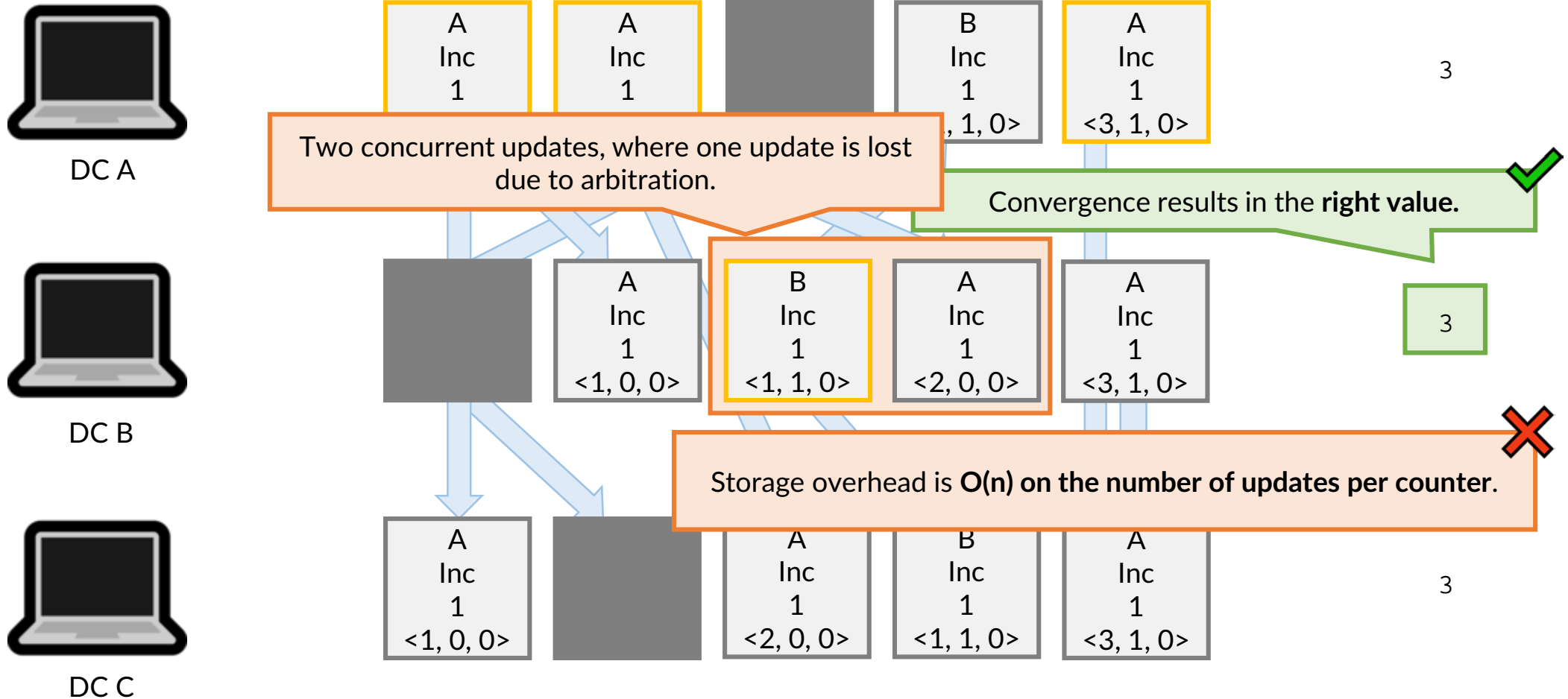
DC C



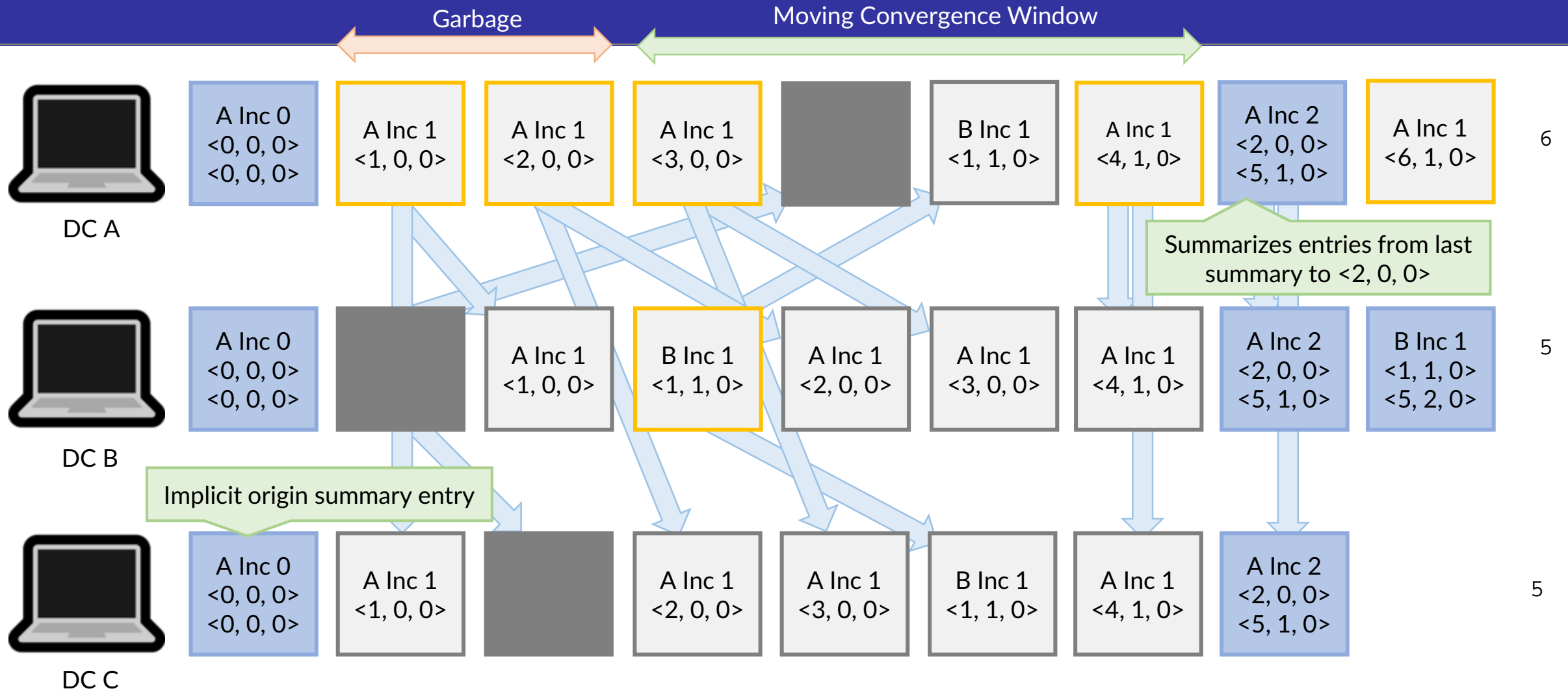
State-based and Transactions



Operation-based and Transactions



Counters



CRDTs, tho?



CRDTs?

Are these **actually** CRDTs? Probably not, need a new name.
Interesting point in the design space (c.f., Riak, Antidote)



Thanks!

