

# High Availability under Eventual Consistency

## Latency Magnitudes

- $\lambda$ , up to 50ms (local region DC)
- $\Lambda$ , between 100ms and 300ms (inter-continental)

### No inter-DC replication

Client writes observe  $\lambda$  latency

### Planet-wide geo-replication

Replication Techniques	Latency Ranges	Description
Consensus/Paxos	$[\Lambda, 2\Lambda]$	with no divergence
Primary-Backup	$[\lambda, \Lambda]$	asynchronous/lazy
Multi-Master	$\lambda$	allowing divergence

## EC and CAP for Geo-Replication

### Eventually Consistent

- In an ideal world there would be only one consistency model: when an update is made all observers would see that update.
- Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.

### High Availability

- Special case of weak consistency
- After an update, if no new updates are made to the object, eventually all reads will return the same value, that reflects the last update. E.g: DNS.

This can later be reformulated to avoid quiescence, by adapting a session guarantee.

### CAP Theorem

Of three properties of shared-data systems – **data consistency, system availability, and tolerance to network partition** – only two can be achieved at any given time.

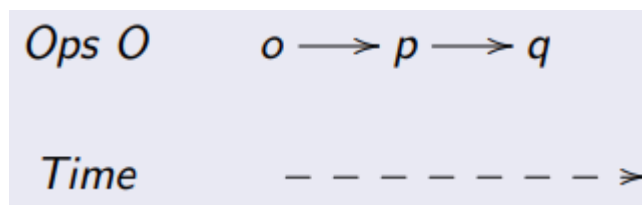
# Session Guarantees

- **Read Your Writes** – read operations reflect previous writes
- **Monotonic Reads** – successive reads reflect a non-decreasing set of writes
- **Writes Follow Reads** – writes are propagated after reads on which they depend.
- **Monotonic Writes** – writes are propagated after writes that logically precede them.

## From Sequential to Concurrent Executions

- Consensus provides illusion of a single replica
- This also preserves (slow) sequential behaviour
- EC Multi-master (or active-active) can expose concurrency

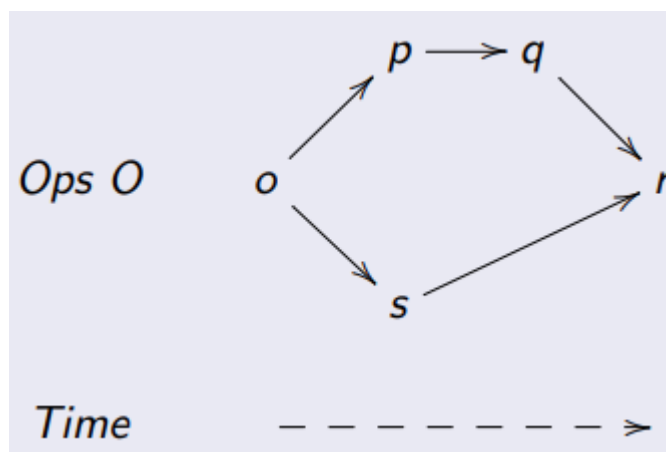
### Sequential Execution



Ordered set  $(O, <)$ .

- $O = \{o, p, q\}$  and  $o < p < q$

### Concurrent Execution



Partially ordered set  $(O, <)$ .  $o < p < q < r$  and  $o < s < r$

Some ops in  $O$  are concurrent:  $p \parallel s$  and  $q \parallel s$

## Conflict-Free Replicated Data Types (CRDTs)

- Convergence after concurrent updates  $\rightarrow$  favor AP under CAP
  - Examples include counters, sets, mv-registers, maps, graphs
- Operation based CRDTs  $\rightarrow$  operation effects must commute

- State based CRDTs are rooted on join semi-lattices

## Operation-based CRDTs, Effect Commutativity

- In some datatypes, all operations are commutative

## State-based CRDTs, Join semi-lattices

$S \rightarrow$  (partial ordered set)

$U \rightarrow$  join, deriving least upper bounds

$\perp \rightarrow$  initial state (usually the least element)

- Join properties in a semilattice  $(S, \leq, U)$ :
  - **Idempotence**,  $a \cup a = a$ ,
  - **Commutativity**,  $a \cup b = b \cup a$ ,
  - **Associative**,  $(a \cup b) \cup c = a \cup (b \cup c)$ .

## Eventual Consistency, non stop

$\text{upds}(a) \subseteq \text{upds}(b) \Rightarrow a \leq b$ .

- This is slightly **weaker** than the previous definition and implies it:  $\text{upds}(a) = \text{upds}(b) \Rightarrow a = b$ .

## Design of Conflict-Free Replicated Data Types

- A partially ordered log (polog) of operations implements any CRDT
- Replicas keep increasing local views of an evolving distributed polog
- Any query, at replica  $i$ , can be expressed from local polog  $O_i$
- CRDTs are efficient representations that follow some general rules

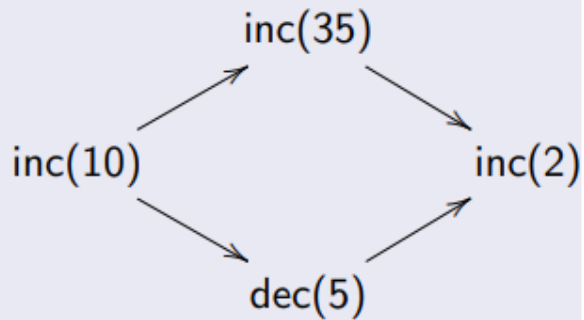
## Principle of permutation equivalence

If operations in sequence can commute, preserving a given result, then under concurrency they should preserve the same result

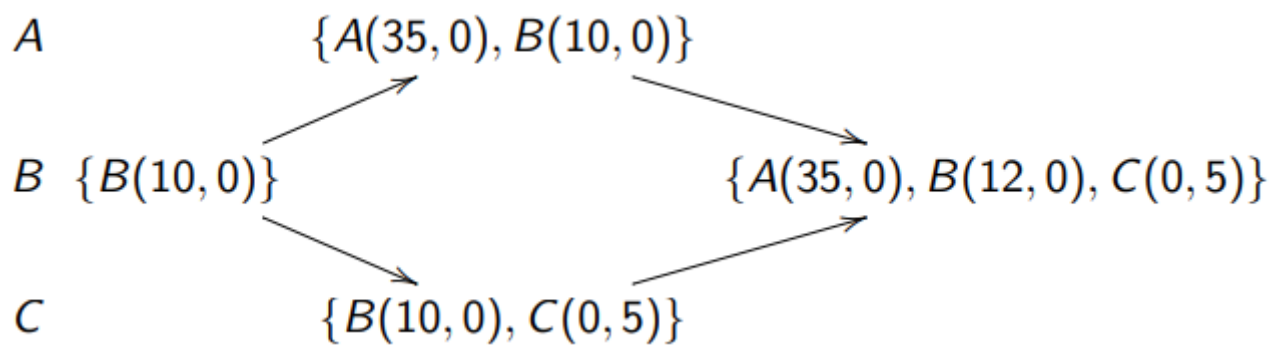
## Sequential

$\text{inc}(10) \longrightarrow \text{inc}(35) \longrightarrow \text{dec}(5) \longrightarrow \text{inc}(2)$   
 $\text{dec}(5) \longrightarrow \text{inc}(2) \longrightarrow \text{inc}(10) \longrightarrow \text{inc}(35)$

## Concurrent



## Implementing Counters



At any time, counter value is sum of incs minus sum of decs

## Registers

Ordered set of write operations

Register value is x, the last written value

## Sequential execution

$A \quad wr(x) \longrightarrow wr(j) \longrightarrow wr(k) \longrightarrow wr(x)$

## Sequential execution under distribution

$A \quad wr(x)$   
 $B \quad wr(j) \longrightarrow wr(k)$   
 $wr(x)$

## Implementing Registers

CRDT register implemented by attaching local wall-clock times

## Sequential execution under distribution

$A \quad (11:00)x$   
 $B \quad (12:02)j \longrightarrow (12:05)k$   
 $(11:30)?$   
 $?$

Problem: Wall-clock on B is one hour ahead of A

Value  $x$  might not be writeable again at A since  $12:05 > 11:30$

- Concurrent semantics should preserve the sequential semantics
  - This also ensures correct sequential execution under distribution

## Multi-value Registers

### Concurrent execution

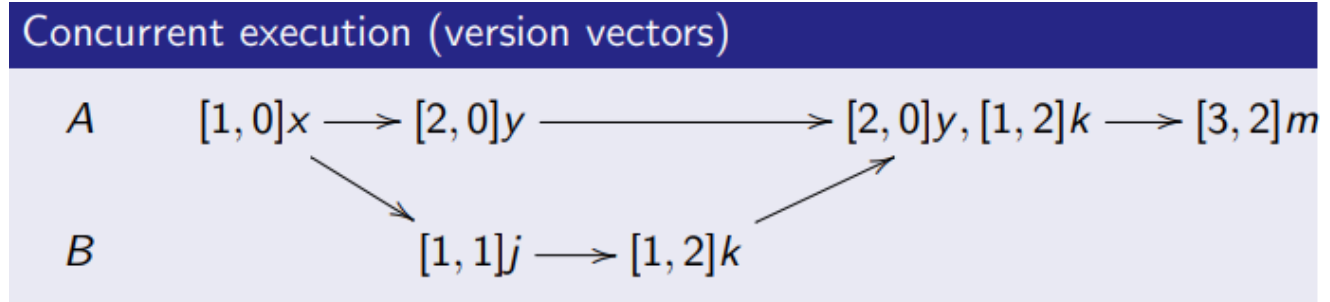
$A \quad wr(x) \longrightarrow wr(y) \longrightarrow \{y, k\} \longrightarrow wr(m) \longrightarrow \{m\}$   
 $B \quad wr(j) \longrightarrow wr(k)$

## Implementation

Concurrency can be precisely tracked with version vectors

Metadata can be compressed with a common causal context and a single scalar per

value (dotted version vectors)



## Registers in Redis

- Multi-value registers allows executions leading to concurrent values
- Presenting concurrent values is at odds with the sequential API
- Redis both tracks causality and registers wall-clock times
- Querying uses Last-Writer-Wins selection among concurrent values
- This preserves correctness of sequential semantics
- A value with clock 12:05 can still be causally overwritten at 11:30

## State-based CRDTs: G-Set

$$\begin{aligned}
 \Sigma &= \mathcal{P}(V) \\
 \sigma_i^0 &= \{\} \\
 \text{apply}_i((\text{add}, v), s) &= s \cup \{v\} \\
 \text{eval}_i(\text{rd}, s) &= s \\
 \text{merge}_i(s, s') &= s \cup s'
 \end{aligned}$$

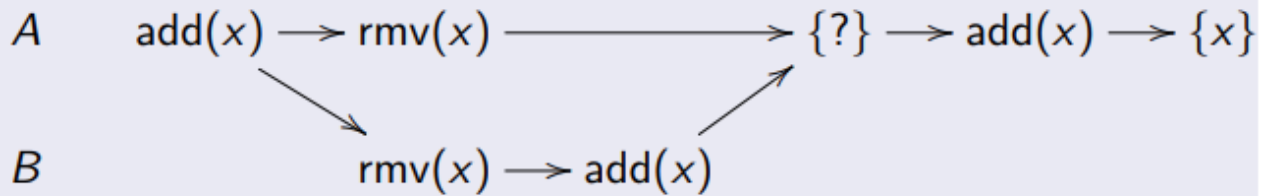
## State-based CRDTs: 2P-Set

$$\begin{aligned}
 \Sigma &= \mathcal{P}(V) \times \mathcal{P}(V) \\
 \sigma_i^0 &= \{\}, \{\} \\
 \text{apply}_i((\text{add}, v), (s, t)) &= s \cup \{v\}, t \\
 \text{apply}_i((\text{rmv}, v), (s, t)) &= s, t \cup \{v\} \\
 \text{eval}_i(\text{rd}, s) &= s \setminus t \\
 \text{merge}_i((s, t), (s', t')) &= s \cup s', t \cup t'
 \end{aligned}$$

## Sets

Problem: Concurrently adding and removing the same element

### Concurrent execution



## Add-Wins Sets

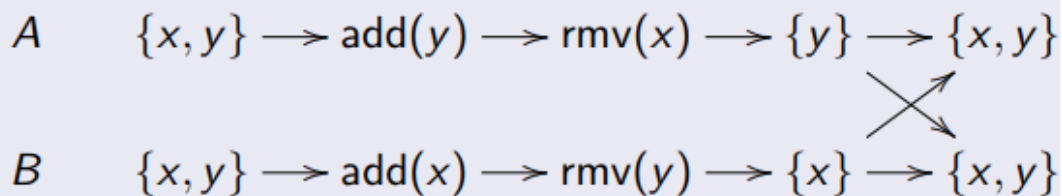
Consider a set of known operations  $O_i$ , at node  $i$ , that is ordered by an happens-before partial order  $<$ . Set has elements

$\{e \mid \text{add}(e) \in O_i \wedge \nexists \text{rmv}(e) \in O_i \cdot \text{add}(e) < \text{rmv}(e)\}$

- Redis CRDT sets are Add-Wins Sets

Can we always explain a concurrent execution by a sequential one?

### Concurrent execution



### Two (failed) sequential explanations

$H1 \quad \{x, y\} \rightarrow \dots \rightarrow \text{rmv}(x) \rightarrow \{\cancel{x}, y\}$

$H2 \quad \{x, y\} \rightarrow \dots \rightarrow \text{rmv}(y) \rightarrow \{x, \cancel{y}\}$

Concurrent executions can have richer outcomes

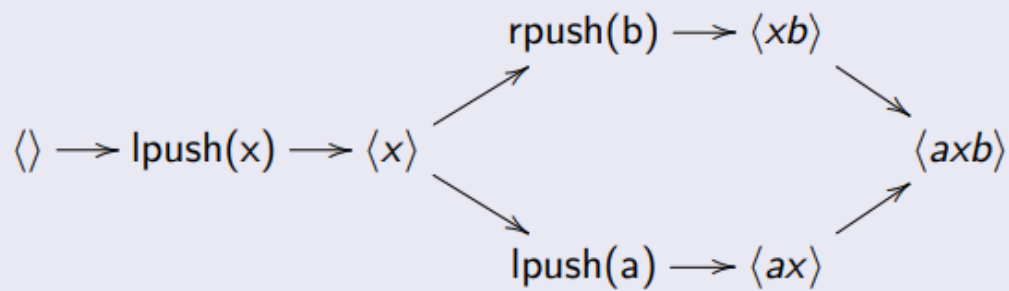
## Remove-Wins Sets

$X_i := \{e \mid \text{add}(e) \in O_i \wedge \forall \text{rmv}(e) \in O_i \cdot \text{rmv}(e) < \text{add}(e)\}$

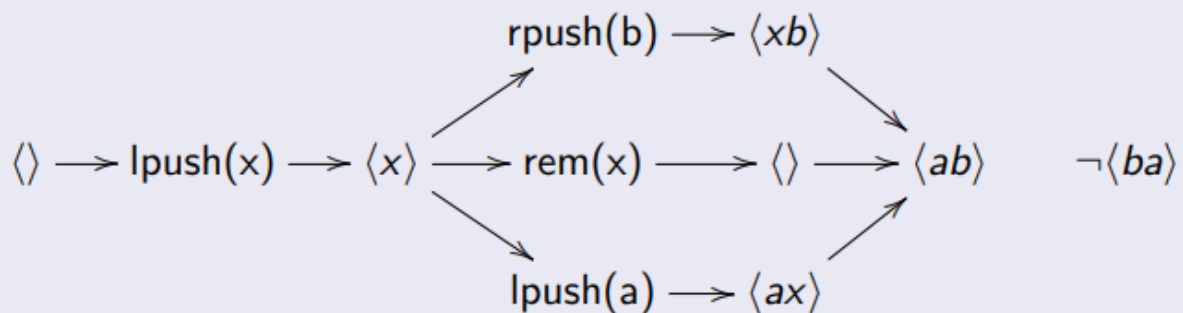
- Remove-Wins requires more metadata than Add-Wins
- Both Add and Remove-Wins have same semantics in a total order
- They are different but both preserve sequential semantics

## Sequence/List

## Element $x$ is kept



## Element $x$ is removed (Redis enforces Strong Specification)



## Summary

- Concurrent executions are needed to deal with latency
- Behaviour changes when moving from sequential to concurrent

Road to accommodate transition:

- Permutation equivalence
- Preserving sequential semantics
- Concurrent executions lead to richer outcomes

CRDTs provide sound guidelines and encode policies