High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto
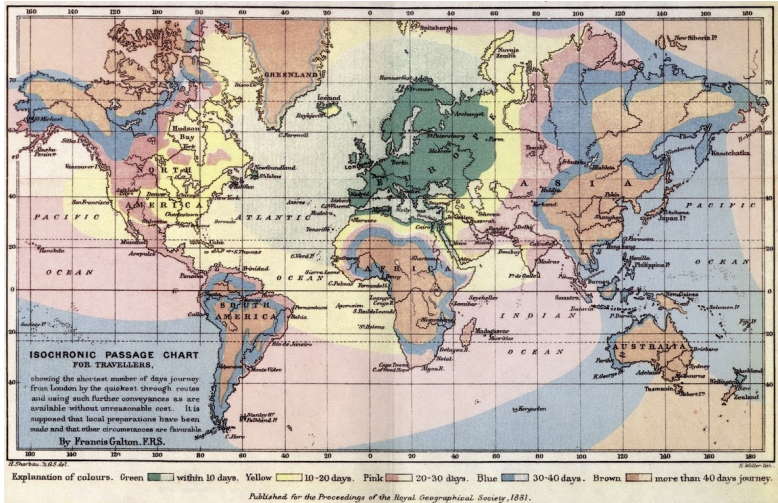
# High Availability under Eventual Consistency

Carlos Baquero
DEI, FEUP, Universidade do Porto

MEIC SDLE 2021

# The speed of communication in the 19th century
Francis Galton Isochronic Map

High Availability
under Eventual
Consistency
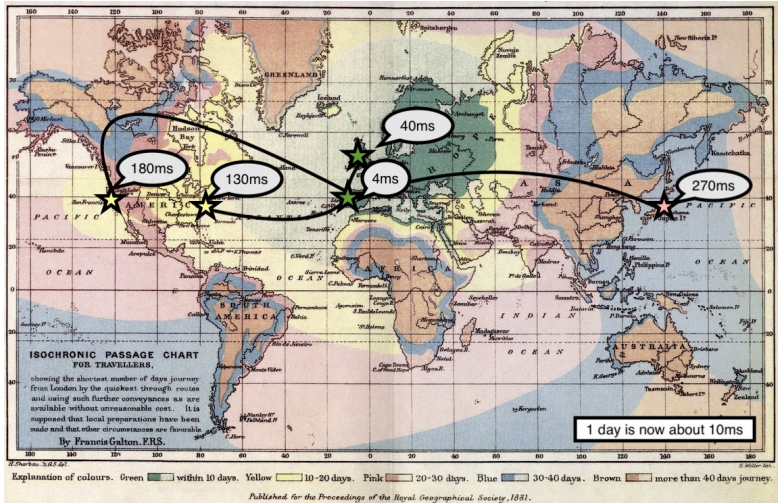
Carlos Baquero
DEI, FEUP,
Universidade do
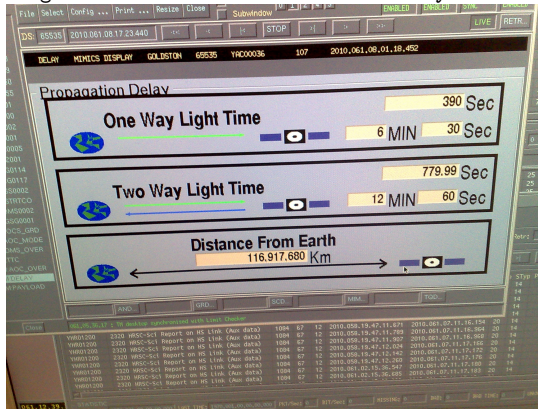Porto

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Time delay between Mars and Earth

`blogs.esa.int/mex/2012/08/05/time-delay-between-mars-and-earth/`



Delay/Disruption Tolerant Networking

`www.nasa.gov/content/dtn`

# Latency magnitudes
Geo-replication

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- $\lambda$, 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 versus client side write latency ranges

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

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

## Eventually Consistent. CACM 2009, Werner Vogels

- 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.

## CAP theorem. PODC 2000, Eric Brewer

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.

We will focus on AP.

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

A 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.

# Session Guarantees [Doug Terry, et al]

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- 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. (Writes made during the session are ordered after any Writes whose effects were seen by previous Reads in the session.)
- Monotonic Writes – writes are propagated after writes that logically precede them. (In other words, a Write is only incorporated into a server's database copy if the copy includes all previous session Writes.)

# From sequential to concurrent executions

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Consensus provides illusion of a single replica

This also preserves (slow) sequential behaviour

## Sequential execution

*Ops O*      $o \longrightarrow p \longrightarrow q$

*Time*      $- - - - - - - \succ$

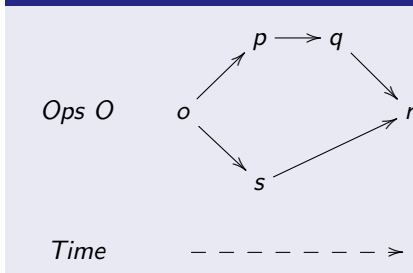We have an ordered set $(O, <)$. $O = \{o, p, q\}$ and $o < p < q$

# From sequential to concurrent executions

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

EC Multi-master (or active-active) can expose concurrency

## Concurrent execution



*Ops O*

*Time*     – – – – – – – ➤

Partially ordered set $(O, \prec)$. $o \prec p \prec q \prec r$ and $o \prec s \prec r$
Some ops in O are concurrent: $p \parallel s$ and $q \parallel s$

# Conflict-Free Replicated Data Types (CRDTs)

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- Convergence after concurrent updates. Favor AP under CAP
- Examples include counters, sets, mv-registers, maps, graphs
- Operation based CRDTs. Operation effects must commute
- State based CRDTs are rooted on join semi-lattices

# Operation-based CRDTs, effect commutativity

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- In some datatypes all operations are commutative.
- PN-Counter: $inc(dec(c)) = dec(inc(c))$
- G-Set: $add_a(add_b(s)) = add_b(add_a(s))$

For more complex examples (e.g. sets with add and remove) operations need to generate "special" commutative effects. Here we will only cover examples of state-based CRDTs as they are more common in practice.

# State-based CRDTs, Join semi-lattices

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- An (partial) ordered set S; $\langle S, \leq \rangle$.
- A join, $\sqcup$, deriving least upper bounds; $\langle S, \leq, \sqcup \rangle$.
- An initial state, usually the least element $\bot$; $\langle S, \leq, \sqcup, \bot \rangle$. ($\forall a \in S, a \sqcup \bot = a$)
- Alternative to a (unique) initial state, is a one time init in each replica assigning any element from $S$.
- Join properties in a semilattice $\langle S, \leq, \sqcup \rangle$:
    - Idempotence, $a \sqcup a = a$,
    - Commutativity, $a \sqcup b = b \sqcup a$,
    - Associative, $(a \sqcup b) \sqcup c = a \sqcup (b \sqcup c)$.
- $\leq$ reflects monotonic state evolution – increase of information.
- Updates must conform to $\leq$.
- In general, queries can return non-monotonic values, and in other domains than $S$. E.g: Returning a set size.

# Eventual Consistency, non stop

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Now convergence can related to known updates, with no need to
stop updates: $\mathrm{upds}(a) \subseteq \mathrm{upds}(b) \Rightarrow a \leq b$.

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

# Design of Conflict-Free Replicated Data Types

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

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 $O_i$

Example: Counter at $i$ is $|\{\text{inc} \mid \text{inc} \in O_i\}| - |\{\text{dec} \mid \text{dec} \in O_i\}|$

CRDTs are efficient representations that follow some general rules

# Principle of permutation equivalence

High Availability
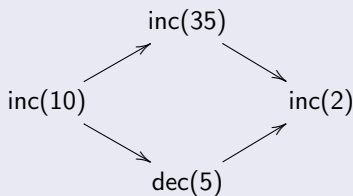under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

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

## Sequential

$$inc(10) \longrightarrow inc(35) \longrightarrow dec(5) \longrightarrow inc(2)$$

$$dec(5) \longrightarrow inc(2) \longrightarrow inc(10) \longrightarrow inc(35)$$
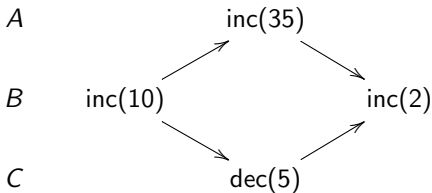
## Concurrent



You guessed: Result is 42

# Implementing Counters
Example: CRDT PNCounters

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Lets track total number of incs and decs done at each replica

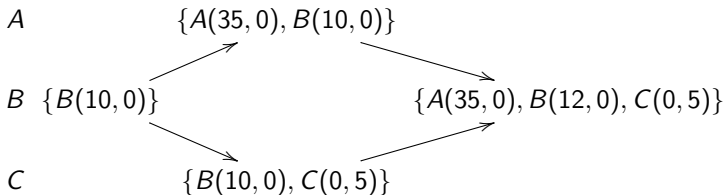$$\{A(\text{incs}, \text{decs}), \ldots, C(\ldots, \ldots)\}$$

# Implementing Counters
Example: CRDT PNCounters

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Separate positive and negative counts are kept per replica

$A$ $\quad\quad\quad\quad \{A(35, 0), B(10, 0)\}$

$B \ \{B(10, 0)\}$ $\quad\quad\quad\quad\quad\quad\quad \{A(35, 0), B(12, 0), C(0, 5)\}$

$C$ $\quad\quad\quad\quad \{B(10, 0), C(0, 5)\}$

Joining does point-wise maximums among entries (semilattice)

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

# State-based CRDTs: PN-Counter

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

$$
\begin{aligned}
\Sigma &= I \to \mathbb{N} \times \mathbb{N} \\
\sigma_i^0 &= \{(r, (0, 0)) \mid r \in I\} \\
\mathsf{apply}_i(\mathsf{inc}, m) &= m\{i \mapsto (\mathsf{fst}(m(i)) + 1, \mathsf{snd}(m(i)))\} \\
\mathsf{apply}_i(\mathsf{dec}, m) &= m\{i \mapsto (\mathsf{fst}(m(i)), \mathsf{snd}(m(i)) + 1)\} \\
\mathsf{eval}_i(\mathsf{rd}, m) &= \sum_{r \in I} \mathsf{fst}(m(r)) - \mathsf{snd}(m(r)) \\
\mathsf{merge}_i(m, m') &= \{(r, \max(m(r), m'(r))) \mid r \in I\}
\end{aligned}
$$

Note: max is pointwise maximum

# Registers

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Registers are an ordered set of write operations

## Sequential execution

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

## Sequential execution under distribution

$A \qquad wr(x) \qquad \qquad \qquad \qquad wr(x)$

$B \qquad \qquad \qquad wr(j) \longrightarrow wr(k)$

Register value is $x$, the last written value

# Last Writer Wins

A simple approach to evolve state without strong coordination, is to adopt a *Last Writer Wins* policy (see also Thomas write rule). Recently popularized in the Cassandra system, this policy uses timestamps to discard *older* writes and attain convergence.

$$
\begin{aligned}
\Sigma &= T \times \mathbb{N} \\
\sigma_i^0 &= (0, 0) \\
\text{apply}_i((\text{wr}, t', n'), (t, n)) &= (t, n) \textbf{ if } t' < t \textbf{ else } (t', n') \\
\text{eval}_i(\text{rd}, (t, n)) &= n \\
\text{merge}_i((t, n), (t', n')) &= (t, n) \textbf{ if } t' < t \textbf{ else } (t', n')
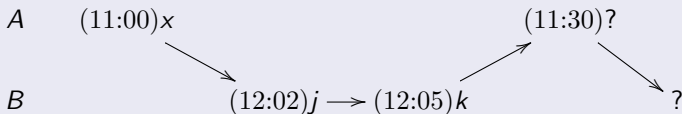\end{aligned}
$$

LWW Integer

# Implementing Registers
*Naive* Last-Writer-Wins

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

CRDT register implemented by attaching local wall-clock times

## Sequential execution under distribution

$A$     $(11:00)x$                                $(11:30)?$

$B$                   $(12:02)j \longrightarrow (12:05)k$           ?

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

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

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Register shows value $v$ at replica $i$ iff

$$\mathrm{wr}(v) \in O_i$$

and

$$\nexists \mathrm{wr}(v') \in O_i \cdot \mathrm{wr}(v) < \mathrm{wr}(v')$$

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Concurrent semantics should preserve the sequential semantics

This also ensures correct sequential execution under distribution
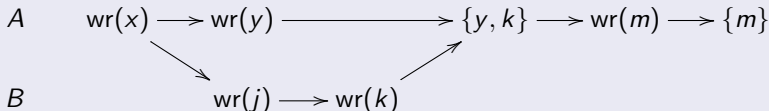
# Multi-value Registers

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Concurrency semantics shows all concurrent values

$$\{v \mid \mathrm{wr}(v) \in O_i \wedge \nexists \mathrm{wr}(v') \in O_i \cdot \mathrm{wr}(v) \prec \mathrm{wr}(v')\}$$

### Concurrent execution



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

$B \quad \mathrm{wr}(j) \longrightarrow \mathrm{wr}(k)$

Dynamo shopping carts are multi-value registers with payload sets

The $m$ value could be an application level merge of values $y$ and $k$
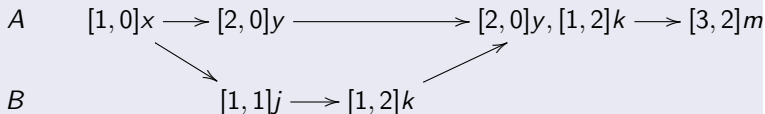
# Implementing Multi-value Registers

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Concurrency can be precisely tracked with version vectors

### Concurrent execution (version vectors)

$A$ $\quad [1,0]x \longrightarrow [2,0]y \longrightarrow [2,0]y, [1,2]k \longrightarrow [3,2]m$

$B$ $\qquad\qquad\qquad [1,1]j \longrightarrow [1,2]k$

Metadata can be compressed with a common causal context and a single scalar per value (dotted version vectors)

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

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

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

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

# State-based CRDTs: 2P-Set

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

# Sets
Sequential Semantics

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Consider add and rmv operations

$$X = \{\ldots\}, \; \text{add(a)} \longrightarrow \text{add(c)} \text{ we observe that } \text{a}, \text{c} \in X$$

$$X = \{\ldots\}, \; \text{add(c)} \longrightarrow \text{rmv(c)} \text{ we observe that } \text{c} \notin X$$

In general, given $O_i$, the set has elements

$$\{e \mid \text{add(e)} \in O_i \land \nexists \text{rmv(e)} \in O_i \cdot \text{add(e)} < \text{rmv(e)}\}$$
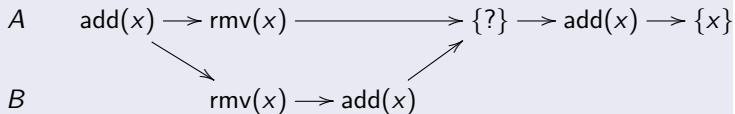
2P-Set breaks sequential semantics

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Problem: Concurrently adding and removing the same element

## Concurrent execution

$A$ $\quad$ add$(x) \longrightarrow$ rmv$(x) \longrightarrow \{?\} \longrightarrow$ add$(x) \longrightarrow \{x\}$

$B$ $\quad\quad\quad\quad\quad$ rmv$(x) \longrightarrow$ add$(x)$

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Let's choose Add-Wins

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

$$\{e \mid \mathsf{add}(e) \in O_i \ \wedge \nexists \, \mathsf{rmv}(e) \in O_i \cdot \mathsf{add}(e) \prec \mathsf{rmv}(e)\}$$

Is this familiar?

The sequential semantics applies identical rules on a total order

Redis CRDT sets are Add-Wins Sets

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

$$
\begin{aligned}
\Sigma &= \mathcal{P}(T \times V) \times \mathcal{P}(T) \\
\sigma_i^0 &= \{\}, \{\} \\
\mathrm{apply}_i((\mathrm{add}, v), (s, t)) &= s \cup \{(utag(), v)\}, t \\
\mathrm{apply}_i((\mathrm{rmv}, v), (s, t)) &= s, t \cup \{u \mid (u, v) \in s\} \\
\mathrm{eval}_i(\mathrm{rd}, s) &= \{v \mid (u, v) \in s \wedge u \notin t\} \\
\mathrm{merge}_i((s, t), (s', t')) &= s \cup s', t \cup t'
\end{aligned}
$$

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

### Concurrent execution

$A \qquad \{x, y\} \longrightarrow \mathsf{add}(y) \longrightarrow \mathsf{rmv}(x) \longrightarrow \{y\} \longrightarrow \{x, y\}$

$B \qquad \{x, y\} \longrightarrow \mathsf{add}(x) \longrightarrow \mathsf{rmv}(y) \longrightarrow \{x\} \longrightarrow \{x, y\}$

### Two (failed) sequential explanations

$H1 \qquad \{x, y\} \longrightarrow \ldots \longrightarrow \mathsf{rmv}(x) \longrightarrow \{\not{x}, y\}$

$H2 \qquad \{x, y\} \longrightarrow \ldots \longrightarrow \mathsf{rmv}(y) \longrightarrow \{x, \not{y}\}$

Concurrent executions can have richer outcomes

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Alternative: Let's choose Remove-Wins

$$X_i \doteq \{e \mid \mathsf{add}(e) \in O_i \ \wedge \forall \, \mathsf{rmv}(e) \in O_i \cdot \mathsf{rmv}(e) \prec \mathsf{add}(e)\}$$

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

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

Thought experiment: A CRDT Set data type could store enough
information to allow a parametrized query that shows both Add-Wins
or Remove-Wins. One can expect a metadata size penalty for the
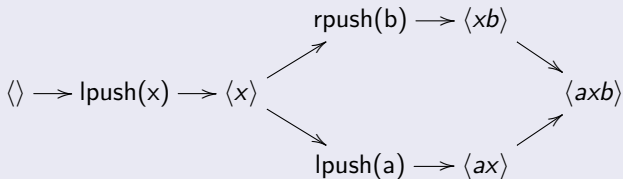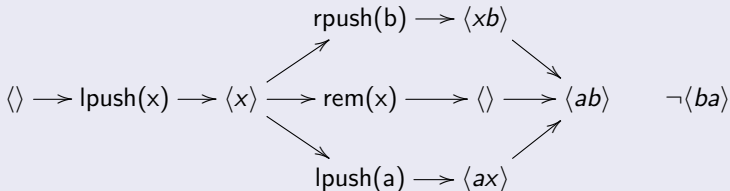flexibility.

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

## Element $x$ is kept



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

High Availability
under Eventual
Consistency

Carlos Baquero
DEI, FEUP,
Universidade do
Porto

- 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