High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

# High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

MIEI SDLE 2021

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

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

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで

# The speed of communication in the 21st century RTT data gathered via http://www.azurespeed.com

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho



・ロット (雪) ・ (日) ・ (日) ・ (日)

# The speed of communication in the 21st century If you really like high latencies ...

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### Time delay between Mars and Earth

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

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @



Delay/Disruption Tolerant Networking www.nasa.gov/content/dtn

### Latency magnitudes Geo-replication

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

- $\lambda$ , up to 50ms (local region DC)
- Λ, 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  $[\Lambda, 2\Lambda]$ Primary-Backup  $[\lambda, \Lambda]$ Multi-Master  $\lambda$  (with no divergence) (asynchronous/lazy)

(allowing divergence)

# EC and CAP for Geo-Replication

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### 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 – Eventual Consistency

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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 Universidade do Minho

- 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.)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# From sequential to concurrent executions

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Consensus provides illusion of a single replica

This also preserves (slow) sequential behaviour

Sequential execution

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

*Time* ---- >

We have an ordered set (0, <).  $O = \{o, p, q\}$  and o

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# From sequential to concurrent executions

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho EC Multi-master (or active-active) can expose concurrency

Concurrent execution



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  $q \parallel s$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

# Conflict-Free Replicated Data Types (CRDTs)

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

- 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 Universidade do Minho

- In some datatypes all operations are commutative.
- PN-Counter: inc(dec(c)) = dec(inc(c))
- G-Set:  $\operatorname{add}_a(\operatorname{add}_b(s)) = \operatorname{add}_b(\operatorname{add}_a(s))$

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.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# State-based CRDTs, Join semi-lattices

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

- An (partial) ordered set S;  $\langle S, \leq \rangle$ .
- A join,  $\Box$ , deriving least upper bounds;  $\langle S, \leq, \Box \rangle$ .
- 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 ≤.
- In general, queries can return non-monotonic values, and in other domains than S. E.g: Returning a set size.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# Eventual Consistency, non stop

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

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

# Design of Conflict-Free Replicated Data Types

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

> 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  $|\{inc | inc \in O_i\}| - |\{dec | dec \in O_i\}|$

CRDTs are efficient representations that follow some general rules

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# Principle of permutation equivalence

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho 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



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

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

(日) (四) (日) (日) (日)



Joining does point-wise maximums among entries (semilattice)

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

# State-based CRDTs: PN-Counter

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Note: max is pointwise maximum



#### High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Registers are an ordered set of write operations

Sequential execution

$$A \qquad \operatorname{wr}(x) \longrightarrow \operatorname{wr}(j) \longrightarrow \operatorname{wr}(k) \longrightarrow \operatorname{wr}(x)$$

Sequential execution under distribution



▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Register value is x, the last written value

# Last Writer Wins

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho 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.

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

### LWW Integer

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

### Implementing Registers Naive Last-Writer-Wins

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

CRDT register implemented by attaching local wall-clock times



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

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Register shows value v at replica i iff

 $wr(v) \in O_i$ 

and

$$eq \mathsf{wr}(v') \in \mathit{O}_i \cdot \mathsf{wr}(v) < \mathsf{wr}(v')$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

# Preservation of sequential semantics

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Concurrent semantics should preserve the sequential semantics

This also ensures correct sequential execution under distribution

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

# Multi-value Registers

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho Concurrency semantics shows all concurrent values

$$\{v \mid \mathsf{wr}(v) \in O_i \land \nexists \mathsf{wr}(v') \in O_i \cdot \mathsf{wr}(v) \prec \mathsf{wr}(v')\}$$



Dynamo shopping carts are multi-value registers with payload sets

The m value could be an application level merge of values y and k

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# Implementing Multi-value Registers

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Concurrency can be preciselly tracked with version vectors

Concurrent execution (version vectors)

$$A \qquad [1,0]x \longrightarrow [2,0]y \longrightarrow [2,0]y, [1,2]k \longrightarrow [3,2]m$$
$$B \qquad [1,1]j \longrightarrow [1,2]k$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

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

Registers in Redis LWW arbitration

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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 Universidade do Minho

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

◆□ ▶ ◆□ ▶ ◆ 臣 ▶ ◆ 臣 ▶ ○ 臣 ○ のへで

# State-based CRDTs: 2P-Set

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Consider add and rmv operations

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

 $X = \{\ldots\}, \ \mathsf{add}(\mathsf{c}) \longrightarrow \mathsf{rmv}(\mathsf{c}) \ \mathrm{we \ observe \ that} \ \mathsf{c} \not\in \mathsf{X}$ 

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

 $\{e \mid \mathsf{add}(e) \in \mathsf{O}_{\mathsf{i}} \land \nexists \mathsf{rmv}(e) \in \mathsf{O}_{\mathsf{i}} \cdot \mathsf{add}(e) < \mathsf{rmv}(e)\}$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

2P-Set breaks sequential semantics

Sets Concurrency Semantics

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### Problem: Concurrently adding and removing the same element



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### 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

 $\{e \mid \mathsf{add}(e) \in \mathsf{O}_i \ \land \nexists \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{add}(e) \prec \mathsf{rmv}(e)\}$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

 $\{e \mid \mathsf{add}(e) \in \mathsf{O}_i \ \land \nexists \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{add}(e) \prec \mathsf{rmv}(e)\}$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Is this familiar?

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

 $\{e \mid \mathsf{add}(e) \in \mathsf{O}_i \ \land \nexists \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{add}(e) \prec \mathsf{rmv}(e)\}$ 

Is this familiar?

The sequential semantics applies identical rules on a total order

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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

 $\{e \mid \mathsf{add}(e) \in \mathsf{O}_i \ \land \nexists \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{add}(e) \prec \mathsf{rmv}(e)\}$ 

Is this familiar?

The sequential semantics applies identical rules on a total order

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

Redis CRDT sets are Add-Wins Sets

# State-based CRDTs: Basic Add-Wins Set

Observed-remove add-wins

#### High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

$$\Sigma = \mathcal{P}(T \times V) \times \mathcal{P}(T)$$

$$\sigma_i^0 = \{\}, \{\}$$

$$\mathsf{apply}_i((\mathsf{add}, v), (s, t)) = s \cup \{(utag(), v)\}, t$$

$$\mathsf{apply}_i((\mathsf{rmv}, v), (s, t)) = s, t \cup \{u \mid (u, v) \in s\}$$

$$\mathsf{eval}_i(\mathsf{rd}, s) = \{v \mid (u, v) \in s \land u \notin t\}$$

$$\mathsf{merge}_i((s, t), (s', t')) = s \cup s', t \cup t'$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

### Equivalence to a sequential execution? Add-Wins Sets

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho Can we always explain a concurrent execution by a sequential one?

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

### 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

Concurrency Semantics Remove-Wins Sets

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### Alternative: Let's choose Remove-Wins

$$X_i \doteq \{e \mid \mathsf{add}(e) \in \mathsf{O}_i \ \land \forall \ \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{rmv}(e) \prec \mathsf{add}(e)\}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Concurrency Semantics Remove-Wins Sets

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

Alternative: Let's choose Remove-Wins

 $X_i \doteq \{e \mid \mathsf{add}(e) \in \mathsf{O}_i \land \forall \mathsf{rmv}(e) \in \mathsf{O}_i \cdot \mathsf{rmv}(e) \prec \mathsf{add}(e)\}$ 

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

# Overview. Delaying choice of semantics

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

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.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

# Sequence/List Weak/Strong Specification [Attiya et al, PODC 16]

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

### Element x is kept



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



# Take home message

High Availability under Eventual Consistency

Carlos Baquero Universidade do Minho

- 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