High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

# High Availability under Eventual Consistency

Carlos Baquero
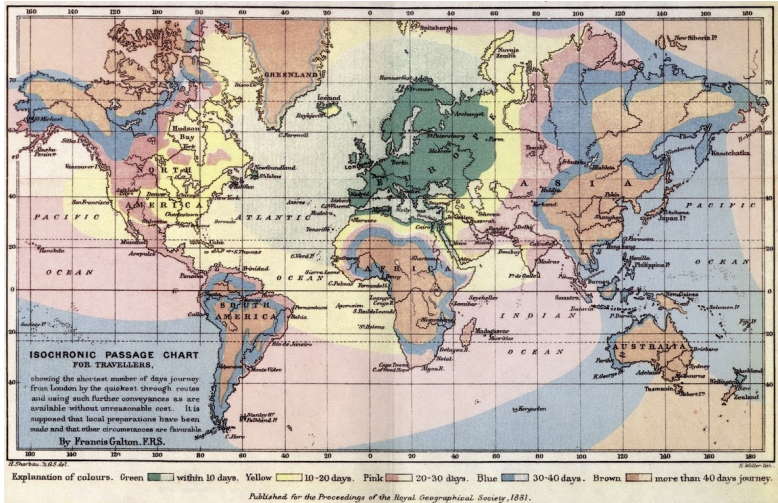Universidade do Minho & INESC TEC

Shonan PL4DS Meeting 2019

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

High Availability
under Eventual
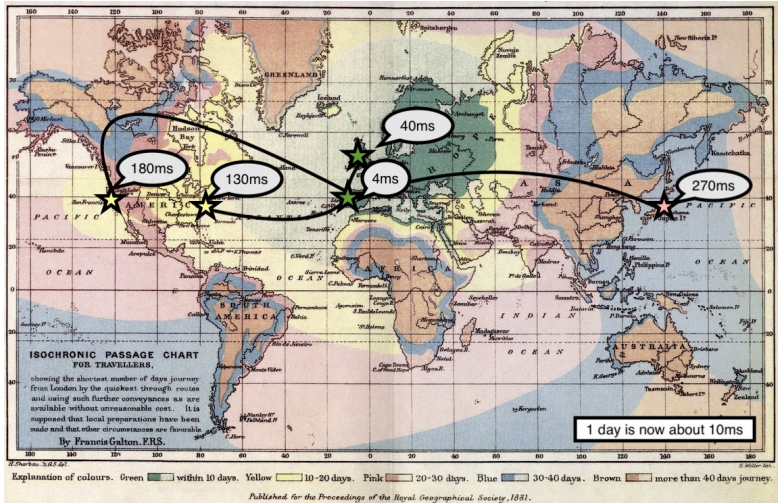Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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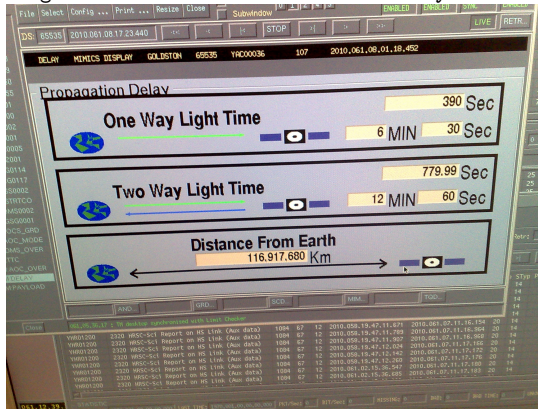
High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

## 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

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

# Latency magnitudes
Geo-replication

- $\lambda$, 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) |

Multi-Master executions maximize availability and response speeds

# From sequential to concurrent executions
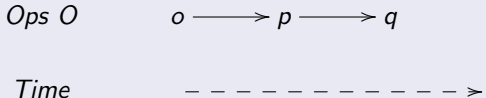
High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Consensus provides illusion of a single replica

This also preserves (slow) sequential behaviour

## Sequential execution

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

*Time*       $- - - - - - - - - - - >$

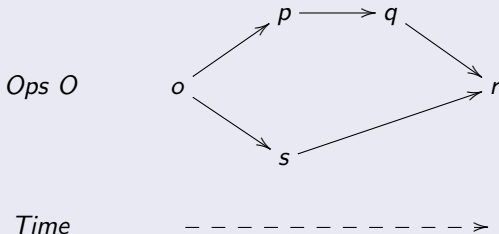We have an ordered set $(O, <)$. $O = \{o, p, q\}$ and $o < p < q$

# From sequential to concurrent executions

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Multi-master 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 $q \parallel s$

# Conflict-Free Replicated Data Types (CRDTs)

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

- 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

# Design of Conflict-Free Replicated Data Types

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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 following some **design principles**

# Principle of permutation equivalence
### E.g. Counters

High Availability
under Eventual
Consistency

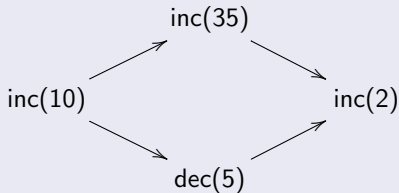Carlos Baquero
Universidade do
Minho & INESC
TEC

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

# Registers

Registers are an ordered set of write operations

## Sequential execution

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

## Sequential execution under distribution

$A \qquad \mathrm{wr}(x) \qquad\qquad\qquad\qquad\qquad \mathrm{wr}(x)$

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

Register value is $x$, the last written value

# Implementing Registers
*Naive* Last-Writer-Wins

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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
Universidade do
Minho & INESC
TEC

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
Universidade do
Minho & INESC
TEC

Concurrent semantics should preserve the sequential semantics

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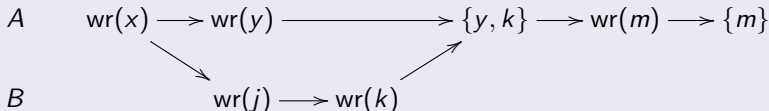
# Multi-value Registers

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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

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

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

Dynamo shopping carts are multi-value registers with payload sets

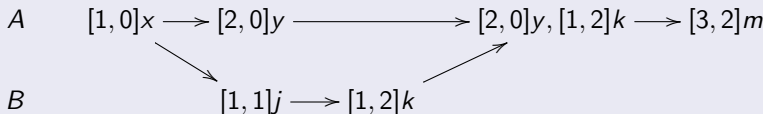The $m$ value could be an application level merge of values $y$ and $k$

# Implementing Multi-value Registers

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Concurrency can be precisely tracked with version vectors

## Concurrent execution (version vectors)

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

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

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
Universidade do
Minho & INESC
TEC

Consider add and rmv operations

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

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

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

$$\{e \mid \mathsf{add(e)} \in \mathsf{O_i} \wedge \nexists \mathsf{rmv(e)} \in \mathsf{O_i} \cdot \mathsf{add(e)} < \mathsf{rmv(e)}\}$$
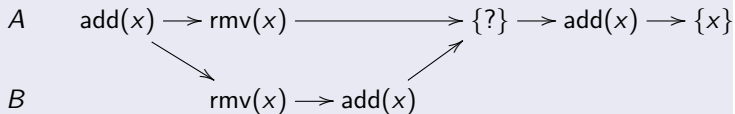
# Sets
Concurrency Semantics

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Problem: Concurrently adding and removing the same element

## Concurrent execution

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

$B$     $\text{rmv}(x) \longrightarrow \text{add}(x)$

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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

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

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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

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

Is this familiar?

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

## Concurrency Semantics
Add-Wins Sets

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

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

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

### Concurrent execution

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

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

### Two (failed) sequential explanations

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

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

Concurrent executions can have richer outcomes

## Concurrency Semantics
Remove-Wins Sets

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Alternative: Let's choose Remove-Wins

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

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

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

# Reference

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

Conflict-Free Replicated Data Types CRDTs.
Encyclopedia of Big Data Technologies, 2019.
Nuno M. Preguiça, Carlos Baquero, Marc Shapiro.

# Open Questions

High Availability
under Eventual
Consistency

Carlos Baquero
Universidade do
Minho & INESC
TEC

- Given operations and queries can CRDTs be synthesized?
- Stateful (dis)order tolerant communication protocols
- Integrity guaranties and privacy in data evolution