# CRDTs: State-based approaches to high availability

Carlos Baquero
Universidade do Porto

cbm@fe.up.pt

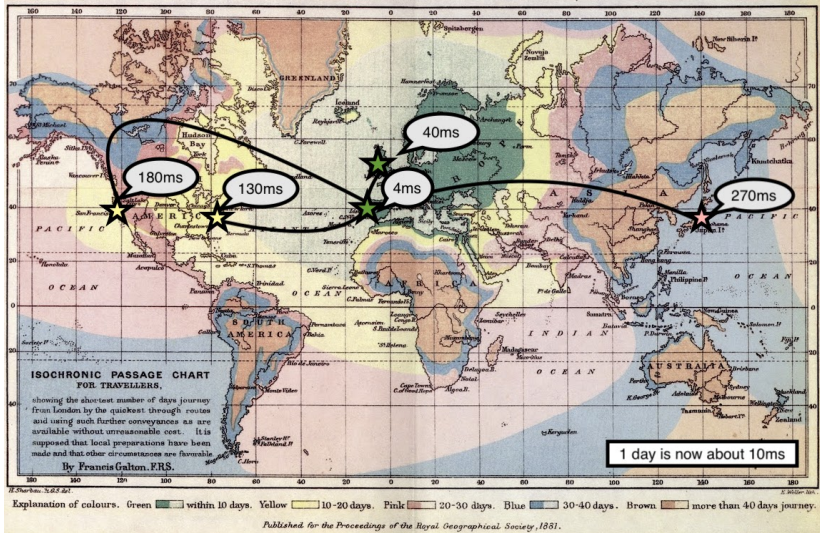DARE

September 12th, 2023. Brussels, BE

# Outline

- CAP trade-off and motivation for CRDTs
- Operations vs State
- State based CRDTs
- Delta state based CRDTs
- Decomposing states
- Update propagation

# The speed of communication in the 19th century

Francis Galton Isochronic Map

# The speed of communication in the 21st century

RTT data gathered via http://www.azurespeed.com

# Communication networks will keep expanding



Earth $\leftrightarrow$ Mars link, depends on planets positions:

- One way: 3 to 21 minutes
- Round trip: 6 to 42 minutes

# 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)

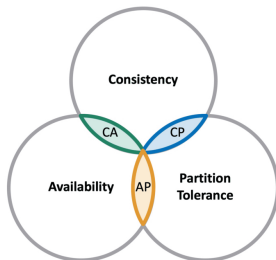Multi-Master $\lambda$ (allowing divergence)

# EC and CAP for Geo-Replication

## 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 the three properties – data consistency, system availability, and tolerance to network partition – only two can be simultaneously achieved.



(Image from https://hazelcast.com/glossary/cap-theorem/)

# High Availability – Eventual Consistency

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.

# From sequential to concurrent executions

Consensus provides illusion of a single replica

This also preserves (slow $\approx$ 200ms) sequential behaviour

## Sequential execution

Ops O $\qquad o \longrightarrow p \longrightarrow q$
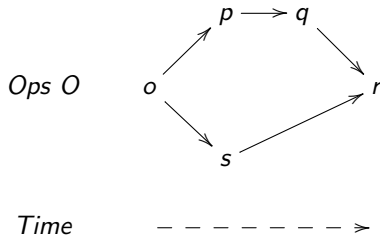
Time $\qquad - - - - - - - \succ$

We have a total order on the set of operations $(O, <)$.
$O = \{o, p, q\}$ and $o < p < q$

# From sequential to concurrent executions

EC Multi-master is fast, any replica can locally accept updates
But 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)

- ▶ 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
- ▶ State based can emulate operation based – and vice-versa
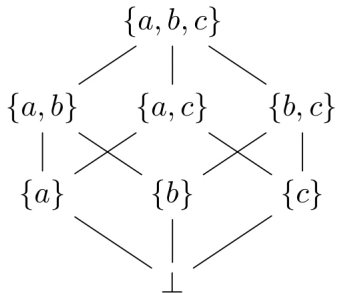
# Operation-based CRDTs, effect commutativity

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

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

- ▶ 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 state evolution/inflation.
- ▶ Updates/mutations $m()$ increase information, $x \leq m(x)$.
- ▶ In general, queries can return non-monotonic values, and in other domains than $S$. E.g: Returning a set size.

# State-based CRDTs, Join semi-lattices



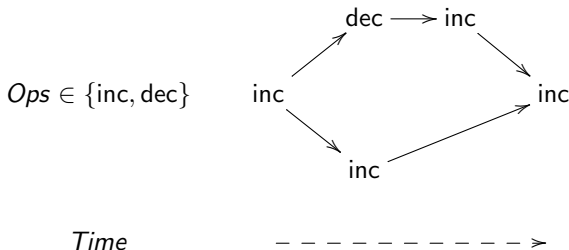(b) $\mathsf{GSet}\langle\{a, b, c\}\rangle$.

# 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 $O_i$

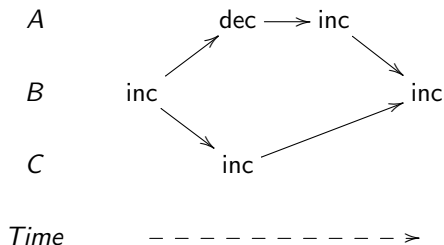Example: Counter at $i$ is $|\{\text{inc} \mid \text{inc} \in O_i\}| - |\{\text{dec} \mid \text{dec} \in O_i\}|$

$$Ops \in \{\text{inc}, \text{dec}\}$$



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

State based CRDTs are efficient specialized representations of pologs
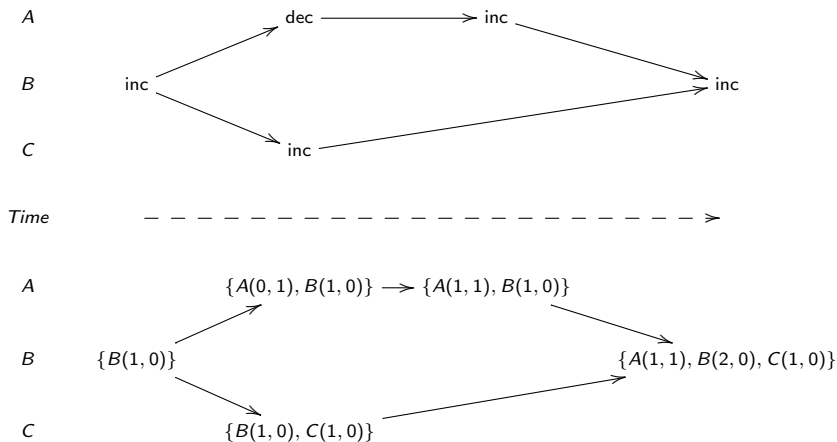
# Implementing Counters
Example: CRDT PNCounters



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

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

This is a semi-lattice, with $\bot = \{A(0,0), B(0,0), C(0,0)\}$

# Implementing Counters
Example: CRDT PNCounters



A          dec ⟶ inc

B    inc                    inc

C          inc

Time   – – – – – – – – – – – – – – – – ⟶

A      $\{A(0,1), B(1,0)\} \rightarrow \{A(1,1), B(1,0)\}$

B   $\{B(1,0)\}$                   $\{A(1,1), B(2,0), C(1,0)\}$

C         $\{B(1,0), C(1,0)\}$

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

# State-based GCounter

$$
\begin{aligned}
\text{GCounter} &= \mathbb{I} \hookrightarrow \mathbb{N} \\
\bot &= \{\} \\
\text{inc}_i(m) &= m\{i \mapsto m[i] + 1\} \\
\text{value}(m) &= \sum_{j \in \text{dom } m} m[j] \\
m \sqcup m' &= \{j \mapsto \max(m[j], m'[j]) \mid j \in \text{dom } m \cup \text{dom } m'\}
\end{aligned}
$$

- State maps replica ids to integers
- inc increments self entry ($i$)
- Join is pointwise max
- Similar in structure to a version vector

# State-based PNCounter

$$
\begin{aligned}
\text{PNCounter} &= \text{GCounter} \times \text{GCounter} \\
\bot &= (\bot, \bot) \\
\text{inc}_i((p, n)) &= (\text{inc}_i(p), n) \\
\text{dec}_i((p, n)) &= (p, \text{inc}_i(n)) \\
\text{value}((p, n)) &= \text{value}(p) - \text{value}(n) \\
(p, n) \sqcup (p', n') &= (p \sqcup p', n \sqcup n')
\end{aligned}
$$

▶ Solved through a pair of GCounters (product composition)
  ▶ increments and decrements tracked separately
  ▶ counter value obtained as difference
▶ In practice a single map to pairs is used (as seen on the example)

# Grow only set GSet$\langle E \rangle$

$$
\begin{aligned}
\mathsf{GSet}\langle E \rangle &= \mathcal{P}(E) \\
\bot &= \{\} \\
\mathsf{add}_i(e, s) &= s \cup \{e\} \\
\mathsf{elements}(s) &= s \\
\mathsf{contains}(e, s) &= e \in s \\
s \sqcup s' &= s \cup s'
\end{aligned}
$$

► Trivial CRDT: state same as sequential datatype
► Add already an inflation
► Merging replica states by set union
► Anonymous CRDT: no need for node ids in the state

How to delete from a set?

# Two phase set 2PSet⟨E⟩

$$
\begin{aligned}
\text{2PSet}\langle E\rangle &= \mathcal{P}(E) \times \mathcal{P}(E) \\
\bot &= \{\} \times \{\} \\
\text{add}_i(e, (s, t)) &= (s \cup \{e\}, t) \\
\text{rmv}_i(e, (s, t)) &= (s, t \cup \{e\}) \\
\text{elements}((s, t)) &= s \setminus t \\
\text{contains}(e, (s, t)) &= e \in (s \setminus t) \\
(s, t) \sqcup (s', t') &= (s \cup s', t \cup t')
\end{aligned}
$$

Problems:

- ▶ Does not comply with usual sequential behaviour
- ▶ Cannot delete and readd elements

Lets try to add unique tags to each added element

# Basic add-wins set AWSet⟨E⟩

$$
\begin{aligned}
\mathsf{AWSet}\langle E\rangle &= \mathcal{P}(E \times U) \times \mathcal{P}(U) \\
\bot &= \{\} \times \{\} \\
\mathsf{add}_i(e,(s,t)) &= (s \cup \{(e, utag())\}, t) \\
\mathsf{rmv}_i(e,(s,t)) &= (s, t \cup \{u \mid (e,u) \in s\}) \\
\mathsf{elements}((s,t)) &= \{e \mid (e,u) \in s \wedge u \notin t\} \\
\mathsf{contains}(e,(s,t)) &= e \in \mathsf{elements}((s,t)) \\
(s,t) \sqcup (s',t') &= (s \cup s', t \cup t')
\end{aligned}
$$

Problems:

- ▶ Removed elements are still kept in s
- ▶ How to generate unique tags? random strings?

# Dots as Causal Histories

We can use the node id $i \in \mathbb{I}$ and a natural counter $\mathbb{N}$, to create tags of the form $\mathbb{I} \times \mathbb{N}$.

- $\{(a, 1), (a, 2), (b, 1), (b, 2), (b, 3), (c, 1), (c, 2), (c, 3)\}$
- $\{a \mapsto 2, b \mapsto 3, c \mapsto 3\}$

In cases with fixed membership one can compress to a vector

- $[2, 3, 3]$

# Causal Context

$$
\begin{aligned}
\mathsf{CausalContext} &= \mathcal{P}(\mathbb{I} \times \mathbb{N}) \\
\max_i(c) &= \max(\{n \mid (i, n) \in c\} \cup \{0\}) \\
\mathsf{next}_i(c) &= (i, \max_i(c) + 1)
\end{aligned}
$$

Example: $\mathsf{next}_a(\{(a, 1), (a, 2), (b, 1), (b, 2), (b, 3)\}) = (a, 3)$

# (causal) Add-wins set AWSet$\langle E \rangle$

$$
\begin{aligned}
\text{AWSet}\langle E \rangle &= \mathcal{P}(E \times \mathbb{I} \times \mathbb{N}) \times \mathcal{P}(\mathbb{I} \times \mathbb{N}) \\
\bot &= \{\} \times \{\} \\
\text{add}_i(e, (s, c)) &= (s \cup \{(e, d)\}, c \cup \{d\}) \textbf{ where } d = \text{next}_i(c) \\
\text{rmv}_i(e, (s, c)) &= (s \setminus \{(e, d) \mid (e, d) \in s\}, c) \\
\text{elements}((s, c)) &= \{e \mid (e, d) \in s\} \\
\text{contains}(e, (s, c)) &= e \in \text{elements}((s, c)) \\
(s, c) \sqcup (s', c') &= (s \cap s' \cup f(s, c') \cup f(s', c), c \cup c') \\
&\qquad \textbf{where } f(a, b) = \{(e, d) \mid (e, d) \in a \wedge d \notin b\}
\end{aligned}
$$

Note 1: For simplicity we assume $(e, (i, n)) \equiv (e, i, n)$.

Note 2: $f(a, b)$ is a filter function that preserves elements in $a$ not known in $b$.

# State based CRDTs

### Pros

- ▶ Permissive communication model: Drops, Reorder, Duplicates
- ▶ Causal consistency by default
- ▶ Decouples operation rate from transmission rate
- ▶ Can compact operations

### Cons

- ▶ All metadata makes part of the state
- ▶ Requires full state transmission

What if could separate in updates the state that inflates the lattice?

# Delta state-based GCounter

$$
\begin{aligned}
\mathsf{GCounter} &= \mathbb{I} \hookrightarrow \mathbb{N} \\
\bot &= \{\} \\
\mathsf{inc}_i(m) &= m\{i \mapsto m[i] + 1\} \\
\mathsf{inc}_i^{\delta}(m) &= \{i \mapsto m[i] + 1\} \\
\mathsf{value}(m) &= \sum_{j \in \mathsf{dom}\, m} m[j] \\
m \sqcup m' &= \{j \mapsto \mathsf{max}(m[j], m'[j]) \mid j \in \mathsf{dom}\, m \cup \mathsf{dom}\, m'\}
\end{aligned}
$$

# Delta State CRDTs

- State: like normal state-based CRDTs
- Messages: based on *delta-states*, that are states, hopefully small
- Essential difference: delta-mutators, which return delta states, to be
  - joined with current state
  - joined with other delta states, forming *delta-groups d*, to send in messages
- For each mutator $m$ we can define a delta-mutator $m^\delta$ such that:

$$m(X) = X \sqcup m^\delta(X)$$

# Basic anti-entropy algorithm

**durable state:**
  $X_i \in S$, CRDT state; initially
    $X_i = \bot$

**volatile state:**
  $D_i \in S$, delta-buffer; initially
    $D_i = \bot$

**on** operation$_i(m^\delta)$
  **let** $d = m^\delta(X_i)$
  $X_i \leftarrow X_i \sqcup d$
  $D_i \leftarrow D_i \sqcup d$

**on** receive$_{j,i}(d)$
  $X_i \leftarrow X_i \sqcup d$
  $D_i \leftarrow D_i \sqcup d$

**periodically**
  **let** $m = $ choose$_i(X_i, D_i)$
  **for** $j$ **in** neighbors$_i$ **do**
    send$_{i,j}(m)$
  $D_i \leftarrow \bot$

# Problems with naive anti-entropy

Naive transitive propagation causes much redundancy

- ▶ deltas re-propagated back to where they came from
- ▶ delta-group wholly joined to delta-buffer, even if mostly reflected in state

Solution

- ▶ avoid back-propagation of delta-groups; tagging origin
- ▶ remove redundant state in received delta-groups; how?

We need tools to compare states

# Join Decompositions

$D$ is a irredundant join decomposition of $s$ :

(decomposition) $$\bigsqcup D = s$$

(irredundant) $$\forall d \in D \cdot \bigsqcup (D \setminus \{d\}) \sqsubset s$$

# Join Decompositions

Given $s = \{a, b, c\}$:

✓ $D = \{\{a, b, c\}\}$

✗ $D = \{\{b\}, \{c\}\}$

✗ $D = \{\{a, b\}, \{b\}, \{d\}\}$

✓ $D = \{\{a, b\}, \{c\}\}$

✓ $D = \{\{a\}, \{b\}, \{c\}\}$

(This still allows join reducible elements, e.g. $\{a, b\}$)

# Join-Irreducible Decompositions

&#10007; $D = \{\{a, b, c\}\}$

&#10007; $D = \{\{a, b\}, \{c\}\}$

&#10003; $D = \{\{a\}, \{b\}, \{c\}\}$

# Join-Irreducible Decompositions

$$\text{GCounter} \quad \Downarrow m = \{\{i \mapsto n\} \mid i \mapsto n \in m\}$$
$$\text{GSet} \quad \Downarrow s = \{\{e\} \mid e \in s\}$$

# Differences and optimal deltas

We knew how to derive standard mutators from delta mutators

## Mutator from deltas
$$\mathsf{m}(x) = x \sqcup \mathsf{m}^\delta(x)$$

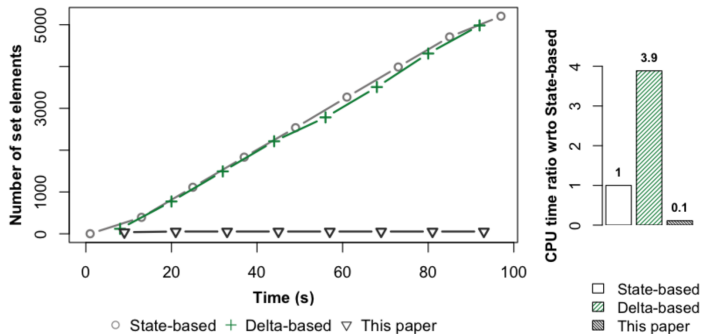Now we can derive optimal delta mutators from standard mutators

## Difference
$$\Delta(a, b) = \bigsqcup\{y \in {\Downarrow}a \mid y \not\sqsubseteq b\}$$

## Optimal delta mutator
$$\mathsf{m}^\delta(x) = \Delta(\mathsf{m}(x), x)$$

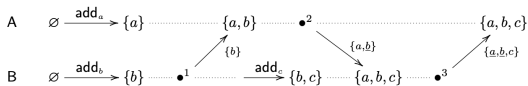# Improving anti-entropy

# Improving anti-entropy



Fig. 4: Delta-based synchronization of a GSet with 2 replicas $A, B \in \mathbb{I}$. Underlined elements represent the BP optimization.
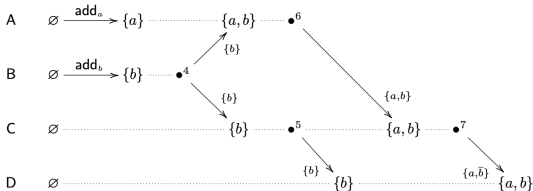


Fig. 5: Delta-based synchronization of a GSet with 4 replicas $A, B, C, D \in \mathbb{I}$. The overlined element represents the RR optimization.

# Improving anti-entropy

Channel can duplicate and reorder, but does not drop (missing Ack step)

**1 inputs:**
**2** $n_i \in \mathcal{P}(\mathbb{I})$, set of neighbors

**3 state:**
**4** $x_i \in \mathcal{L}$, $x_i^0 = \bot$
**5** $B_i \in \mathcal{P}(\mathcal{L})$, $B_i^0 = \varnothing$ $\quad\quad$ $B_i \in \mathcal{P}(\mathcal{L} \times \mathbb{I})$, $B_i^0 = \varnothing$

**6 on** operation$_i(\mathsf{m}^\delta)$
**7** $\delta = \mathsf{m}^\delta(x_i)$
**8** store$(\delta, i)$

**9 periodically** // synchronize
**10** **for** $j \in n_i$
**11** $d = \bigsqcup B_i$ $\quad\quad\quad\quad$ $d = \bigsqcup \{s \mid \langle s, o \rangle \in B_i \wedge o \neq j\}$
**12** send$_{i,j}$(delta, $d$)
**13** $B_i' = \varnothing$

**14 on** receive$_{j,i}$(delta, $d$)
**15** $\quad\quad\quad\quad\quad\quad\quad$ $d = \Delta(d, x_i)$
**16** **if** $d \not\sqsubseteq x_i$ $\quad\quad\quad\quad$ **if** $d \neq \bot$
**17** store$(d, j)$

**18 fun** store$(s, o)$
**19** $x_i' = x_i \sqcup s$
**20** $B_i' = B_i \cup \{s\}$ $\quad\quad\quad\quad$ $B_i' = B_i \cup \{\langle s, o \rangle\}$

**Algorithm 1:** Delta-based synchronization algorithms at replica $i \in \mathbb{I}$: classic version and version with BP and RR optimizations.
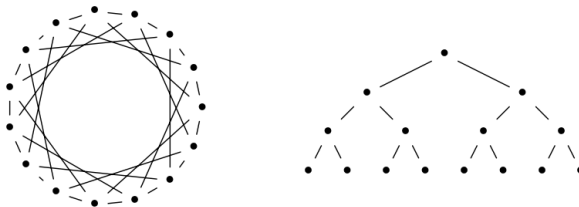
# Improving anti-entropy



Fig. 6: Network topologies employed: a 15-node partial-mesh (to the left) and a 15-node tree (to the right).

# Summary

- CAP trade-off and motivation for CRDTs
- Operations vs State
- State based CRDTs
- Delta state based CRDTs
- Decomposing states
- Update propagation

# Bibliography

▶ N. Preguiça, C. Baquero, M. Shapiro, "Conflict-free Replicated Data Types (CRDTs)," in Encyclopedia of Big Data Technologies, Springer International Publishing, 2018, https://doi.org/10.48550/arXiv.1805.06358

▶ P. Almeida, A. Shoker and C. Baquero, "Delta state replicated data types," in Journal of Parallel and Distributed Computing, Volume 111, 2018, Pages 162-173, ISSN 0743-7315, https://doi.org/10.1016/j.jpdc.2017.08.003

▶ V. Enes, P. Almeida, C. Baquero and J. Leitão, "Efficient Synchronization of State-Based CRDTs," in 2019 IEEE 35th International Conference on Data Engineering (ICDE), Macao, Macao, 2019 pp. 148-159. doi: 10.1109/ICDE.2019.00022

# Topic 1

### Asymmetric Merge Functions

- ▶ Deltas and classic states are combined by merge functions
- ▶ With state growth, deltas become relatively smaller that states
- ▶ Merge arguments commute, but assume *merge(big,small)*
- ▶ Program merge to be more efficient
- ▶ E.g. iterate on the appropriate state
- ▶ Build example data types and asymmetric merges
- ▶ Measure the complexity/performance gains

# Topic 2

### Topology Matters

- ▶ Deltas can improve data transmission of CRDT states
- ▶ Topology of node connections is important
- ▶ Ideally topology/overlay formation should be decentralized
- ▶ Topology can be static or dynamic
- ▶ As long as eventually updates reach all nodes
- ▶ Compare different overlay construction strategies
- ▶ Measure the impact on data transmission size and speed

# Topic 3

## Deflation in Isolation

- ▶ Mutations in state CRDTs are always inflations
- ▶ In PNCounters decrements are separately stored
- ▶ Could we increment by deflating increment count?
- ▶ Maybe, if we do it very privately
- ▶ Keep separate operations or deltas not yet communicated
- ▶ Explore mutual cancellations: `inc` vs `dec`; `add(x)` vs `rmv(x)`
- ▶ Measure potential metadata size improvements

# Practical

### Implementing State, Deltas, Join Decompositions, Differences

- ▶ You will use your language of choice
- ▶ Pick one or more state CRDTs
- ▶ Implement them with:
    - ▶ Both state and delta mutators
    - ▶ A classic join
    - ▶ A method/function for textual dumps
    - ▶ A method/function for extracting a set with state decompositions
    - ▶ A join that accepts a set of states
    - ▶ A method/function for computing the state difference