

Agenda

- Optimistic replication
 - CVS, Subversion
 - Duplicated databases
 - Operational transformation
 - Conflict-free Replicated Data Types (CRDT)

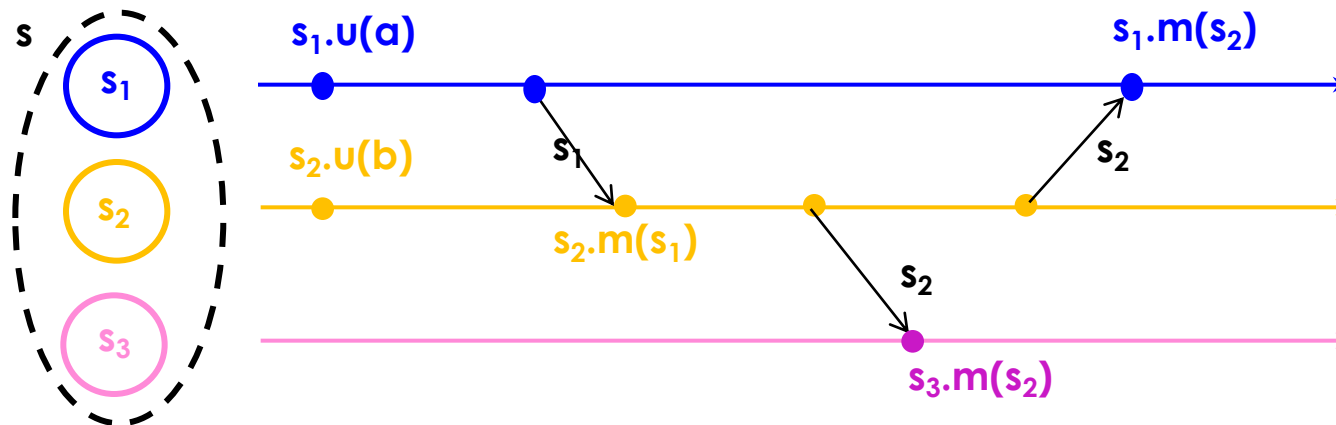
Conflict-free Replicated Data Types (CRDT) (*)

- Design operations to be commutative by construction
- Abstract data types
 - Designed to be replicated at multiple sites
 - Any replica can be modified without coordination
 - State convergence is guaranteed
- State-based and operation-based approaches

(*) M. Shapiro, N. Preguica, C. Baquero, M. Zawirski. A comprehensive study of Convergent and Commutative Replicated Data Types, Research Report, RR-7506, INRIA, 2011

Conflict-free Replicated Data Types

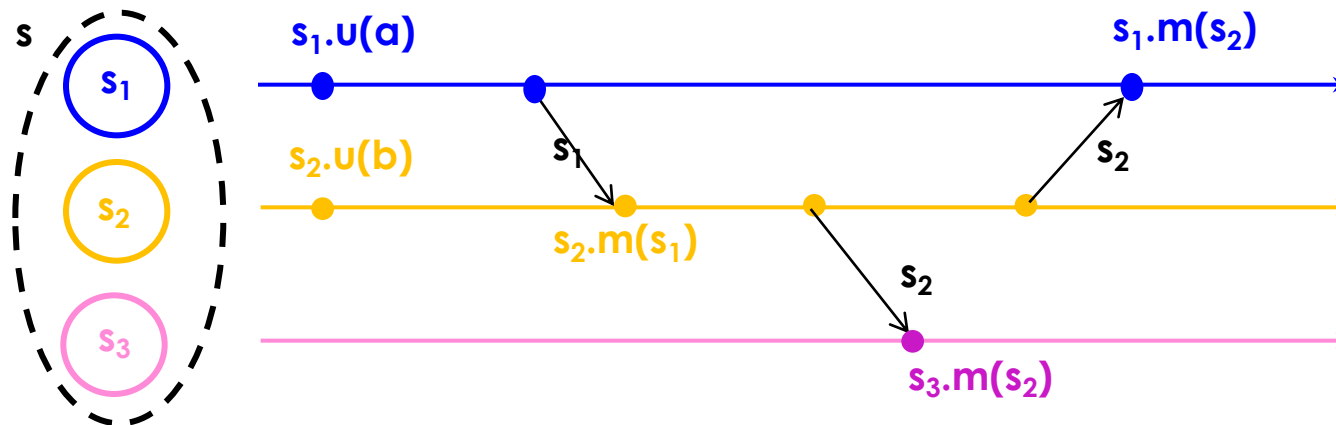
State-based Replication



- Replicated object: a tuple (S, s_0, q, u, m)
 - S : state domain
 - Replica at process p_i has state $s_i \in S$
 - s_0 : initial state
- Each replica can execute one of following commands
 - q : query object's state
 - u : update object's state
 - m : merge state from a remote replica

Conflict-free Replicated Data Types

State-based Replication



- Algorithm
 - Periodically, replica at p_i sends its current state to p_j
 - Replica p_j merges received state into its local state by executing m
- After receiving all updates (irrespective of order), each replica will have same state

Conflict-free Replicated Data Types

Semi-lattice

- Partial order \leq set S with a least upper bound (LUB), denoted \sqcup
 - $m = x \sqcup y$ is a LUB of $\{x, y\}$ under \leq if and only if
$$\forall m', x \leq m' \wedge y \leq m' \Rightarrow x \leq m \wedge y \leq m \wedge m \leq m'$$
- It follows that \sqcup is:
 - commutative: $x \sqcup y = y \sqcup x$
 - idempotent: $x \sqcup x = x$
 - associative: $(x \sqcup y) \sqcup z = x \sqcup (y \sqcup z)$

Conflict-free Replicated Data Types

Semi-lattice – Example on integers

- Partial order \leq on set of integers
- Least upper bound \sqcup : max (maximum function)
- Therefore, we have:
 - commutative: $\max(x, y) = \max(y, x)$
 - idempotent: $\max(x, x) = x$
 - associative: $\max(\max(x, y), z) = \max(x, \max(y, z))$

Conflict-free Replicated Data Types

Semi-lattice – Example on sets

- Partial order \subseteq on sets
- Least upper bound \sqcup : \cup (set union)
- Therefore, we have:
 - commutative: $A \cup B = B \cup A$
 - idempotent: $A \cup A = A$
 - associative: $(A \cup B) \cup C = A \cup (B \cup C)$

Conflict-free Replicated Data Types

Monotonic Semi-lattice Object

- A state-based object with partial order \leq , noted (S, \leq, s_0, q, u, m) , that has the following properties, is called a monotonic semi-lattice:
 1. Set S of values forms a semi-lattice ordered by \leq
 2. Merging state s with remote state s' computes the LUB of the two states, i.e., $s \bullet m(s') = s \sqcup s'$
(delivery order is not important)
 3. State is monotonically non-decreasing across updates, i.e., $s \leq s \bullet u$
(updates have effect, no rollback)

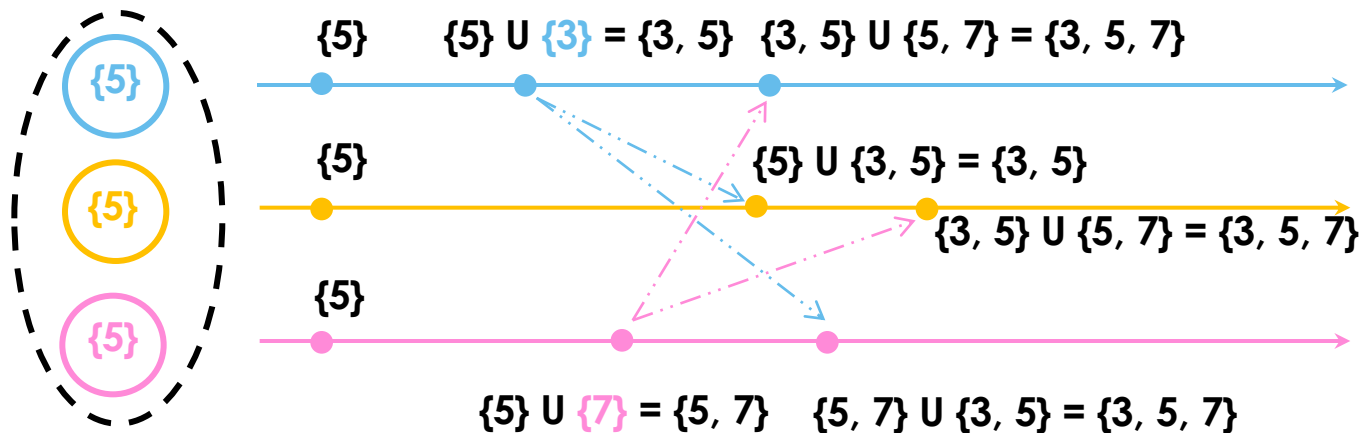
Conflict-free Replicated Data Types

Convergent Replicated Data Type (CvRDT)

- Theorem: Assuming eventual delivery and termination, any replicated state-based object that satisfies the monotonic semi-lattice property is SEC
- Since:
 - Merge is both commutative and associative
 - We do not care about order
 - Merge is idempotent
 - We do not care about delivering more than once

Convergent Replicated Data Types

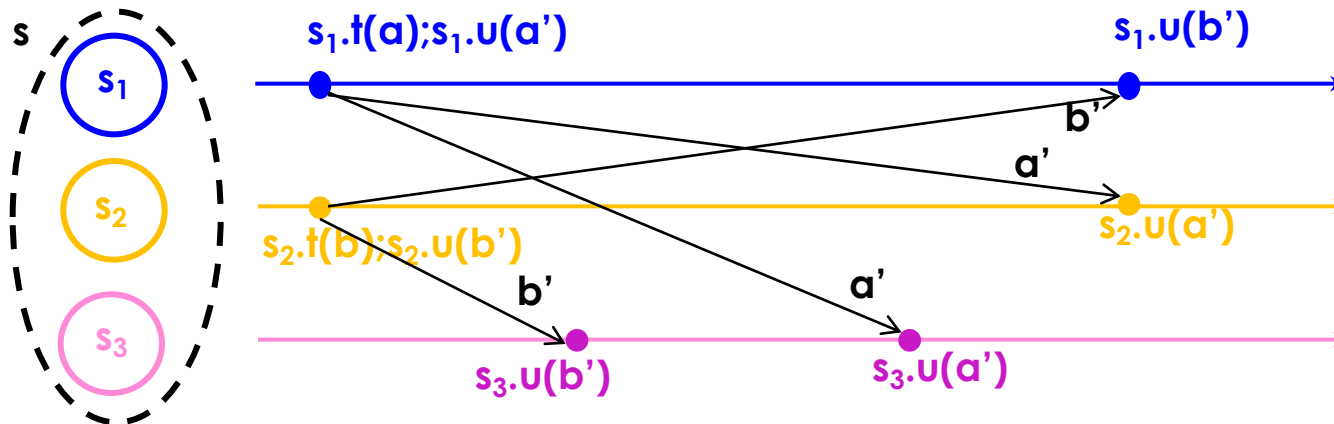
Example



- Each replica can execute one of following commands
 - query q : returns entire set
 - update u : adds new element (e, α) to local set
 - merge m : compute unions between local set and remote set

Conflict-free Replicated Data Types

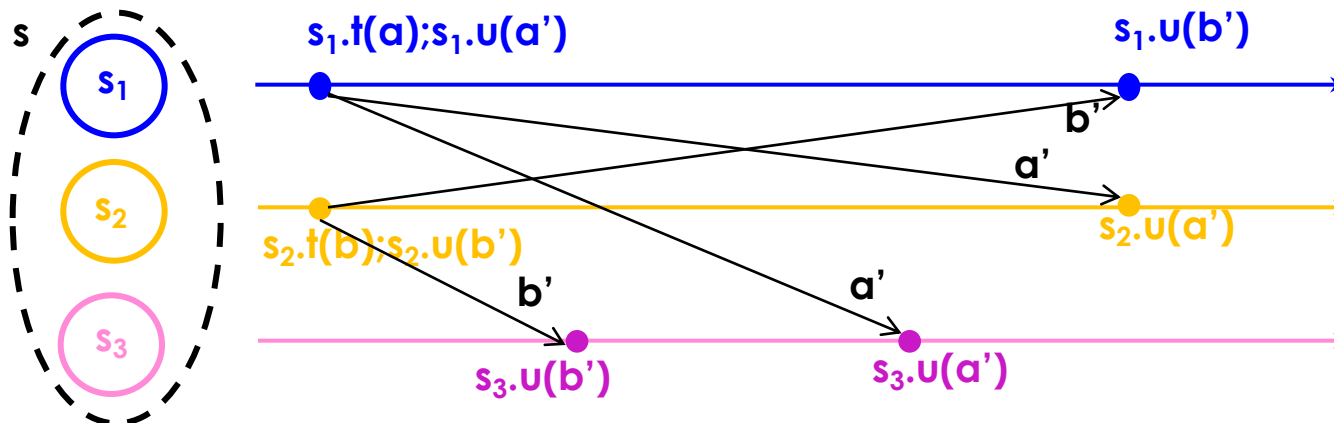
Operation-based Replication



- Replicated object: a tuple (S, s_0, q, t, u, P) .
 - S : state domain
 - Replica at process p_i has state $s_i \in S$
 - s_0 : initial state
- Each replica can execute one of following commands
 - q : query object's state
 - t : side-effect-free prepare-update method (at local replica)
 - u : effect-free update method (at all replicas)
 - P : delivery precondition

Conflict-free Replicated Data Types

Operation-based Replication



- Algorithm
 - Updates are delivered to all replicas
 - Use causally-ordered broadcast communication protocol, i.e., deliver every message to every node exactly once, w.r.t. happen-before order
 - Happen-before: updates from some replica are delivered to all recipients in the order they happened (effectively delivery precondition, P ; delivery of u at a replica may be delayed until P is true)
 - Note: concurrent updates can be delivered in any order

Conflict-free Replicated Data Types

Commutativity Property

- Updates (t, u) and (t', u') commute, if and only if for any reachable replica state s where both u and u' are enabled:
 - u (resp. u') remains enabled in state $s \cdot u'$ (resp. $s \cdot u$)
 - $s \cdot u \cdot u' \equiv s \cdot u' \cdot u$
- Commutativity holds for concurrent updates

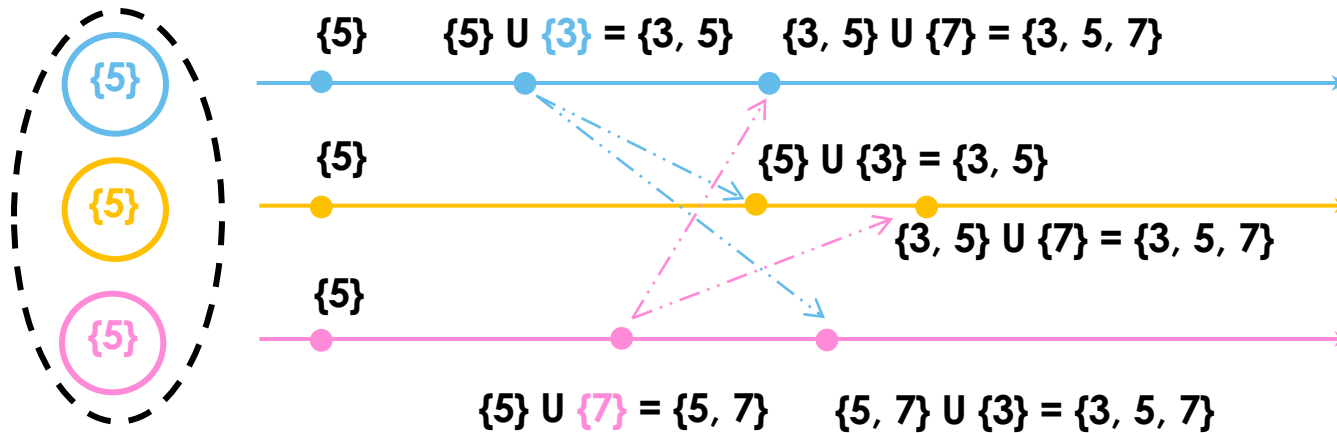
Conflict-free Replicated Data Types

Commutative Replicated Data Type (CmRDT)

- **Theorem:** Assuming causal delivery of updates and method termination, any replicated op-based object that satisfies the commutativity property for all concurrent updates is SEC

Commutative Replicated Data Types

Example



- query q : returns entire set
- prepare method t : adds new element (e, α) to local set
- update u : add delta to any remote replica

Consistency Maintenance

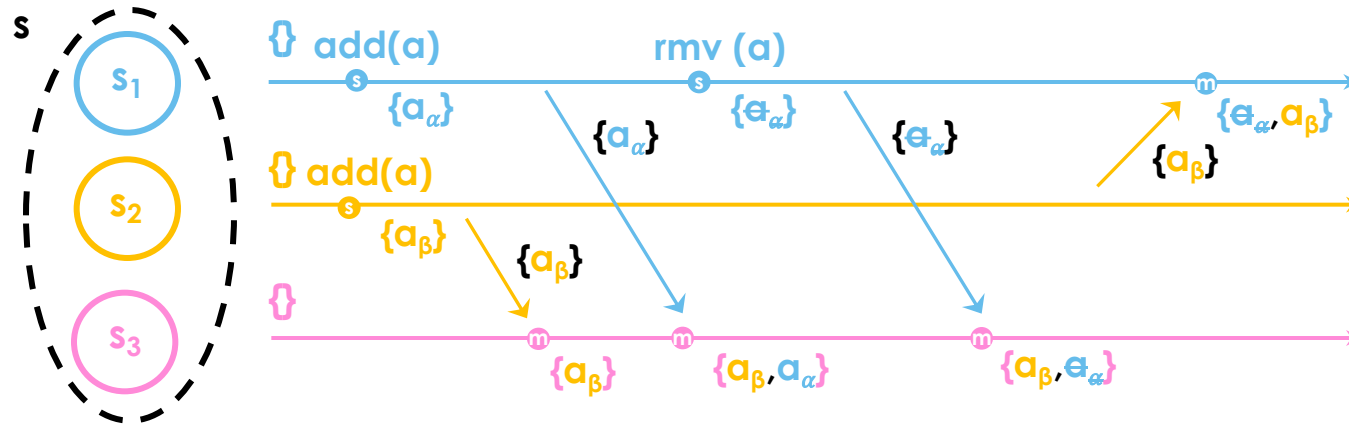
Conflict-free Replicated Data Types (CRDT)

- Register
 - Last-Writer Wins
 - Multi-Value
- Set
 - Grow-Only
 - 2-Phase
 - **Observed-Remove**
 - Observed-Update-Remove
- Map
- **Counter**
- Graph
 - Directed
 - Monotonic DAG
 - Edit graph
- **Sequence**



Conflict-free Replicated Data Types

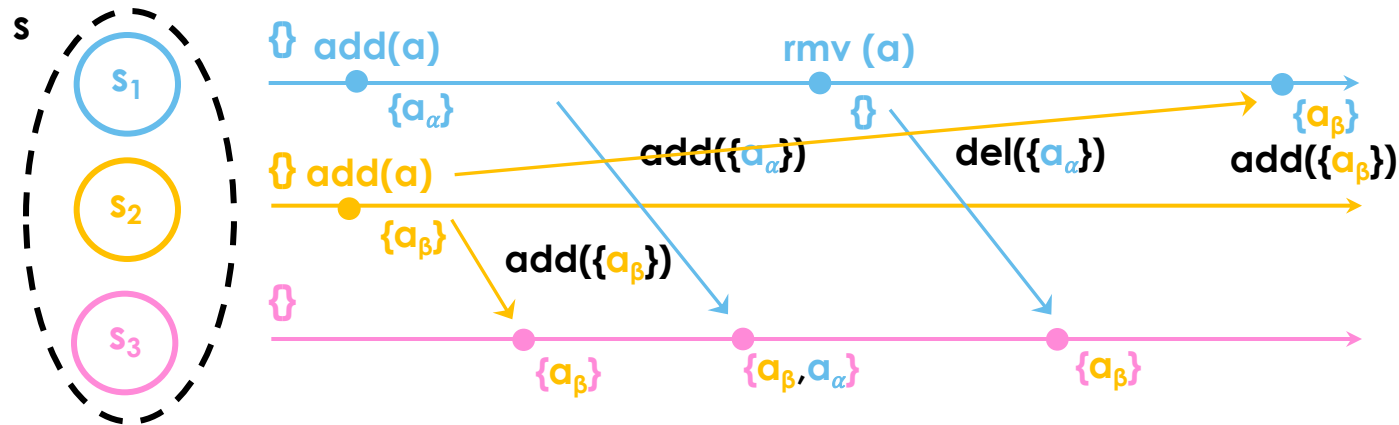
Observed-Remove Set (CvRDT)



- Payload: (A, R) - added/removed sets of $(element, unique-token)$
- Operations:
 - $add(e): A := A \cup \{(e, \alpha)\}$
 - $remove(e): R := R \cup \{(e, -) \in A\}$ remove all unique elements observed
 - $lookup(e): \exists (e, -) \in A \setminus R$
 - $merge(S, S'): (A \cup A', R \cup R')$
- $\{true\} add(e) \parallel remove(e) \{e \in S\}$

Conflict-free Replicated Data Types

Observed-Remove Set (CmRDT)



- Payload:
 $S = \{(e, \alpha), (e, \beta), (e, \gamma), \dots\}$ where $\alpha, \beta, \gamma, \dots$ are unique tokens
- Operations:
 $\text{add}(e): S := S \cup \{(e, \alpha)\}$ where α is a fresh unique token
 $\text{lookup}(e): \exists \alpha: (e, \alpha) \in S$
 $\text{remove}(e): R := \{(e, \alpha) \mid \exists \alpha (e, \alpha) \in S\}$ (at source) no tombstones
 $S := S \setminus R$
 $\{true\} \text{ add}(e) \parallel \text{remove}(e) \{e \in S\}$

Conflict-free Replicated Data Types

P-Counter (CvRDT)

- Payload:
 - $P = [\text{int}, \text{int}, \dots]$
- Operations:
 - $\text{value}(): \sum_i P[i]$
 - $\text{increment}(): P[\text{MyID}]++$
 - $\text{merge}(S, S'): S \sqcup S' = [\dots, \max(s.P[i], s'.P[i]), \dots]_i$
- Positive

Conflict-free Replicated Data Types

PN-Counter (CvRDT)

- Payload:
 - $P = [\text{int}, \text{int}, \dots]$,
 - $N = [\text{int}, \text{int}, \dots]$
- Operations:
 - $\text{value}()$: $\sum_i P[i] - \sum_i N[i]$
 - $\text{increment}()$: $P[\text{MyID}]++$
 - $\text{decrement}()$: $N[\text{MyID}]++$
 - $\text{merge}(S, S')$: $S \sqcup S' = ([\dots, \max(s.P[i], s'.P[i]), \dots]_i, [\dots, \max(s.N[i], s'.N[i]), \dots]_i)$
- Positive or negative

Conflict-free Replicated Data Types (CRDT)

CvRDT vs. CmRDT

- Both approaches are equivalent
 - A state-based object can emulate an operation-based object, and vice-versa
- **Operation-based:**
 - More efficient since you only ship small updates
 - But require exactly once causally-ordered broadcast
- **State-based:**
 - Only require reliable broadcast
 - Communication overhead of shipping the whole state
- **Delta State-based:**
 - Small messages
 - Dissemination over unreliable communication channels

Conflict-free replicated data type (CRDT) (Text) Sequence

- Document = linear sequence of elements
- Unique position identifiers
 - Each element has a unique identifier
 - Identifier remains constant for the lifetime of the document
 - Dense total order of identifiers consistent with element order:
$$\forall id_x, id_y: id_x < id_y \Rightarrow \exists id_z : id_x < id_z < id_y$$
- Real numbers require an infinite precision
- Different approaches for generating identifiers: Treedoc, Logoot

Logoot (*)

- Identifier:

$\langle p_1, s_1, h_1 \rangle . \langle p_2, s_2, h_2 \rangle \dots \langle p_k, s_k, h_k \rangle$

p_i integer

s_i site identifier

h_i logical clock at site s_i

$\langle 0, NA, NA \rangle$

$\langle 87, 1, 0 \rangle$

$\langle 87, 1, 0 \rangle . \langle 111, 6, 7 \rangle$

$\langle 89, 4, 5 \rangle$

$\langle MAX, NA, NA \rangle$

CSCW conference

March 19-23, 2011

Hangzhou, China

(*) Stéphane Weiss, Pascal Urso and Pascal Molli. *Logoot : a Scalable Optimistic Replication Algorithm for Collaborative Editing on P2P Networks*. In ICDCS , Montreal, Quebec, Canada , June 2009

Logoot

<0,NA,NA>

<1,1,1>.<5,2,4>.<4,6,2>

<4,2,1>.<2,4,6>.<3,3,6>

<6,3,4>

<MAX,NA,NA>

CSCW conference

March 19-23, 2011

Hangzhou, China

- Site 9 wants to insert 5 lines between first 2 lines
- Take first components <1,1,1> and <4,2,1>
 - 3 identifiers: <1,9,h>, <2,9,h> and <3,9,h>
- Take first 2 components <1,1,1>.<5,2,4> and <4,2,1>.<2,4,6>
 - 26 identifiers
 - <1,1,1>.<{5-9},9,h>
 - <2,9,h>.<{1-9},9,h>
 - <3,9,h>.<{1-9},9,h>
 - <4,2,1>.<{1},9,h>

Choose 5 identifiers

Logoot

<0,NA,NA>

<1,1,1>.<5,2,4>.<4,6,2>

<4,2,1>.<2,4,6>.<3,3,6>

<6,3,4>

<MAX,NA,NA>

CSCW conference

March 19-23, 2011

Hangzhou,China

- Remote insertion: Insert(<2,9,1>.<2,9,2>, “Computer Supported Cooperative Work”)
- Use binary search algorithm

<0,NA,NA>

<1,1,1>.<5,2,4>.<4,6,2>

<2,9,1>.<2,9,2>

<4,2,1>.<2,4,6>.<3,3,6>

<6,3,4>

<MAX,NA,NA>

CSCW conference

Computer Supported Cooperative Work

March 19-23, 2011

Hangzhou,China

Logoot

<0,NA,NA>

<1,1,1>.<5,2,4>.<4,6,2>

<2,9,1>.<2,9,2>

<4,2,1>.<2,4,6>.<3,3,6>

<6,3,4>

<MAX,NA,NA>

CSCW conference

Computer Supported Cooperative Work

March 19-23, 2011

Hangzhou, China

- Remote deletion: Delete(<4,2,1>.<2,4,6>.<3,3,6>)
- Use binary search algorithm

<0,NA,NA>

<1,1,1>.<5,2,4>.<4,6,2>

<2,9,1>.<2,9,2>

<6,3,4>

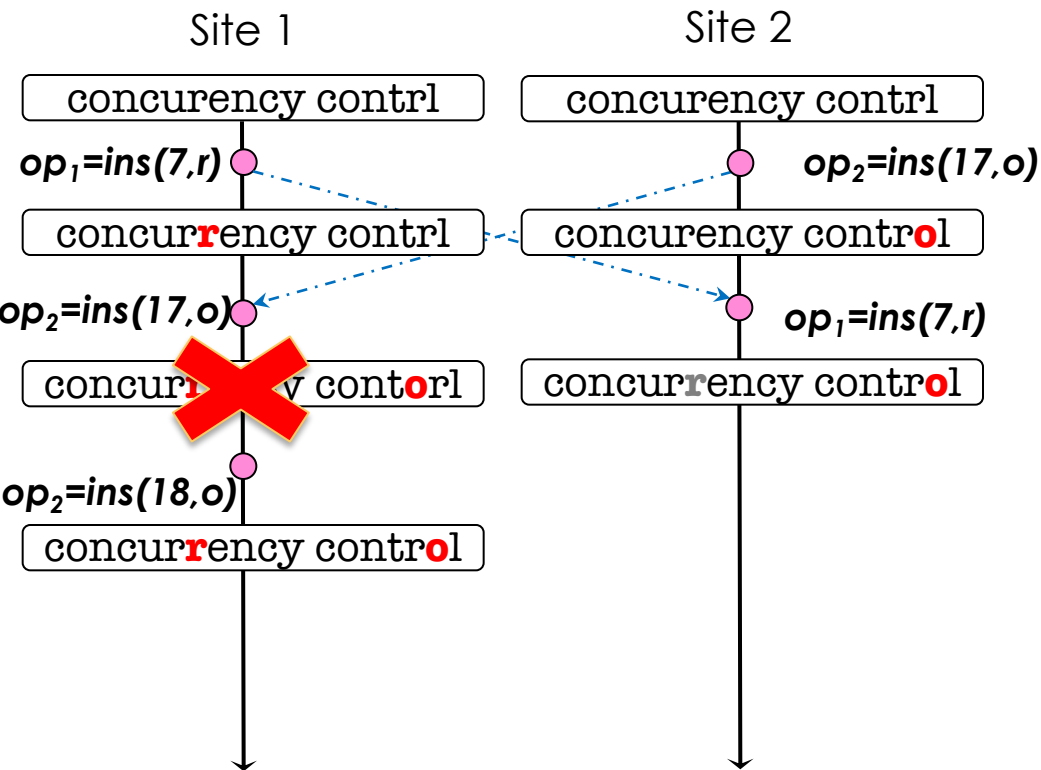
<MAX,NA,NA>

CSCW conference

Computer Supported Cooperative Work

Hangzhou, China

Operational Transformation (OT)



- Transforms non commuting operations to make them commute
- Genericity
- Time complexity
 - Average: $O(H.c)$ H : #ops
 - Worst case: $O(H^2)$ c : avg. #conc. ops
- Difficult to write correct transformation functions
- State vectors used for detecting concurrency \Rightarrow scalability limitations
- **Not very suitable for large scale peer-to-peer collaboration**

Conflict-free Replicated Data Types (CRDT)

- Time complexity

Average: $O(k \cdot \log(n))$

Worst case: $O(H \cdot \log(H))$

H: #ops

n: doc. size (non deleted chars.)

k: avg. size of Logoot identifier

- No need for concurrency detection

- Identifiers storage cost

- New design for each data type

- **Suitable for large-scale collaboration**

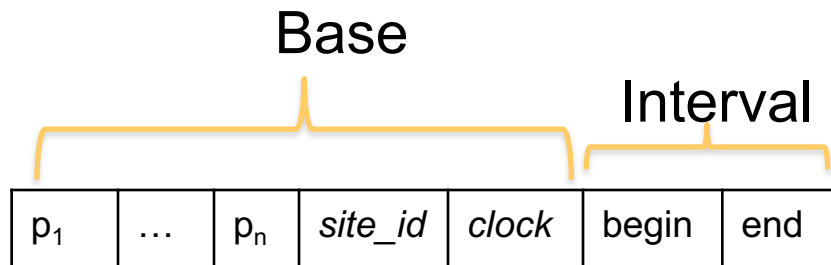
<1,2,1>	c
<1,2,2>	o
<2,1,2>	n
<3,1,3>	c
<3,1,3><8,4,5>	u
<3,2,5>	r
<4,1,7>	e
<4,1,7><9,2,6>	n
<7,2,8>	c
<9,1,7>	y
<10,2,8>	
<12,3,1>	c
<12,3,1><6,5,1>	o
<12,3,1><7,8,2>	n
<12,3,1><7,8,2><12,3,5>	t
<12,3,1><7,8,2><13,3,6>	r
<12,3,1><7,8,2><14,3,7>	l

ins(<3,2,5><13,1,7>, r)

ins(<12,3,1><7,8,2><13,3,6><7,2,9>, o)

Conflict-free Replicated Data Types (CRDT) LogootSplit

LogootSplit identifiers



1,1,[0,16] concurrency contrl

Insert r between "concur" and "ency contrl"

1,1,[0,5]	concur
1,1,5,2,1,[0,0]	r
1,1,[6,16]	ency contrl

Insert o between "ency contrl" and "l"

1,1,[0,5]	concur
1,1,5,2,1,[0,0]	r
1,1,[6,15]	ency contrl
1,1,15,3,1,[0,0]	o
1,1,[16,16]	l

André, L. et al., "Supporting Adaptable Granularity of Changes for Massive Scale Collaborative Editing" CollaborateCom 2013

LogootSplit

Site 3

ABCDEF

2,1,**[0,5]**



insert XY between B and C

Site 4

ABCDEF

2,1,**[0,5]**



LogootSplit

Site 3

ABCDEF

2,1,[0,5]



insert XY between B and C

AB

XY

CDEF

2,1,[0,1]

3,1

2,1,[2,5]

Site 4

ABCDEF

2,1,[0,5]



LogootSplit

Site 3

ABCDEF

2,1,[0,5]



insert XY between B and C

AB

2,1,[0,1]

XY

2,1,**1,3,1**,[0,1]

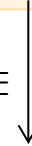
CDEF

2,1,[2,5]

Site 4

ABCDEF

2,1,[0,5]



insert ZT between D and E

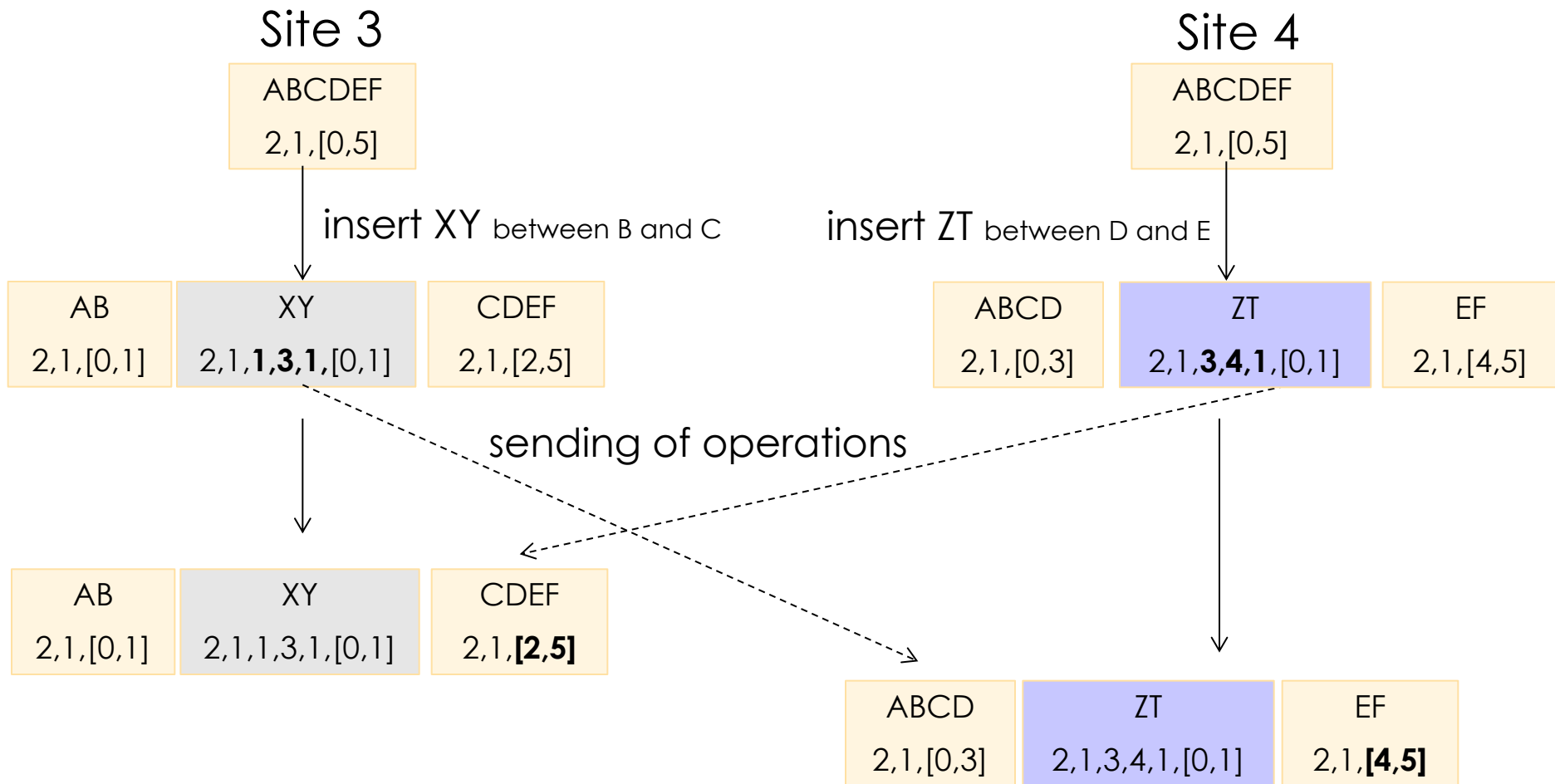
ABCD

2,1,[0,3]

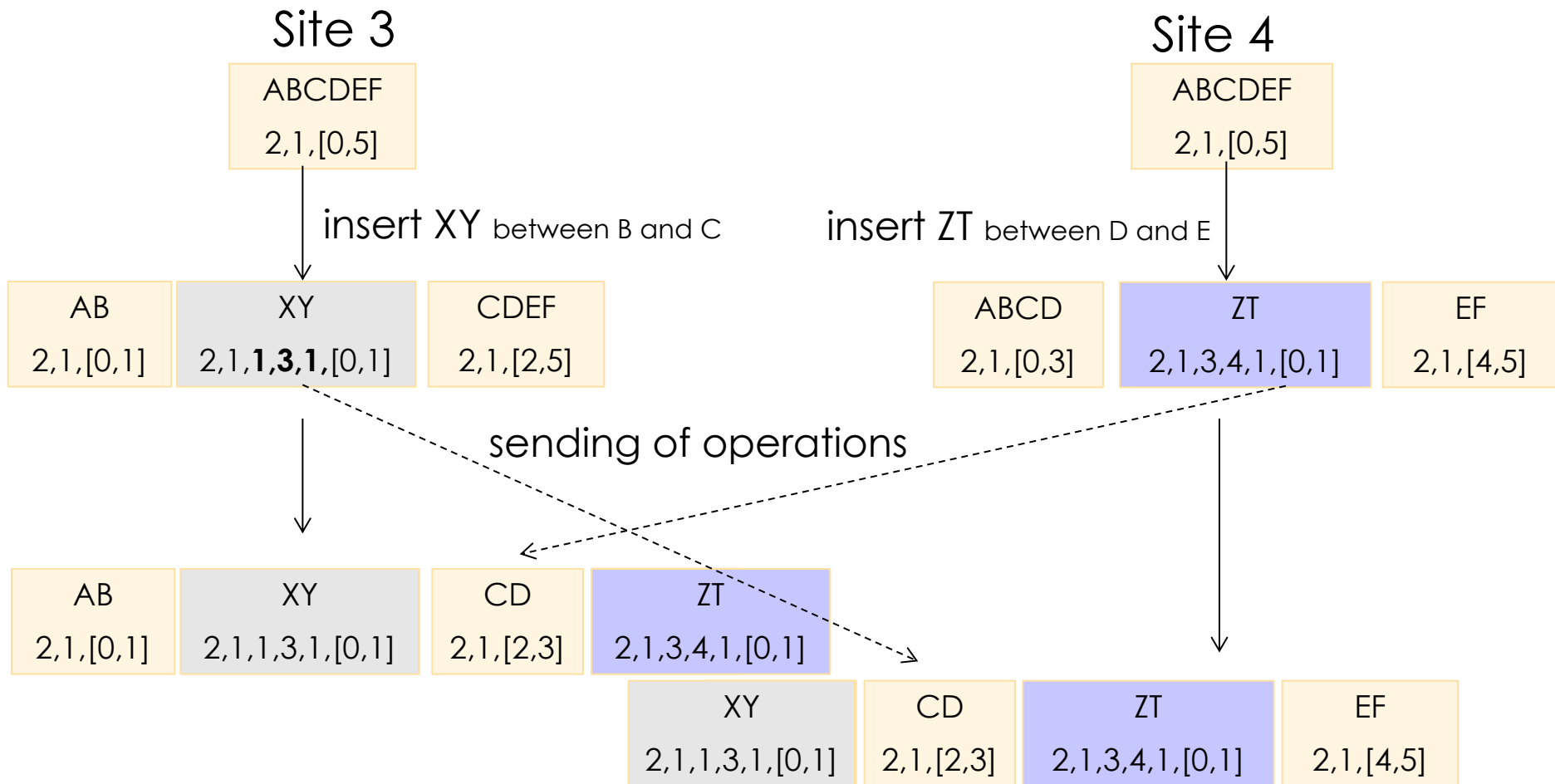
EF

2,1,[4,5]

LogootSplit

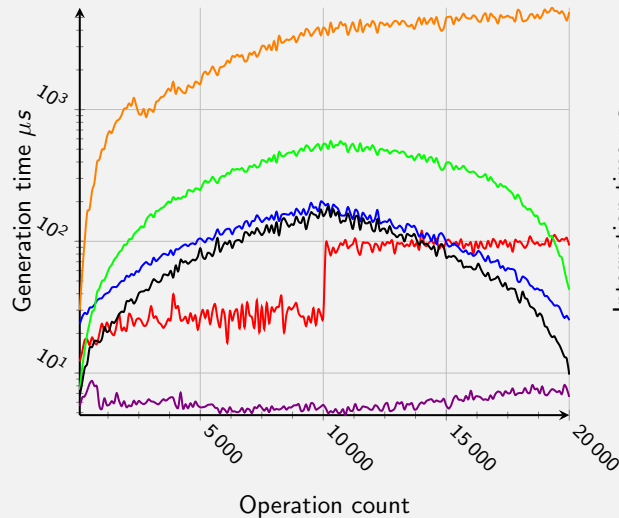


LogootSplit

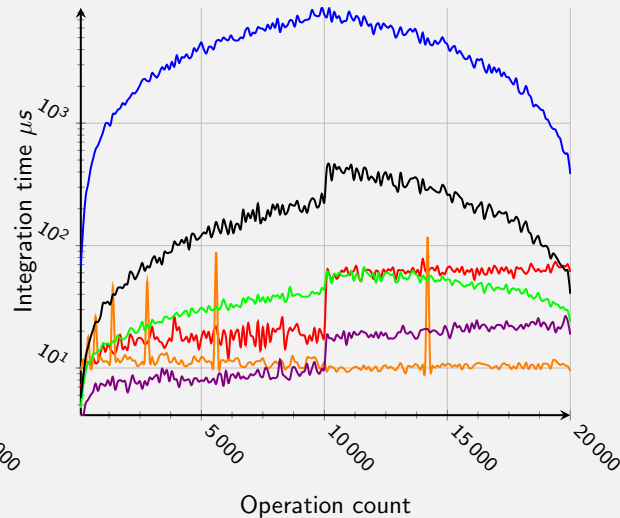


LogootSplit Performance Comparison

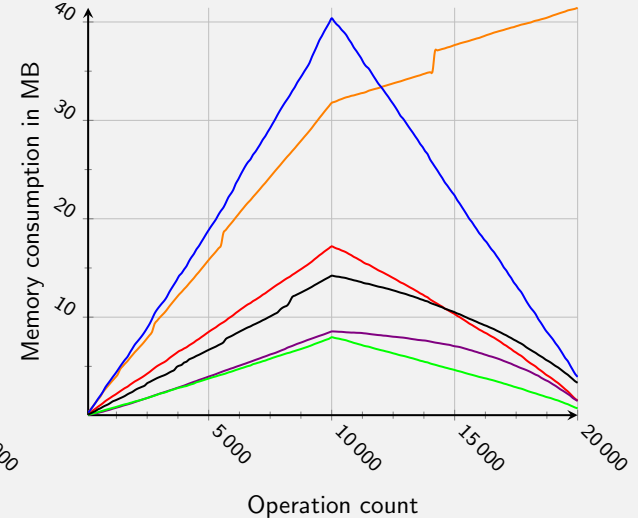
Treedoc ——— Woot ——— Logoot ———
LogootSplitAVL ——— LogootSplitNaive ——— LogootSplitString ———



Generation Time



Integration Time



Memory Consumption

Random block (avg. 50 char.) insertion/deletion

First 80% insertions, then 20% deletions

Delays in MUTE



Delays in GoogleDocs

