

# Feel the Rush: CRDTs in Riak

Zeeshan Lakhani

Software Engineer at Basho Technologies, Inc | Founder/Organizer Papers We Love  
@zeeshanlakhani

5-22-2015 (LambdaConf)

# Partial Order of Events

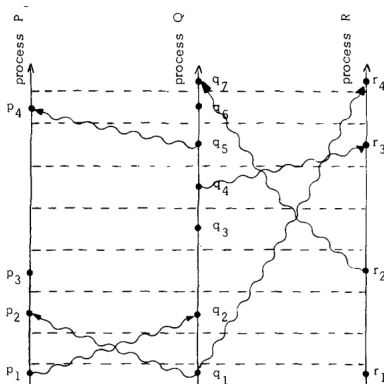
- Time, Clocks, and the Ordering of Events never gets old
- Can only release say that something **happened before**<sup>1</sup>
- Events
  - q1 happens before p2 and before p3
  - p3 and q2 are concurrent -  $p3 \parallel q2$ <sup>2</sup>
  - (figure on next slide)

---

<sup>1</sup>Time, Clocks, and the Ordering of Events in a Distributed System [Lamport]  
<http://bit.ly/1kQ3xhL>

<sup>2</sup>A Brief History of Time in Riak [Cribbs] <http://bit.ly/1Apto80>

# Lamport's "space-time"/process diagram<sup>3</sup>



<sup>3</sup>Time, Clocks, and the Ordering of Events in a Distributed System [Lamport]  
<http://bit.ly/1kQ3xhL>

# Replication + Coordination (CRDTs)

- A kind of safety called **Strong Eventual Consistency**
- Data Structures reducing the need for coordination

# State vs Operations-Based<sup>6</sup>

- CvRDT (convergent)
  - apply change locally; propagate entire state<sup>4</sup>
  - partial order to values
  - grows state monotonically
- CmRDT (commutative)
  - needs reliable broadcast to guarantee operations are delivered in partial order<sup>5</sup>
  - partial order to operations
  - replicas received updates

---

<sup>4</sup>Consistency and Riak [Meiklejohn] <http://bit.ly/1PBykxQ>

<sup>5</sup>CRDT Notes [pfraze] <http://bit.ly/1dpxMdn>

<sup>6</sup>A comprehensive study of Convergent and Commutative Replicated Data Types [Shapiro, et al.] <http://bit.ly/1PBC4zc>

# "Hello World" of a distributed ordering/replication issue

- When two writes occur concurrently, the next read returns their union. Concurrent updates even on unrelated elements, a remove may be undone.<sup>7</sup>
- Using (semantic reconciliation) mechanism, an 'add to cart' operation is never lost. However, deleted items can resurface."<sup>8</sup>

---

<sup>7</sup>An Optimized Conflict-free Replicated Set [Bieniusa, et al.]

<http://bit.ly/1ITUF48>

<sup>8</sup>Dynamo [DeCandia, et al.] <http://bit.ly/13QWj5Y>

# Foundations - Lattice

A bounded **join-semilattice**<sup>9</sup> :

$$\langle S, \sqcup, \perp \rangle$$

smallest element in  $S$ :

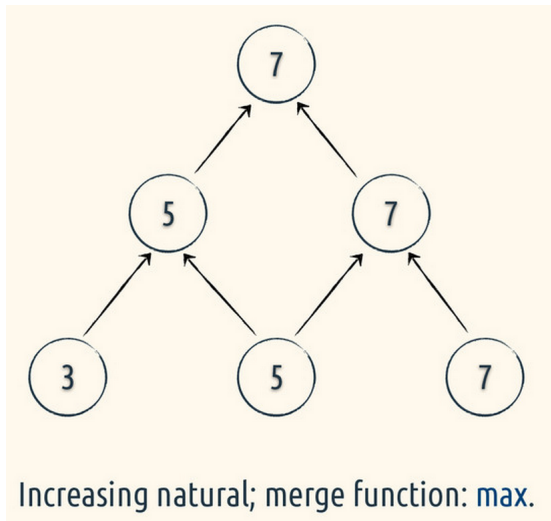
$$\perp$$

binary merge operator:

$$\sqcup$$

---

<sup>9</sup>Logic and Lattices for Distributed Programming [Conway, et al.]  
<http://bit.ly/1IQ6ppV>



<sup>10</sup>Logic and Lattices for Distributed Programming [Conway, et al.]  
<http://bit.ly/1IQ6ppV>



# Foundations - LUB

- a merge **function/operation** that produces a Least Upper Bound (LUB) over a join-semilattice
  - Replica A has a Haskell book - insert into cart
  - Replica B has a Scheme book - insert into cart
  - LUB would be the smallest cart state that's greater than or equal to both elements in the ordering:  $\{A, B\}$ <sup>12</sup>
  - If LUB exists, then it's unique & conflict resolution is deterministic<sup>12</sup>

---

<sup>12</sup>What's the deal w/ lvars and crdts [Kuper] <http://bit.ly/1Kpz4Cq>

## 2P-Set Example<sup>13</sup>

```
{  
  'type': '2p-set',  
  'a': ['a', 'b'],  
  'r': ['b']  
}
```

- only **a** exists

---

<sup>13</sup>meangirls [Kingsbury] <http://bit.ly/1cSrWQQ>

# Regular Or-Set Example<sup>14</sup>

```
{  
  'type': 'or-set',  
  'e': [  
    ['a', [1]],  
    ['b', [1], [1]],  
    ['c', [1, 2], [2, 3]]  
  ]  
}
```

- unique tags associated w/ each assertion
- **a** exists
- **b** != exist ... insertion was deleted,
- **c** exists ... two insertions, one deleted

---

<sup>14</sup>meangirls [Kingsbury] <http://bit.ly/1cSrWQQ>

# Comparing Vector Versions - (VVs) track divergence

## COMPARING VERSION VECTORS

- Descends :  $A \geq B$
- Dominates :  $A > B$
- Concurrent :  $A \mid B$



# Descends, Dominates, Concurrent<sup>15</sup>

- descends  $A \geq B$ 
  - A summarizes at least the same history as B (as seen all the events in B)
- dominates  $A > B$  (good for discarding)
  - A is strictly greater than B (has seen all events B and at least 1 more event)
- concurrent  $A \parallel B$ 
  - A contains at least 1 event unseen by B and B contains at least 1 event unseen by A

---

<sup>15</sup>A Brief History of Time in Riak [Cribbs] <http://bit.ly/1Apto80>

# Riak's implementation - ORSWOT

- Optimized Observe-Remove Set<sup>16, 17</sup>
- Two-Way Comparison
- E.g. 2 Replicas **Merge**. . . one contains element in the set, other does not

---

<sup>16</sup>An Optimized Conflict-free Replicated Set [Bieniusa, et al.]

<http://bit.ly/1ITUF48>

<sup>17</sup>`riak_dt_orswot.erl` [Russell Brown] <http://bit.ly/1ca8Wgb>

# Causality to the Rescue<sup>18</sup>

## Sets as (dotted) version vectors

$[\{\text{actor1, count}\}=\text{Dot1}, \{\text{actor2, count}\}=\text{Dot2}, \dots]$  - minimal clock

*We take all the elements that are in Set A and not in Set B and compare their minimal clocks to Set B's set version vector. Every element whose minimal clock is dominated has been removed from Set B, does not make it into the merged set. Also drop **dots** dominated by Set B's clock<sup>a</sup>*

- Then repeat the other way (compare to Set A's VV)

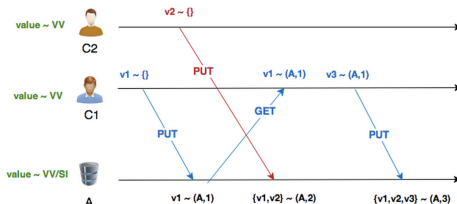
---

<sup>a</sup>[Russell Brown] <http://bit.ly/1AptLzI>

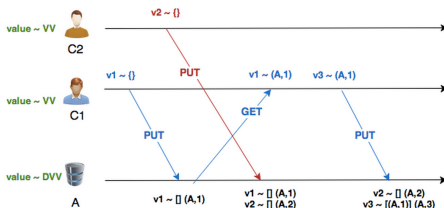
---

<sup>18</sup>[Russell Brown] <http://bit.ly/1AptLzI>

- VVs merge



- DVV





# Riak Data Types<sup>20, 21</sup>

- Maps - riak\_dt\_map
- Sets - riak\_dt\_orswot
- Registers - riak\_dt\_lwwreg
- Flags - riak\_dt\_od\_flag
- Counters - riak\_dt\_emcntr

---

<sup>20</sup>Riak DT [Russell Brown] <http://bit.ly/1Eoj0J4>

<sup>21</sup>Riak Data Types <http://bit.ly/1Q1S7RG>

# Erlang client-side Example

22

```
Map4 = riakc_map:update({<<"interests">>, set},  
    fun(S) ->  
        riakc_set:add_element(<<"robots">>, S)  
    end,  
    Map3),  
Map5 = riakc_map:update({<<"interests">>, set},  
    fun(S) ->  
        riakc_set:add_element(<<"opera">>, S)  
    end,  
    Map4).
```

---

<sup>22</sup>Riak Data Types Map <http://bit.ly/1BhfnEI>