

RainbowFS

Modular Consistency and Co-designed Massive File system

M. Shapiro, UPMC & Inria, Paris
 B. King, Scality SA, Paris
 V. Quéma, CNRS-LIG, Grenoble
 P. Sutra, Télécom ParisSud, Évry
 S. Monnet, U. Savoie-Mont Blanc, Annecy



Background

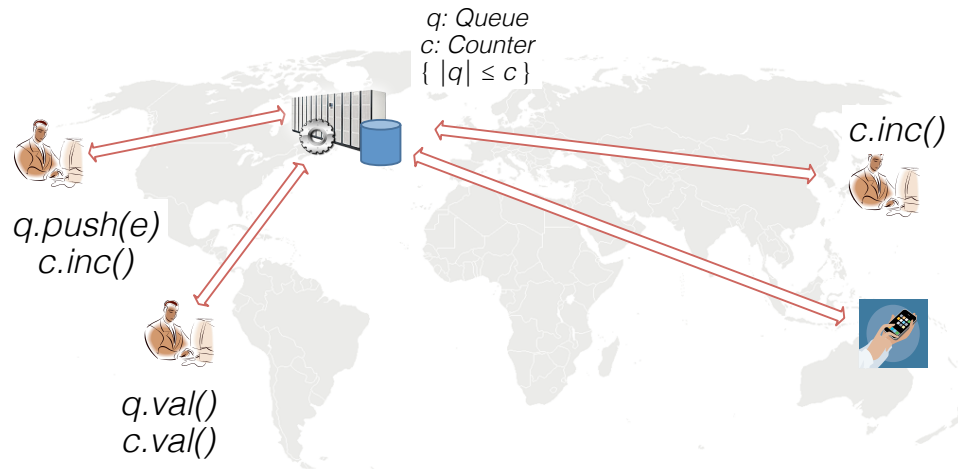
Distributed application \Rightarrow shared data

1. Large scale (cloud) \Rightarrow data replication \Rightarrow consistency issues
 - Strong : dependable, not available, inefficient.
 - Available : parallel, anomalies.
 - Pre-defined models
 2. Complex : configuration, analysis, control, decomposition
 3. System, data access API
- Theory of replication and consistency, tools

[RainbowFS]

2

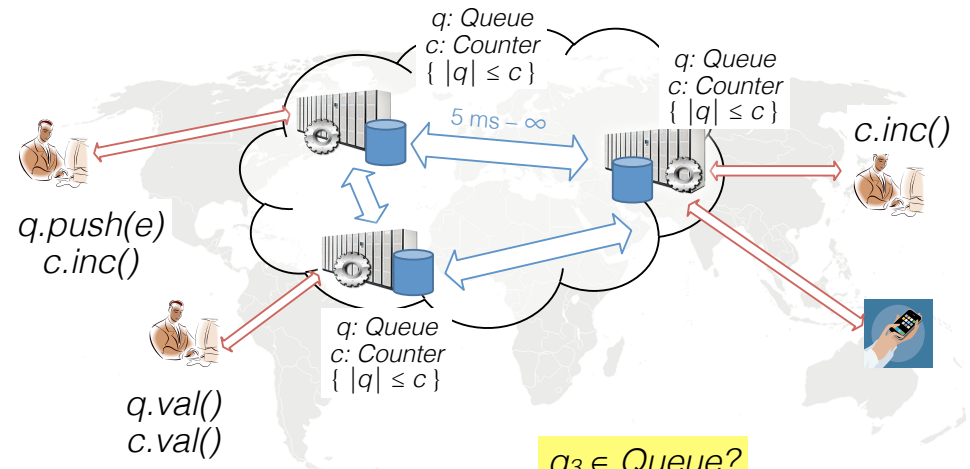
Shared data



[RainbowFS]

3

Geo-replicated sharing

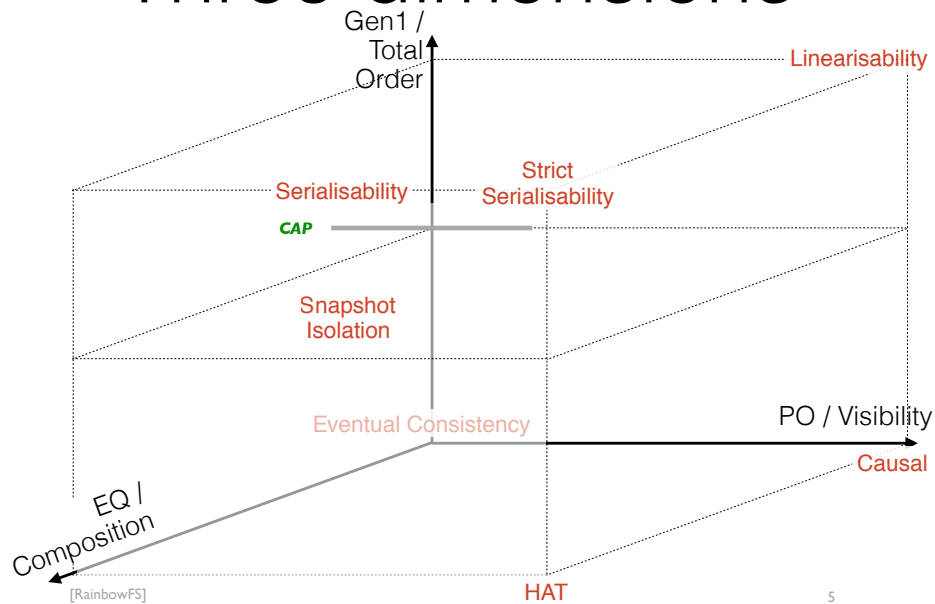


$q_3 \in \text{Queue?}$
 $q_1 = q_2 ?$
 $|q_1| \leq c_4 ?$

[RainbowFS]

4

Three dimensions



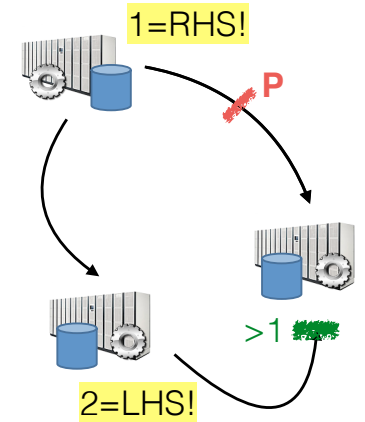
Relative order pattern

mkdir before *creat*

Relative-order invariant pattern:

- “Directory references valid file”
- $x \text{ valid} \wedge x \text{ points to } y \Rightarrow y \text{ valid}$
- Pattern RHS!; LHS!
- Make visible in same order

AP-compatible: Causal Consistency



Joint update pattern

write updates content, metadata

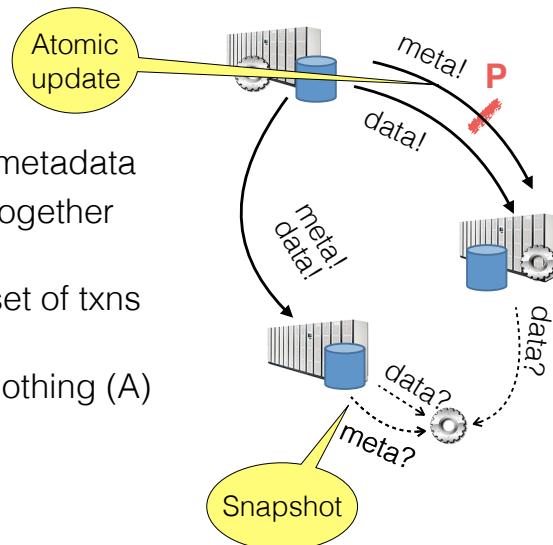
Transmit joint updates together

- write-atomic

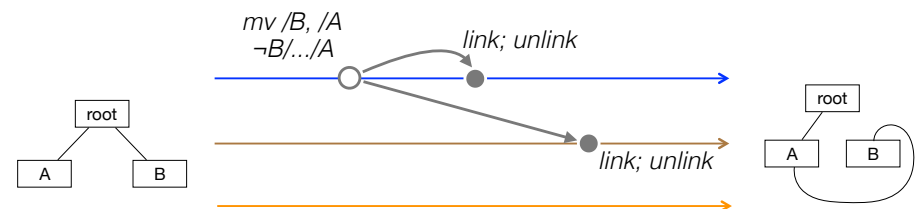
+ Read from common set of txns

- snapshot property

AP-compatible: All-or-Nothing (A)



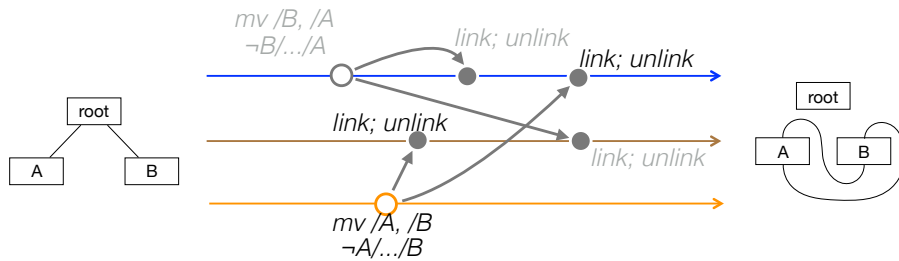
CAP-sensitive pattern



```

mv (node, ddir) { // tree
  if  $\rightarrow$  ancestor (node, ddir) // precondition: at source
    link (ddir, node) // at every ...
    unlink (sdir, node) // ... replica
} // tree
    
```

CAP-sensitive pattern



```

mv (node, ddir) { // tree
  if  $\neg$  ancestor (node, ddir) // precondition: at source
    link (ddir, node) // at every ...
    unlink (sdir, node) // ... replica
} // tree

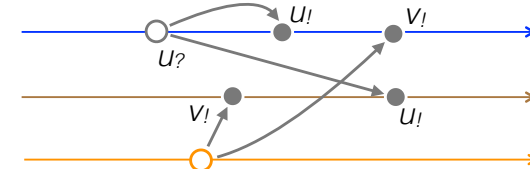
```

- Precondition *not stable* w.r.t. concurrent *mv*
- Forbid concurrency? Synchro, CP.
 - Or remove invariant? AP, degraded semantics

[RainbowFS]

9

A data model for AP



Concurrent, asynchronous updates

- Standard register model: assignments \Rightarrow CP
- AP \Rightarrow concurrent updates merged

CRDT: register, counter, set, map, sequence

- Extends sequential type
- Encapsulates convergent merge

[RainbowFS]

10

geoFS

Scality—UPMC

Tao Thanh Vinh

Posix-like replicated file system

- Available under Partition
- CRDTs: merge concurrent updates
- \parallel *write* same file: merge or rename
- *delete* \parallel *write*: file path survives
- *mv* \parallel *mv* \rightarrow copy; delete

mv \parallel *mv* \Rightarrow unstable precondition

- Either CP or anomalous
- CP: minimal synch footprint

Mahsa Najafzadeh

[RainbowFS]

11

RainbowFS

ANR project 2017–2021

UPMC + LIG + TSP + Scality + USMB

1. Application/system co-design
2. Modular replication
3. Multi-consistency file system

[RainbowFS]

12

1. Application / consistency co-design

Just-Right Consistency:

- Most efficient consistency...
- ...under which *my* application is correct

Static and dynamic verification

File system:

- Maintain tree invariant; others?
- Asynchronous: *mkdir*, *rmdir*, *creat*, *rm*, etc.
- Design options: *write* || *write*, *write* || *rm*
- Synchronised: *mv*

[RainbowFS]

13

3. Geo-replicated massive file system

Build a file system:

- Posix API
- Multi-consistency, tailor to application
- Scales to O(petabyte)
- Apply tools & methods from Tasks 1 & 2

Challenges:

- Geo-distributed, elastic
- Massive performance, partial replication
- Consistency, security, fault tolerance
- Layer above object store

[RainbowFS]

15

2. Modular geo-replication

From application skeleton to running system

Tools for large-scale deployment, monitoring, analysis, diagnosis

- Check correctness & efficiency
- Stress test
- Diagnose root cause

Protocol building blocks, compose for application

- 3D model, variants/strength
- Modular fault tolerance (Abstract)

[RainbowFS]

14



Creative Commons Attribution-ShareAlike 4.0 Intl. License

You are free to:

- Share — copy and redistribute the material in any medium or format
- Adapt — remix, transform, and build upon the material

for any purpose, even commercially, under the following terms:



Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



ShareAlike — If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.

[RainbowFS]

16