Causality for the Masses:

Offering Fresh Data, Low Latency, and High Throughput

Luís Rodrigues





Causality for the Masses:

Offering Fresh Data, Low Latency, and High Throughput

Manuel Bravo,

Luís Rodrigues, Chathuri Gunawardhana, Peter van Roy



UCL Université catholique de Louvain

Limitations of Highly-Available Eventually-Consistent Data Stores

Hagit Attiya Computer Science Department Technion

Faith Ellen Department of Computer Science University of Teronto

Adam Morrison ce Computer Science Department Technion

ABSTRACT

Modern replicated data stores aim to provide high availability, by immediately responding to client requests, often by implementing objects that expose concurrency. Such objects, for example, multi-valued registers (MVRs), do not have sequential specifications. This paper explores a recent model for replicated data stores that can be used to precisely specily cousal consistency for such objects, and liveness properties like eventual consistency, without revealing details of the underlying implementation. The model is used to prove the following results:

- An eventually consistent data store implementing MVRs cannot satisfy a consistency model strictly stronger than observable causal consistency (OCC). OCC is a model somewhat stronger than causal consistency, which captures executions in which client observations can use causality to infer concurrency of operations. This result holds under certain assumptions about the data store.
- Under the same assumptions, an eventually consistent and causally consistent replicated data store must send messages of unbounded size: If s objects are supported by n replicas, then, for every k > 1, there is an execution in which an Ω(min{n, s}k)-bit message is sent.

of data (i.e. accesses to data return without delay), and its consistency, while tolerating message delays. The CAP theorem [8,18] demonstrates the difficulty of achieving this balance, showing that strong Consistency (i.e. atomicity) cannot be satisfied together with high Availability and Partition tolerance.

One aspect of data consistency is a safety property restricting the possible values observed by clients accessing different replicas. The set of possible executions is called a consistency model. For example, cousal consistency [2] ensures that the causes of an operation are visible at a replica no later than the operation itself. (A precise definition appears in Section 3.) A smaller set of possible executions means that there is less uncertainty about the data. So, one consistency model is strictly stronger than another if its executions are a proper subset of the executions of the other.

Some weak consistency models can be trivially satisfied by never updating the data. Therefore, another aspect of data consistency is a liveness property, ensuring that updates are applied at all replicas. The designers of many systems, e.g., Dynamo [13] and Cassandra [1], opt for a very weak liveness property called, somewhat confusingly, eventual consistency [5, 9, 10, 29]. Eventual consistency only ensures that each replica eventually observes all updates to the object, a property also referred to as update propagation [11]. Strongest without compromising availability

Parallel Snapshot Isolation [SOSP'11]

Key ingredient of several consistency criteria

RedBlue Consistency [OSDI'12]

Explicit Consistency

[EuroSys'15]

Session guarantees

[SOSP'97]































Ţ,

Data center should delay the visibility of inconsistent operations























Problem

If causal dependencies are not accurately tracked metadata may generate false positives

False dependencies!




















more metadata

less metadata



more metadata

less metadata





less metadata

One vector per item. One entry in each vector per DC.

more metadata

less metadata



Lamport's clocks

more metadata





expensive



Lamport's clocks

more metadata

less metadata



Problems of the previous state-of-the-art Throughput vs. data staleness tradeoff

GentleRain [SoCC' 14]: Optimizes throughput Compresses metadata into a scalar

Cure [ICDCS' 16]: Optimizes data freshness Relies on a vector clock with an entry per data center



Problems of the previous state-of-the-art Throughput vs. data staleness tradeoff

GentleRain [SoCC' 14]: Optimizes throughput Compresses metadata into a scalar

Cure [ICDCS' 16]: Optimizes data freshness Relies on a vector clock with an entry per data center



Problems of the previous state-of-the-Throughput vs. data staleness tradeoff

GentleRain [SoCC' 14]: Optimizes throughput Compresses metadata into a scalar

Cure [ICDCS' 16]: Optimizes data freshness Relies on a vector clock with an entry per data center



False dependencies damage data freshness Problems of the previous state-of-the-art Throughput vs. data staleness tradeoff

Partial replication aggravates the problem

Visibility latencies have direct impact on client response latency





+

Saturn











God in ancient Roman religion, that become the god of time

╋





Saturn



Greek goddess of law and legislation

┿



Saturn



If this couple cannot fix the problem, nobody can...

┿



÷



Saturn







Saturn Eunomia Orders events across datacenters

+







Distributed metadata service

pluggable to existing geo-distributed data services

handles the dissemination of operations among data centers

Ensures that

clients always observe a causally consistent state

with a negligible performance overhead when compared to an eventually consistency system





Requires a **constant and small** amount of metadata regardless of the system's scale (servers, partitions, and locations)









Mitigates the impact of false dependencies

by relying on a tree-based dissemination









Implements **genuine partial replication** data centers only manage data and metadata of the items replicated locally





Decoupling data and metadata


Decoupling data and metadata



Decoupling data and metadata



Decoupling data and metadata































Metadata propagation b C causa a J a concurr

C causally depends on b A concurrent to both b and C



Metadata propagation C causally depends on b D а a concurrent to both b and c











b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C



b C causally depends on b a concurrent to both b and C






















Metadata propagation

Causal consistency is a partial order

Saturn exploits this fact **by serving possibly different linear extensions** of the causal order to each data center

Each served serialization aims at optimising data freshness, and thus, reducing the impact of false dependencies



Metadata propagation



Causal consistency is a partial order

Saturn exploits this fact **by serving possibly different linear extensions** of the causal order to each data center

Each served serialization aims at optimising data freshness, and thus, reducing the impact of false dependencies



Saturn leverages a set of cooperating servers, namely **serializers**, forming a tree

Causal consistency is trivially enforced by the tree assuming

FIFO links among serializers

Serializers propagate metadata preserving the observed order

Serializers are geographically distributed to optimize data freshness



Metadata dissemination graph





Metadata dissemination graph





Metadata dissemination graph





Optimal dissemination graph

The goal is to build the tree such that metadatapaths latencies (through the tree) **match** data-paths

Weighted Minimal Mismatch

mismatch_{i,j} =
$$|\Delta^M(i,j) - \Delta(i,j)|$$

 $min \sum_{\forall i,j \in V} c_{i,j} \cdot mismatch_{i,j}$



Optimal dissemination graph

The goal is to build the tree such that metadatapaths latencies (through the absolute difference between label-paths and data paths Weighted Minimal Misma $mismatch_{i,j} = \left|\Delta^M(i,j) - \Delta(i,j)\right|$ $min \sum_{\forall i, j \in V} c_{i,j} \cdot mismatch_{i,j}$



Optimal dissemination graph

The goal is to build the tree such that metadatapaths latencies (through the tree) **match** data-paths





Finding the optimal tree is modelled as a **constraint optimization** problem

Input

Data-paths average latencies Candidate locations for serializers (an latencies among them) Access-patterns: to minimize the impact of mismatches



b C causally depends on b a concurrent to both b and C





b C causally depends on b a concurrent to both b and C





b C causally depends on b a concurrent to both b and C





Time: 0

b C causally depends on b a concurrent to both b and C







Time: 02

b C causally depends on b a concurrent to both b and C

а





Time: 02

b C causally depends on b a concurrent to both b and C

а





Time: 02

b C causally depends on b a concurrent to both b and C

а





Time: 023

b C causally depends on b C a concurrent to both b and C

а





Time: 023

b C causally depends on b C a concurrent to both b and C

а





Time: 0234

b C causally depends on b a concurrent to both b and C

а





Time: 02345

b C causally depends on b a concurrent to both b and C

а





Time: 02345



D

а

C causally depends on b

a concurrent to both b and c



Time: 02345

C causally depends on b D a concurrent to both b and c





Time: 02345

C causally depends on b D a concurrent to both b and c





Time: 023456

C causally depends on b a concurrent to both b and c



D



Time: 0234567

C causally depends on b a concurrent to both b and c



D



C causally depends on b A concurrent to both b and C

Time: 0234567



D



C causally depends on b A concurrent to both b and C

Time: 0234567



D



C causally depends on b a concurrent to both b and C

Time: 0234567



D



b C causally depends on b a concurrent to both b and C

Time: 0234567...12





C causally depends on b A concurrent to both b and C

Time: 0234567...12



D



C causally depends on b A concurrent to both b and C

Time: 0234567...12



D


C causally depends on b A concurrent to both b and C

Time: 0234567...12



D



C causally depends on b A concurrent to both b and C

Time: 0234567...12



D



b C causally depends on b a concurrent to both b and C

Time: 0234567...12...14





b C causally depends on b a concurrent to both b and C

Time: 0234567...12...14





b C causally depends on b a concurrent to both b and C

Time: 0234567...12...14





C causally depends on b A concurrent to both b and C

Time: 0234567...12...14



b





Time: 0234567...12...14...16



b





Time: 0234567...12...14...16



b





Time: 0234567...12...14...16



b





Time: 0234567...12...14...16



b











each datacenter has many machines

how to serialize updates performed on different machines without loosing performance?

usually one relies on **sequencers** that limit concurrency





conceived to **replace sequencers** used to serialize updates in replicated storage systems

totally orders—consistently with causality local updates, before shipping them to Saturn and other dcs

the ordering is done in the background, out of client's critical path





































Evaluation

+





Evaluation: goals

- Can Saturn+Eunomia optimize both throughput and data freshness simultaneously?
- Why **Saturn's** genuine partial replication matters?

• How important is the configuration of **Saturn**?

• How **Eunomia** compares to sequencers?

Evaluation: goals

- Can Saturn+Eunomia optimize both throughput and data freshness simultaneously?
- Why **Saturn's** genuine partial replication matters?

• How important is the configuration of **Saturn**?

• How **Eunomia** compares to sequencers?
Saturn vs state-of-the-art

Public Facebook dataset

[B. Viswanath et al. On the evolution of user interaction in Facebook. WOSN '09]

Partitioning among 7 data centers

Replication factor min. 2 max. [2,5]

[J. M. Pujol et al. The little engine(s) that could: Scaling online social networks. SIGCOMM '10]

Workload based on real social network workloads

[F. Benevenuto et al. Characterizing user behavior in online social networks. ICM '09]

Saturn vs state-of-the-art Throughput



max replication degree

Saturn vs state-of-the-art Throughput





Saturn vs state-of-the-art

Saturn vs state-of-the-art Data freshness: average



Saturn vs state-of-the-art Data freshness: average



Evaluation: goals

- Can Saturn+Eunomia optimize both throughput and data freshness simultaneously?
- Why Saturn's genuine partial replication matters?

• How important is the configuration of **Saturn**?

• How **Eunomia** compares to sequencers?

Saturn Genuine partial replication matters



- Sydney receives updates from Ireland
- Ireland updates depend on N. Virginia updates
- Sydney **does not** replicate all N. Virginia data

Saturn Genuine partial replication matters

Latency of updates from Ireland when applied at Sydney



Saturn Genuine partial replication matters

Latency of updates from Ireland when applied at Sydney



Evaluation: goals

- Can Saturn+Eunomia optimize both throughput and data freshness simultaneously?
- Why **Saturn's** genuine partial replication matters?

How important is the configuration of Saturn?

• How **Eunomia** compares to sequencers?

Saturn configuration matters







Evaluation: goals

- Can Saturn+Eunomia optimize both throughput and data freshness simultaneously?
- Why **Saturn's** genuine partial replication matters?

• How important is the configuration of **Saturn**?

• How **Eunomia** compares to sequencers?

maximum throughput achievable by Eunomia vs a classical sequencer



maximum throughput achievable by Eunomia vs a classical sequencer





Other features...

- Support for **fault-tolerance**
- On-line **reconfiguration** protocols
- Efficient migration of clients among data centers
- Impact of stragglers
- More **experiments**



Check the papers for more details.

Take-away messages



- It is possible to preserve concurrency inside a DC, and to capture causality with a single scalar, using an ordering service that is out the the client's critical path
- These scalar clocks can be propagated using tree-based dissemination, that can effectively mitigate the impact of false dependencies
- Combined, they allow the optimisation of both throughput and data freshness