# UNOBTRUSIVE DEFERRED UPDATE STABILIZATION FOR EFFICIENT GEO-REPLICATION

Chathuri Gunawardhana, Manuel Bravo, Luís Rodrigues



In a geo-replicated system, how to apply **efficiently**, remote **updates**, in **causal** order?

# Why causal consistency?

#### Limitations of Highly-Available Eventually-Consistent Data Stores

Hagit Attiya Computer Science Department Technion

Faith Ellen Department of Computer Science University of Toronto

Adam Morrison e 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 specify *causal 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 ( $\mathcal{OCC}$ ).  $\mathcal{OCC}$  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, *causal 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

#### **Explicit Consistency**

[EuroSys'15]

#### **Session guarantees**

[SOSP'97]







more metadata

less metadata



Two main ways to compress and manage metadata

# **Global stabilization procedures**



# **Global stabilization procedures**

Updates are propagated concurrently and later ordered at remote datacenters

### **Serializers**

Updates are ordered before being applied at the origin datacenter

#### GentleRain and Cure use global stabilization!



**Serializer** (typically one per dc): SwiftCloud [Middleware'15], ChainReaction [EuroSys'13], ...



# Metadata compression



# Our goal



# aims at finding a **novel** way of compressing metadata that allows to pick a **better spot** in the throughput-visibility tradeoff

## conceived to **replace sequencers** in georeplicated storage systems

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

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





































**Eunomia** is implemented in C++ (200 LOC). At its core, it uses a **red-black tree** 

**EunomiaKV** is built as a **variant of** the open source version of **Riak** (100 lines of Erlang code)

Our implementation uses **hybrid logical clocks** (HLC) becoming resilient to both clock and workload skew Three datacenter deployment, emulating Amazon EC2 Virginia, Oregon, and Ireland regions

**baseline:** system that adds no overhead due to consistency management

comparison to: global stabilisation solutions, **GentleRain** (single scalar) and **Cure** (vector as EunomiaKV)

# Evaluation



# Evaluation



# **maximum throughput achievable** by Eunomia vs a classical **sequencer**



# Going back to the beginning





# Thanks!

### take-away message

by taking the coordination with an ordering service out the the client's critical path, one can pick a sweet-spot in the throughput vs. visibility tradeoff

check the paper!!!

fault-tolerant version of Eunomia

- impact of **stragglers**
- more **experiments**