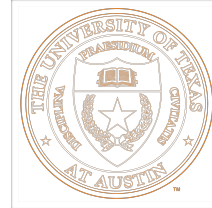# Reasoning with MAD distributed systems

Lorenzo Alvisi
The University of Texas at Austin

## Tolerating arbitrary faults
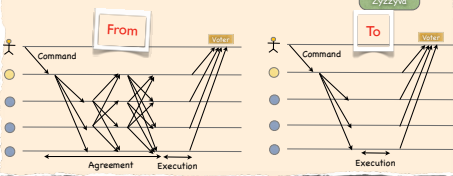
Bad things do happen to good systems …

**Why San Francisco's network admin went rogue**
An inside source reveals details of missteps and misunderstanding in the curious case of Terry Childs, network kidnapper

**Amazon S3 Issues: Load Balancers and MD5**

**Amazon S3 Availability Event: July 20, 2008**
We wanted to provide some additional detail about the problem we experienced on files using MD5 to perform integrity checks. After some investigation, Amazon confirmed the problem and identified the cause:

Amazon's S3 storage system had some issues last week with data corruption on files using MD5 to perform integrity checks. After some investigation, Amazon S3's total requests in the US. Intermittently, under load, it was corrupting single bytes in the byte stream. … Based on our investigation with both internal

**Gmail Disaster: Reports Of Mass Email Deletions**

### Zyzzyva
**Speculative Byzantine Fault Tolerance**
- **Simplifies** the design of BFT replication
  - One protocol to rule them all
  - ✓ latency  ✓ throughput  ✓ cost of replication

BFT?  →  Yes  →  Zyzzyva

From / To

Command — Agreement — Execution
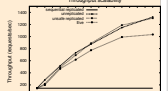Command — Execution

### UpRight
**Practical and configurable BFT replication**
- Applying BFT replication to **real** services
  - BFT HDFS
  - BFT ZooKeeper
- Refining the architecture and APIs
  - Introduce Request Quorum stage before agreement
  - Clean application API for:
    - processing requests
    - taking application state snapshots
- **Configurable** replication
  - u: services are **Up** (live) despite u failures
  - r: services are **Right** (safe) despite r commission failures
  - Replication costs expressed as a function of **u** and **r**

Byzantine (FT) Empire — before UpRight / after UpRight
- Synchronous, no failures (high performance)
- Synchronous (minimal liveness guarantees)
- Synchronous, with faults (minimal liveness guarantees)
- Asynchronous with or without failures (high performance)
- Asynchronous (minimal liveness guarantees)

### EVE
**Replicating multithreaded servers**
- Replication state-of-the-art
  - **Agree** on order of requests, then **execute** them
  - Requires **deterministic** execution
    - Practically this means **single-threaded** execution
- Eve
  - First **execute** requests nondeterministically, without agreeing on order
  - Then **verify** if state and responses match among replicas
  - Efficient rollback on divergence
- Benefits
  - Allows multithreaded execution
  - Up to 12x speedup on a 16-core machine
  - 25% slower than an unreplicated server

## MAD Systems

### What is a MAD system?
Any system that spans **M**ultiple **A**dministrative **D**omains (e.g. peer-to-peer services, cloud/outsourced storage, Internet routing, and wireless mesh routing.

### What is so special about MAD systems?
Traditional threshold FT does **not** apply!
- Nodes can be selfish: cooperation requires incentives
- Sybil attacks can overwhelm any threshold mechanism

If this were not enough, each domain is a black box to its peers: what basis is there for trust?

### The BAR Model
Three classes of MAD nodes:
- **B**yzantine: deviate arbitrarily, for any reason
- **A**cquiescent: follow the assigned protocol obediently
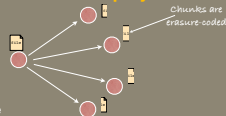- **R**ational: deviate iff doing so increases their utility

#### BAR Tolerant Systems
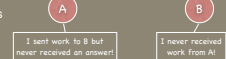- No more than $n/3$ Byzantine nodes
- No bound on number of rational nodes
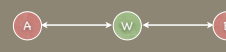
### BAR-B
**A BAR-tolerant cooperative backup system**
- Peers assign each other *chunks* to store on their behalf
- Assignment need not be symmetric
- Deterministic retrieval guarantee despite Byzantine and Rational peers

Chunks are erasure-coded

When assigning work, challenge is handling ``he said / she said''

I sent work to B but never received an answer! / I never received work from A!

Problem could be solved by interposing an acquiescent witness W between A and B

But, just in FT distributed computing it is not prudent to *assume* that any particular node will be correct, we don't want to assume that any W that we may use will not either fail or turn selfish.

**Solution:** use BAR-tolerant State Machine Replication to build the abstraction of an acquiescent W out of node each of which may be Byzantine or selfish
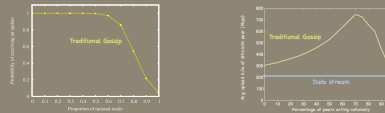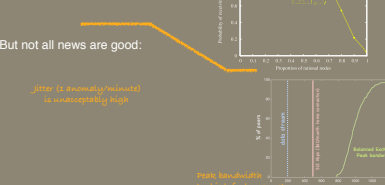
### BAR-Gossip
**BAR-tolerant Nash for P2P live streaming**
- Gossip is attractive infrastructure for P2P live streaming
- But gossip protocols perform poorly if many peers behave selfishly

Traditional Gossip / Data stream

- BAR Gossip relies on Balance Exchange, a provably incentive compatible protocol: no selfish node has unilateral incentives to deviate from it

Balanced exchange

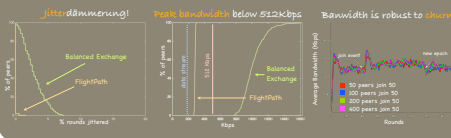- Reliability with BAR Gossip is way up…
- But not all news are good:

jitter (1 anomaly/minute) is unacceptably high

Peak bandwidth too high for home use!
Balanced Exchange Peak bandwidth

### Flightpath
**Approximate equilibria for practical live streaming**
- The price for proving BAR Gossip a Nash equilibrium is lack of flexibility:
  - peers cannot join streaming mid-way
  - communication patterns are inflexible
  - extra overhead
- Flightpath balances obedience with choice through approximate equilibria
  - not Nash, but $\varepsilon$-Nash: selfish node deviate only if doing so increases their utility by more than a factor of $\varepsilon$
- Flightpath supports dynamic membership; provides stable performance despite flash crowds; minimizes jitter; and lowers peak bandwidth below home-use threshold

jitter dämmerung! / Peak bandwidth below 512Kbps / Bandwidth is robust to churn

Balanced Exchange / FlightPath

### Just in! Local social defenses against Sybil attacks
- Current sybil defenses can distinguish honest from forged identities in social graphs that are fast mixing (equivalently, have constant conductance)
- Alas, many social graphs are not fast mixing!
- We are developing a new approach that provides better protection without relying on global graph properties, such as constant conductance, but rather leverages the social graph's community structure.
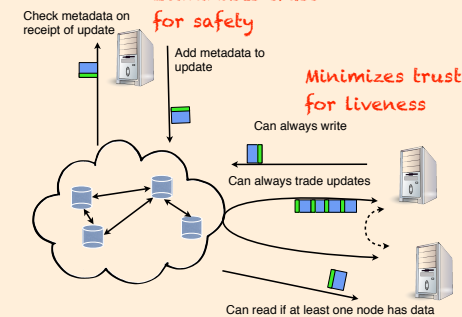
## Dependable Storage

### Depot
**Cloud storage with minimal trust**
- Removes trust from providers
- Not the same thing as making providers more trustworthy!

*Eliminates trust for safety*
*Minimizes trust for liveness*

Check metadata on receipt of update
Add metadata to update
Can always write
Can always trade updates
Can read if at least one node has data

### Teapot
**Minimal trust for today's cloud**
- Same guarantees of Depot, but using unmodified Amazon S3 servers

*In progress!*

Lorenzo Alvisi is a Professor in the Department of Computer Science at UT Austin, where he is a co-director of the Laboratory for Advanced Systems Software (LASR). He holds a Ph.D. and M.S. in Computer Science from Cornell University, and a Laurea summa cum laude in Physics from the University of Bologna, Italy. He is a Fellow of the ACM and the recipient of an Alfred P. Sloan Fellowship and of the NSF CAREER Award.