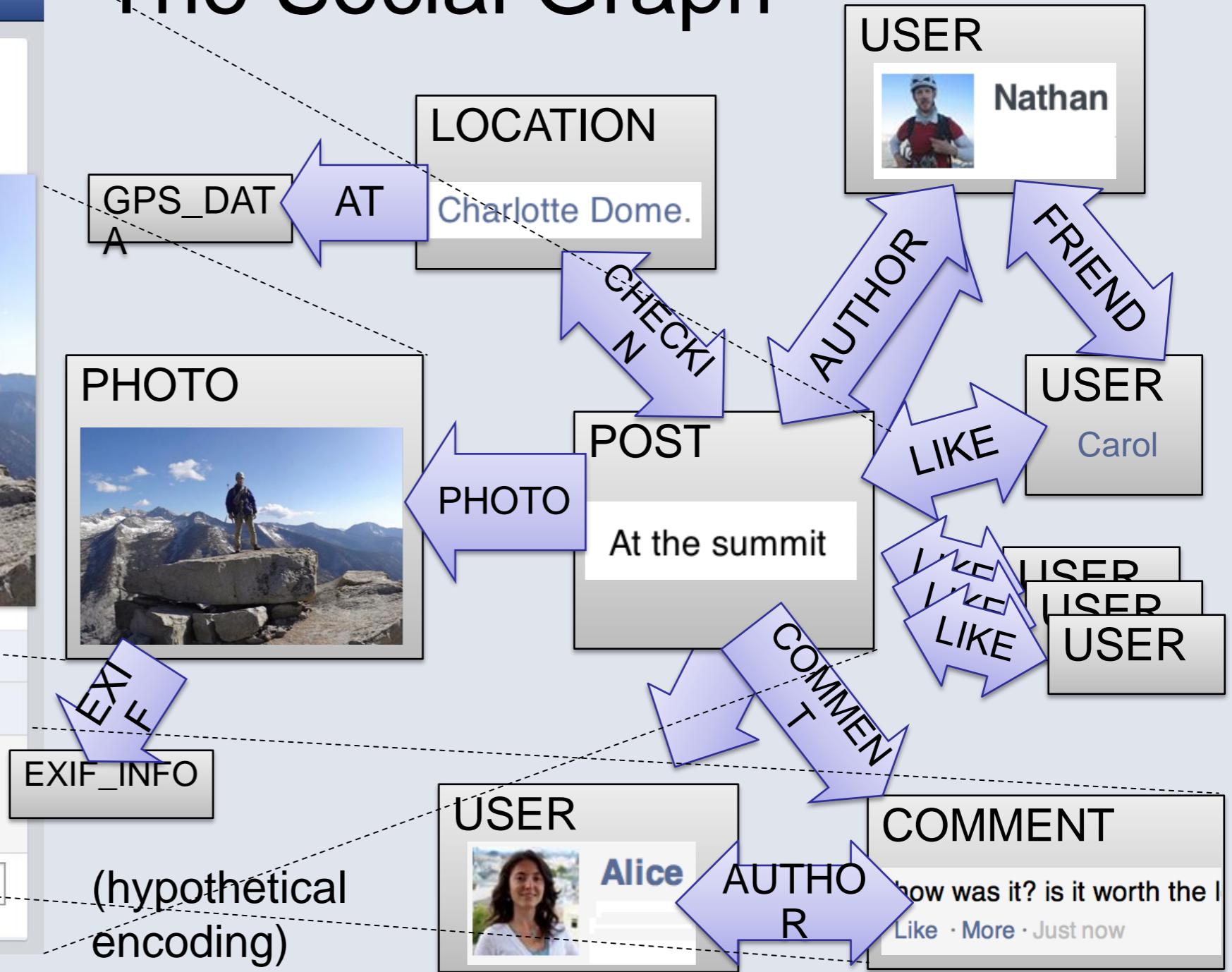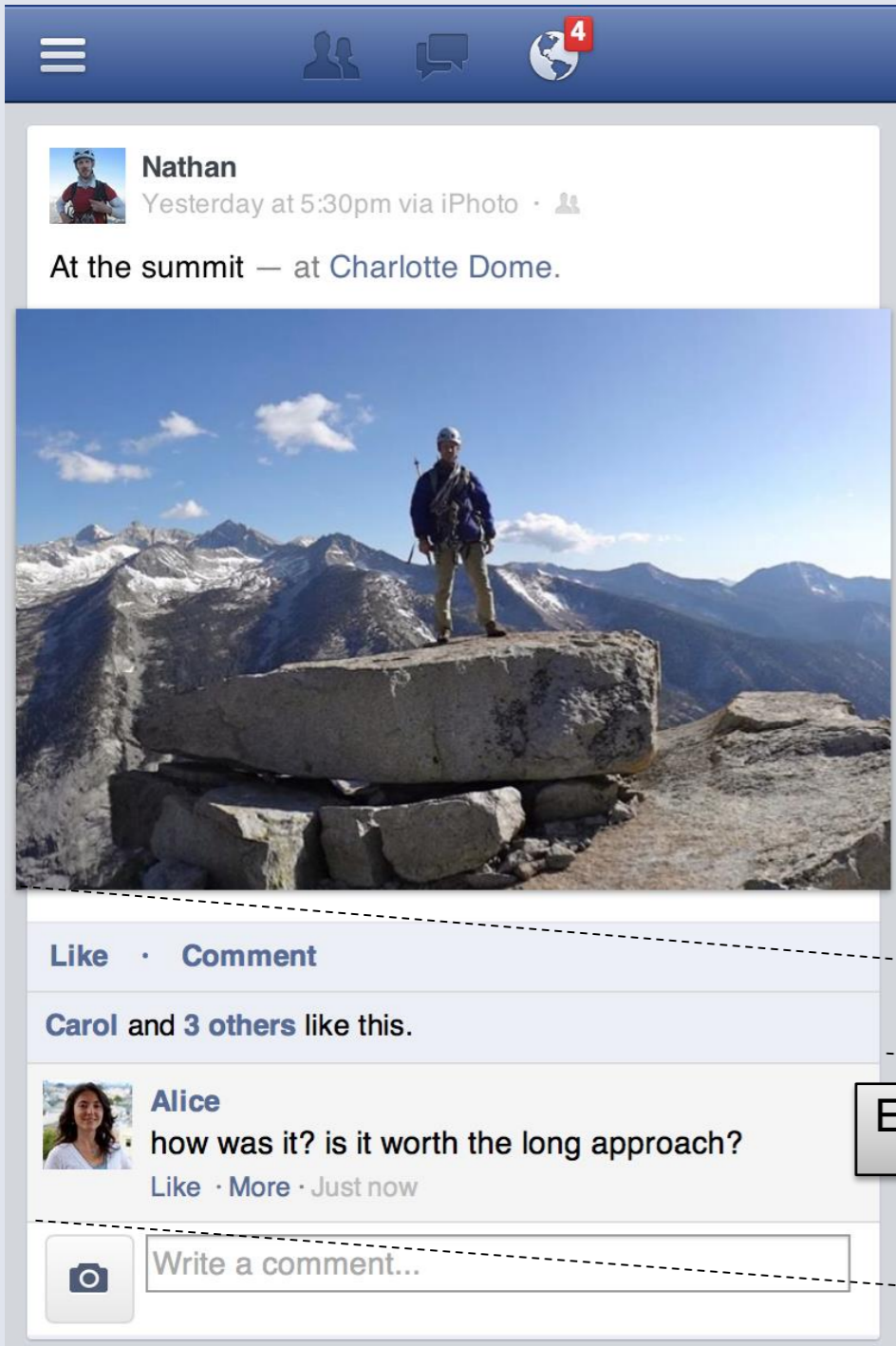**facebook**

# TAO
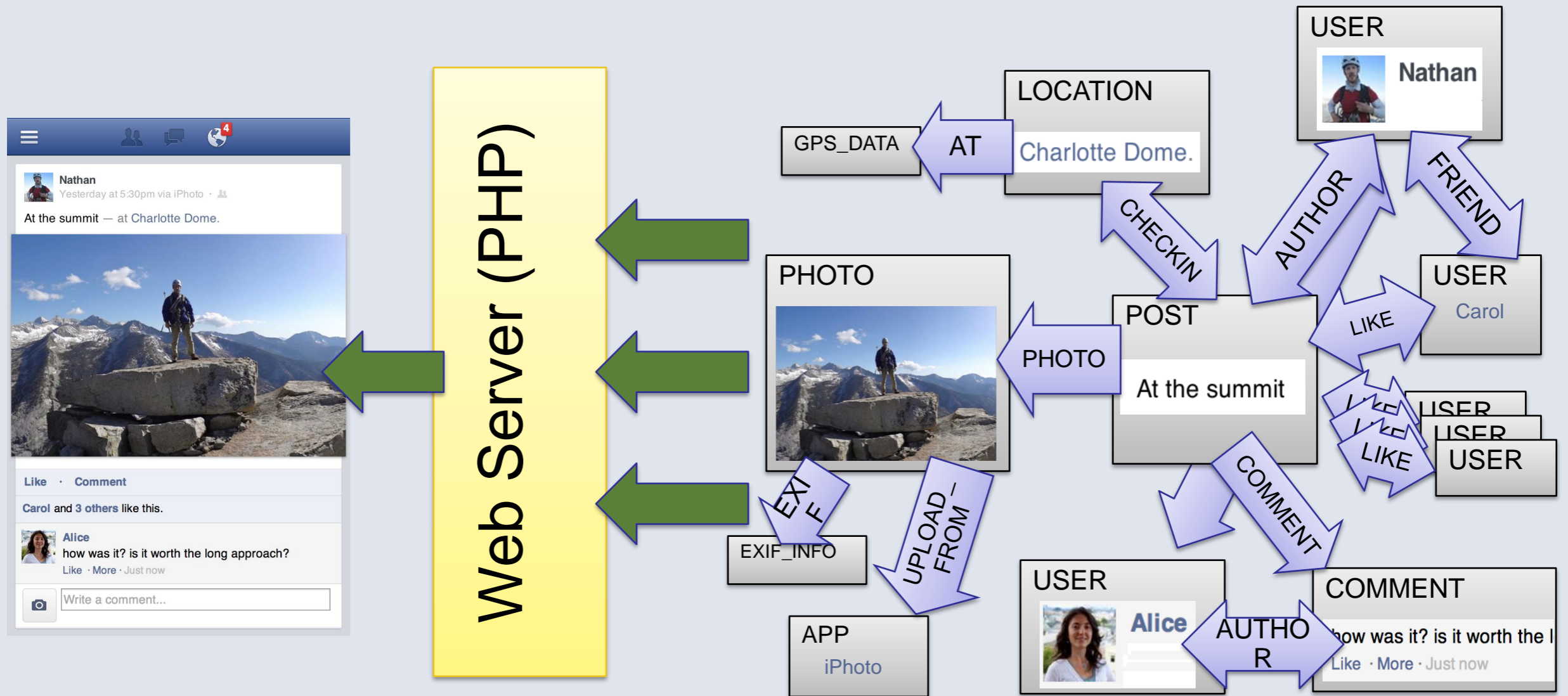Facebook's Distributed Data Store for the Social Graph

**Nathan Bronson**, Zach Amsden, George Cabrera, Prasad Chakka, Peter Dimov, Hui Ding, Jack Ferris, Anthony Giardullo, Sachin Kulkarni, Harry Li, Mark Marchukov, Dmitri Petrov, Lovro Puzar, Yee Jiun Song, Venkat Venkataramani
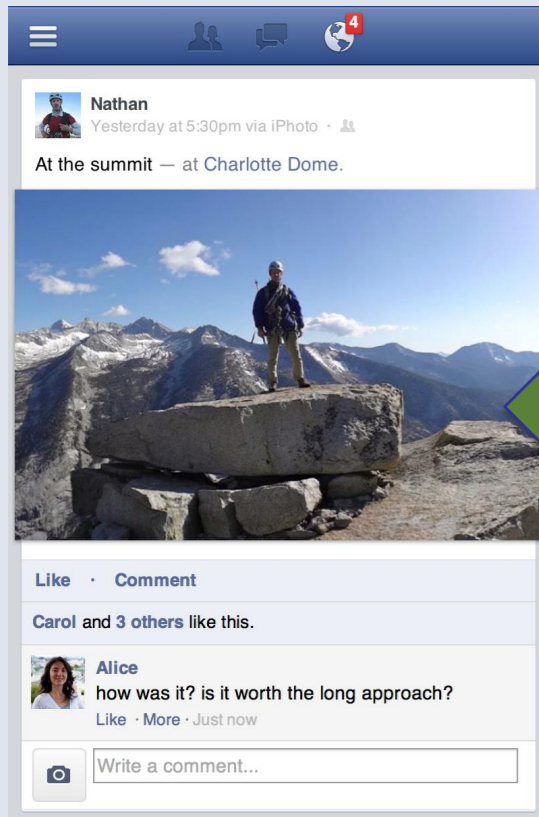
Presented at USENIX ATC – June 26, 2013

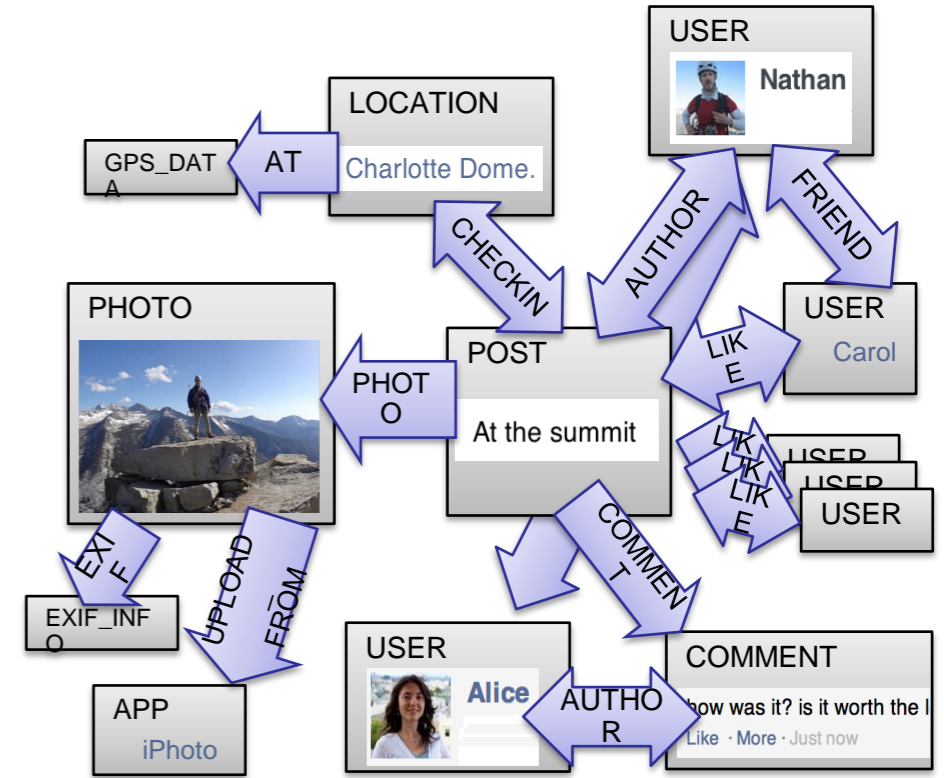# The Social Graph

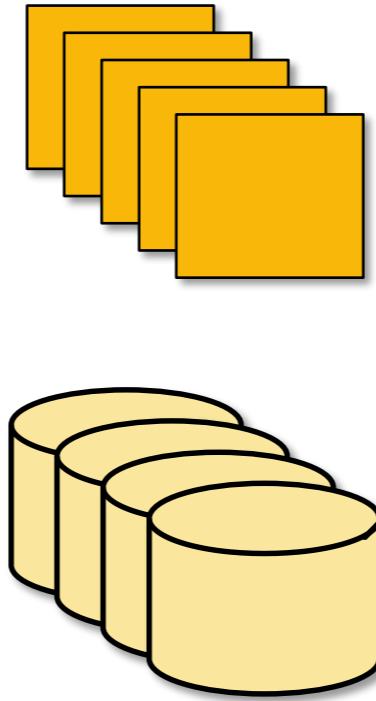# Dynamically Rendering the Graph

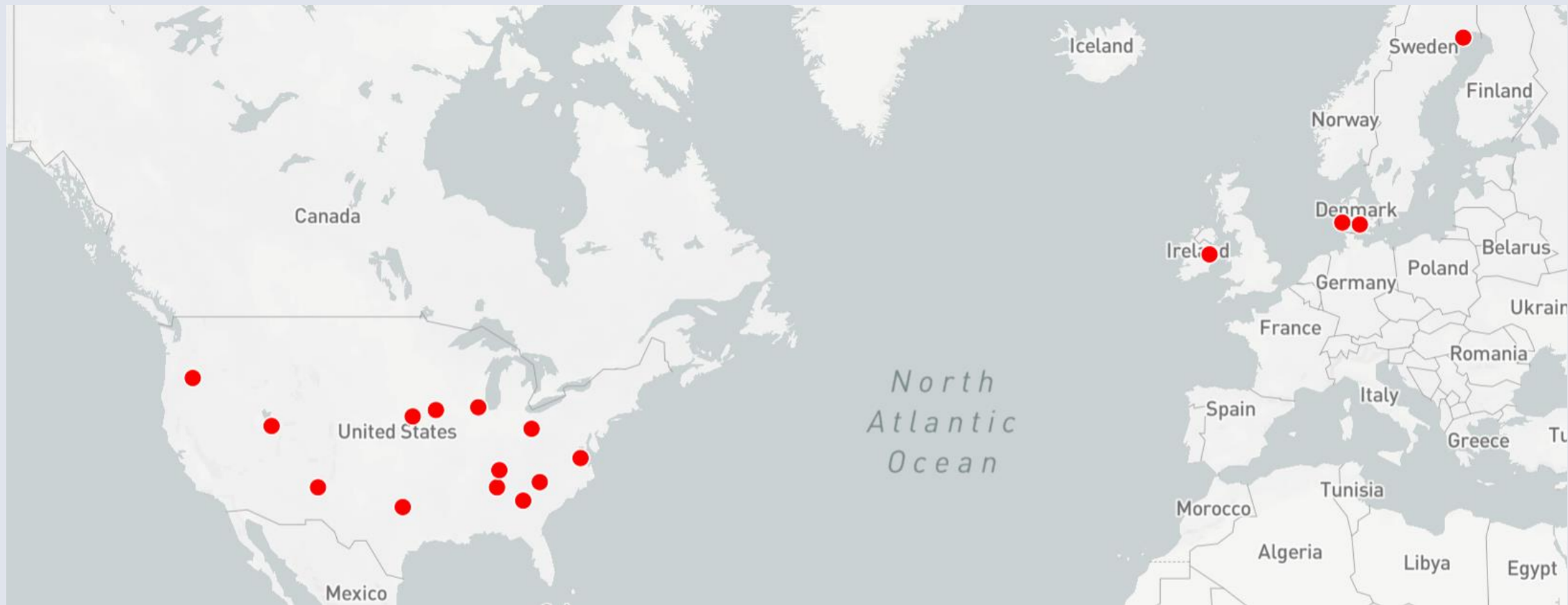# Dynamically Rendering the Graph



Web Server (PHP)

## TAO

- 1 billion queries/second
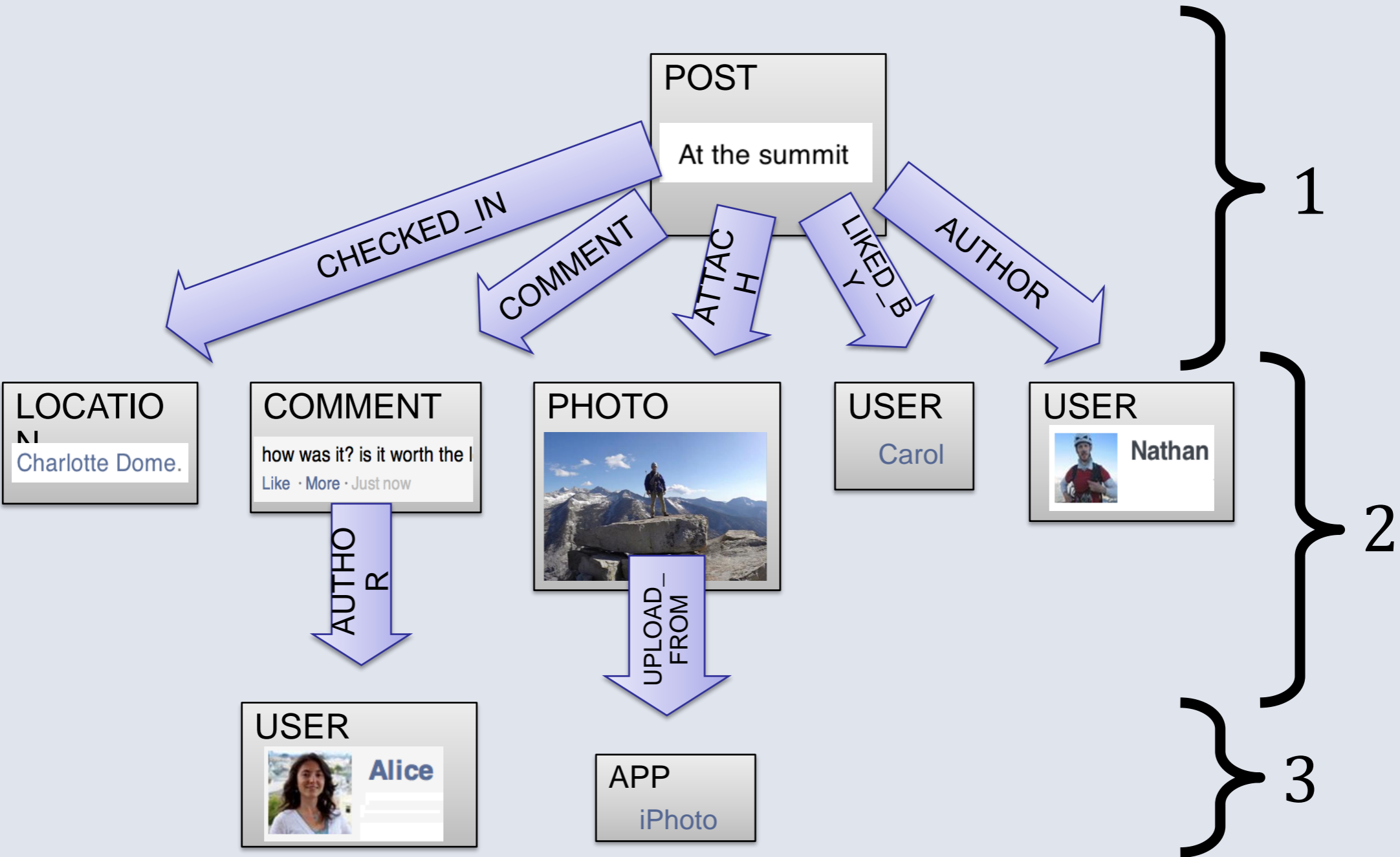- many petabytes of data

# facebook data centers

# What Are TAO's Goals/Challenges?

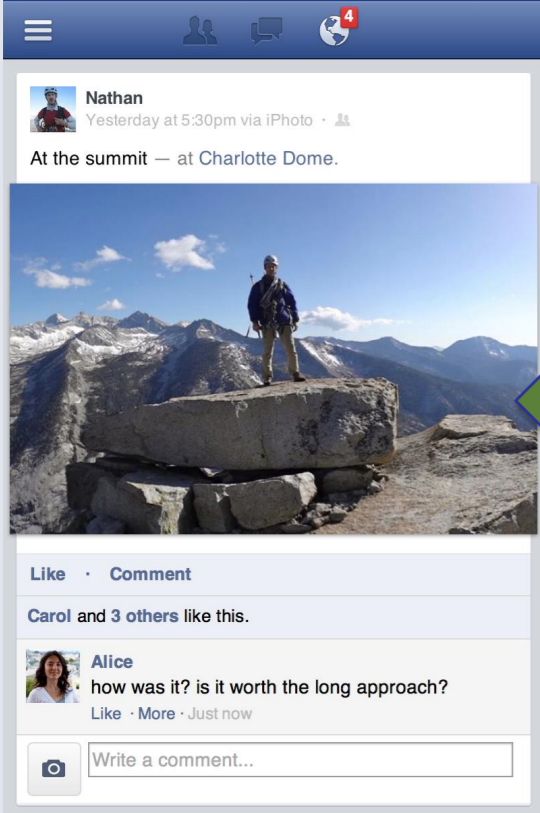- Efficiency at scale

December 2010

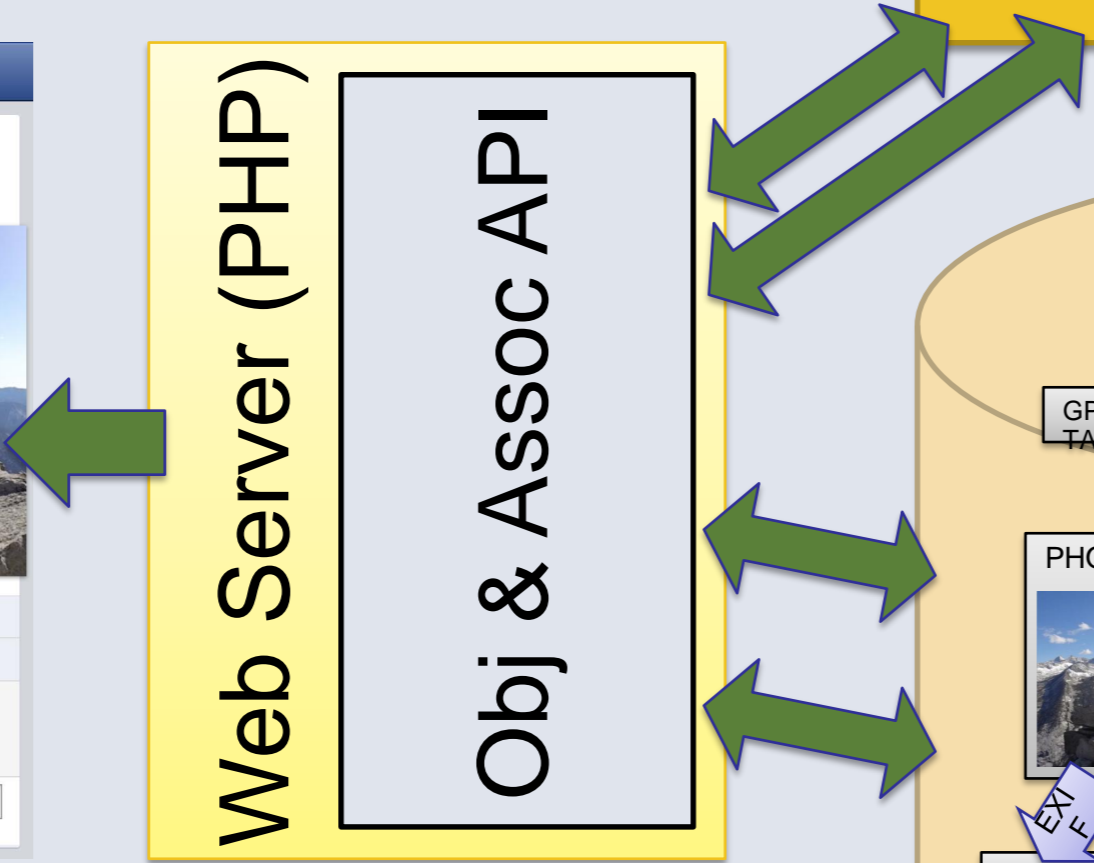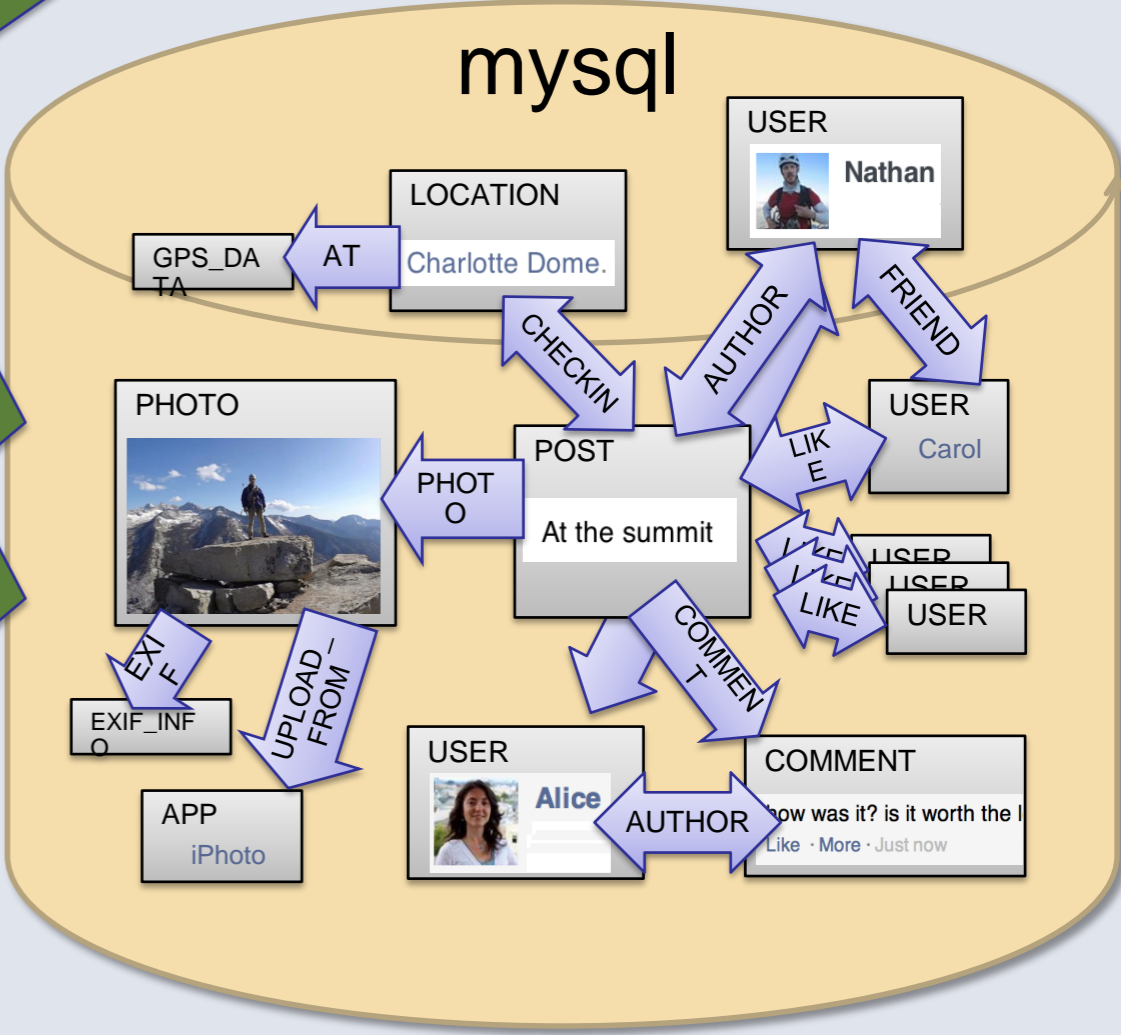# Dynamic Resolution of Data Dependencies

# What Are TAO's Goals/Challenges?

- Efficiency at scale
- Low read latency
- Timeliness of writes
- High Read Availability

facebook

December 2010

# Graph in Memcache

# Objects = Nodes

- Identified by unique 64-bit IDs

- Typed, with a schema for fields

# Associations = Edges

- Identified by <id1, type, id2>

- Bidirectional associations are two edges, same or different type

id: 1807 =>
  type: POST
  str: "At the summ…

<1807,COMMENT,200
3>
time: 1,371,704,655

id: 2003 =>
  type: COMMENT
  str: "how was it …

<308,AUTHORED,2003>
time: 1,371,707,355

<2003,AUTHOR,308>
time: 1,371,707,355

id: 308 =>
  type: USER
  name: "Alice"

# Association Lists

- <id1, type, *>

- Descending order by time

- Query sublist by position or time

- Query size of entire list

id: 1807 =>
  type: POST
  str: "At the summ…

<1807,COMMENT,4141>
time: 1,371,709,009

<1807,COMMENT,8332>
time: 1,371,708,678

<1807,COMMENT,2003>
time: 1,371,707,355

id: 4141 => type: COMMENT
  str: "Been wanting to do …

id: 8332 => type: COMMENT
  str: "The rock is flawless, …

id: 2003 => type: COMMENT
  str: "how was it, was it w…

newer

older

# Objects and Associations API

## Reads – 99.8%

- Point queries
  - **obj_get** _____ 28.9%
  - **assoc_get** _____ 15.7%
- Range queries
  - **assoc_range** _____ 40.9%
  - **assoc_time_range** ____ 2.8%
- Count queries
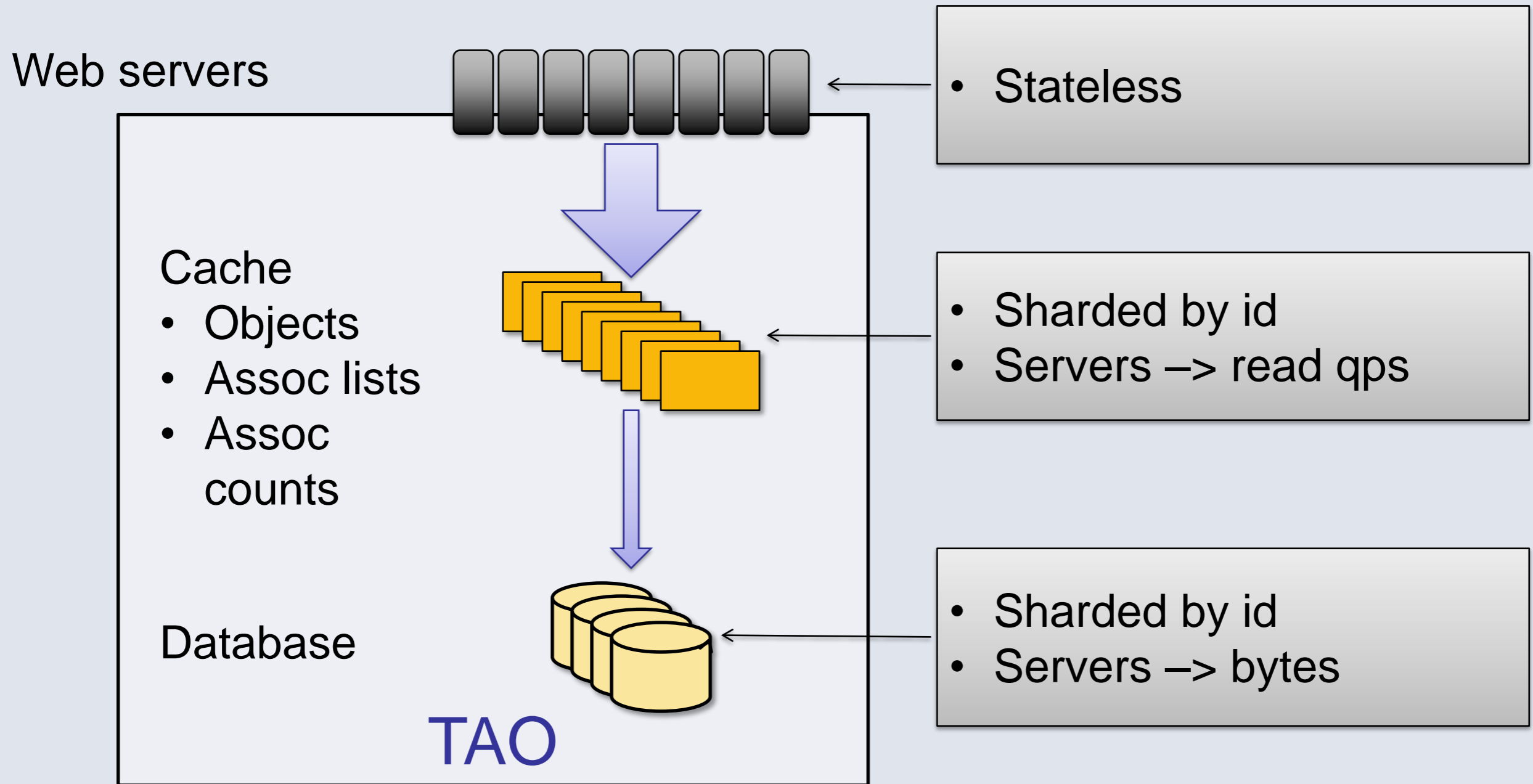  - **assoc_count** _____ 11.7%

## Writes – 0.2%

- Create, update, delete for objects
  - **obj_add** _____ 16.5%
  - **obj_update** _____ 20.7%
  - **obj_del** _____ 2.0%
- Set and delete for associations
  - **assoc_add** _____ 52.5%
  - **assoc_del** _____ 8.3%

# Independent Scaling by Separating Roles

Web servers

- Stateless

Cache
- Objects
- Assoc lists
- Assoc counts

- Sharded by id
- Servers –> read qps

Database

TAO

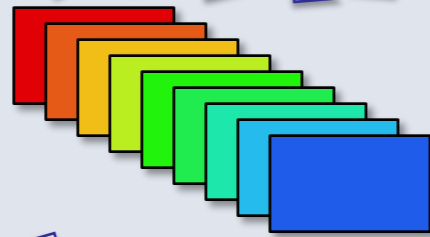- Sharded by id
- Servers –> bytes
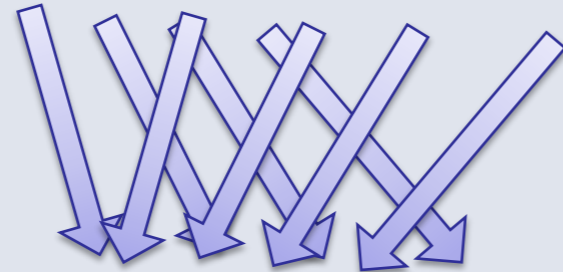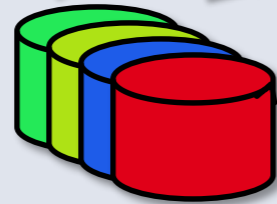
# Subdividing the Data Center

Web servers



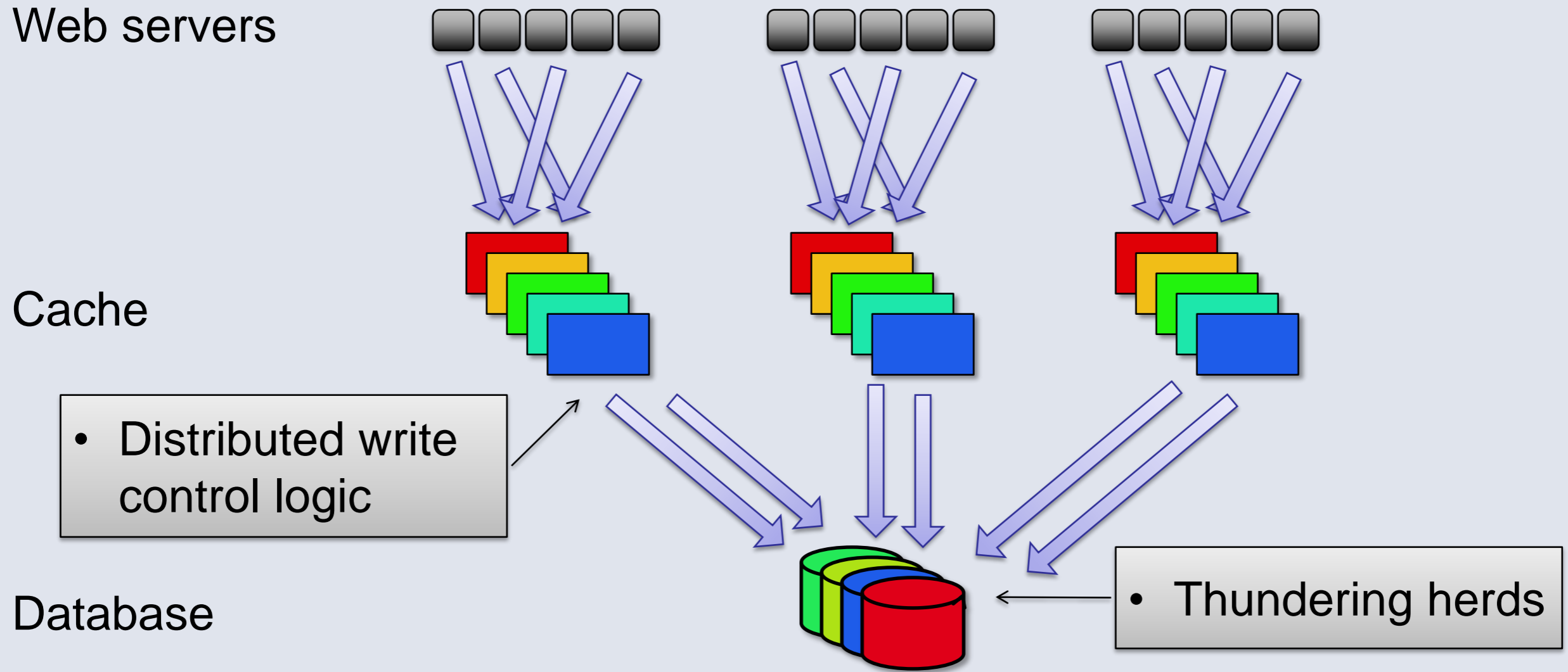- Inefficient failure detection
- Many switch traversals

Cache

- Many open sockets
- Lots of hot spots

Database

# Subdividing the Data Center

Web servers

Cache
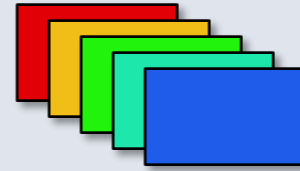
- Distributed write control logic
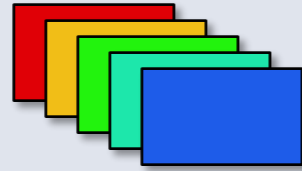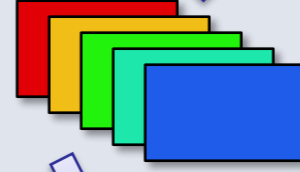
Database

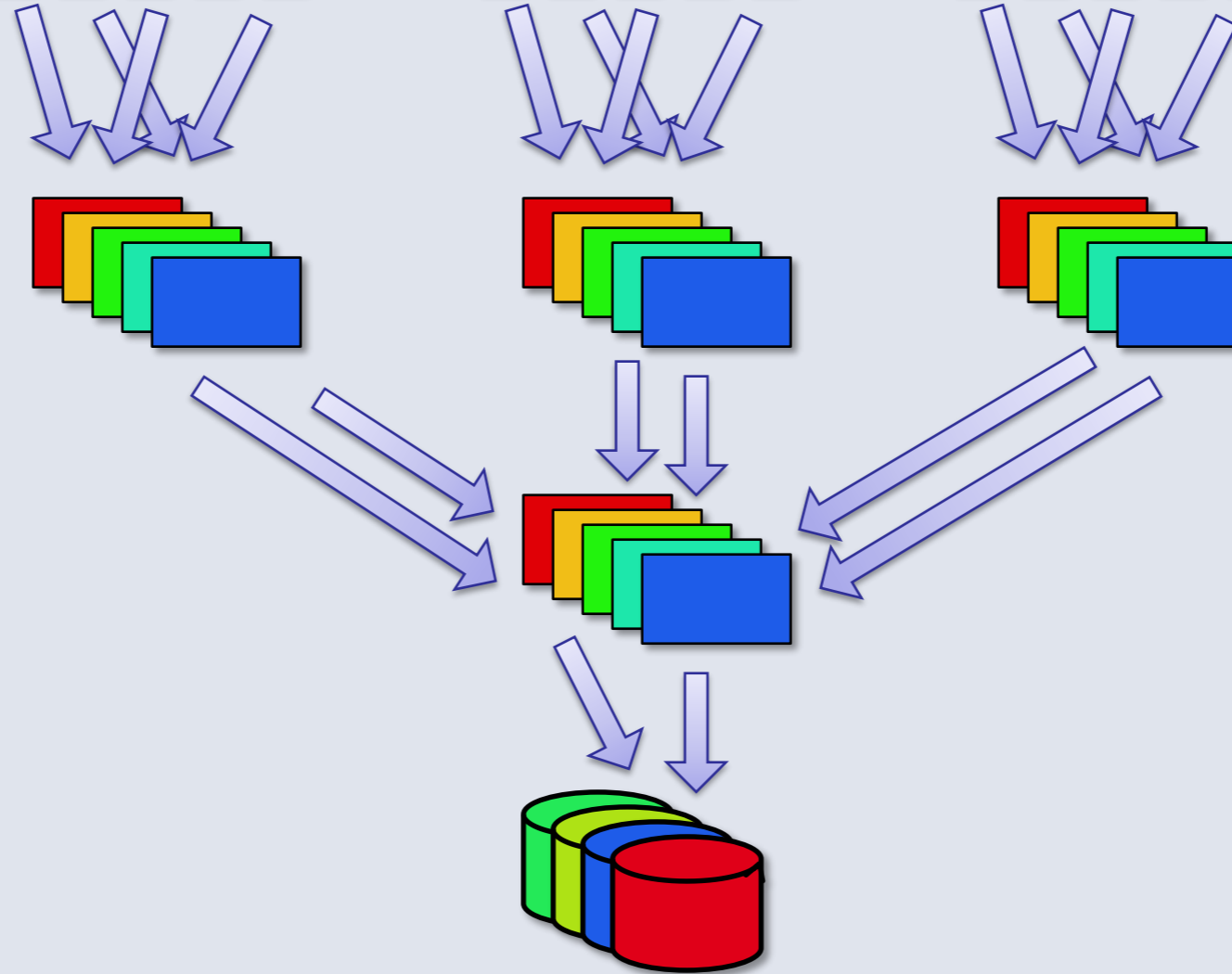- Thundering herds

# Follower and Leader Caches
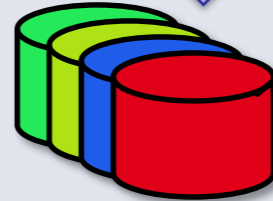
Web servers

Follower cache

Leader cache

Database

# What Are TAO's Goals/Challenges?

- Efficiency at scale

- Low read latency

- Timeliness of writes

- High Read Availability

facebook

December 2010

# Write-through Caching – Association Lists

Web servers

Follower cache

Y,A,B,C          Y,A,B,C

X –> Y    ok

refill X

range get                                           X –> Y    ok                    refill X

Leader cache

Y,A,B,C

X –> Y    ok

Database

Y,...

# Asynchronous DB Replication

Master data center | Replica data center

Web servers

Writes forwarded to master
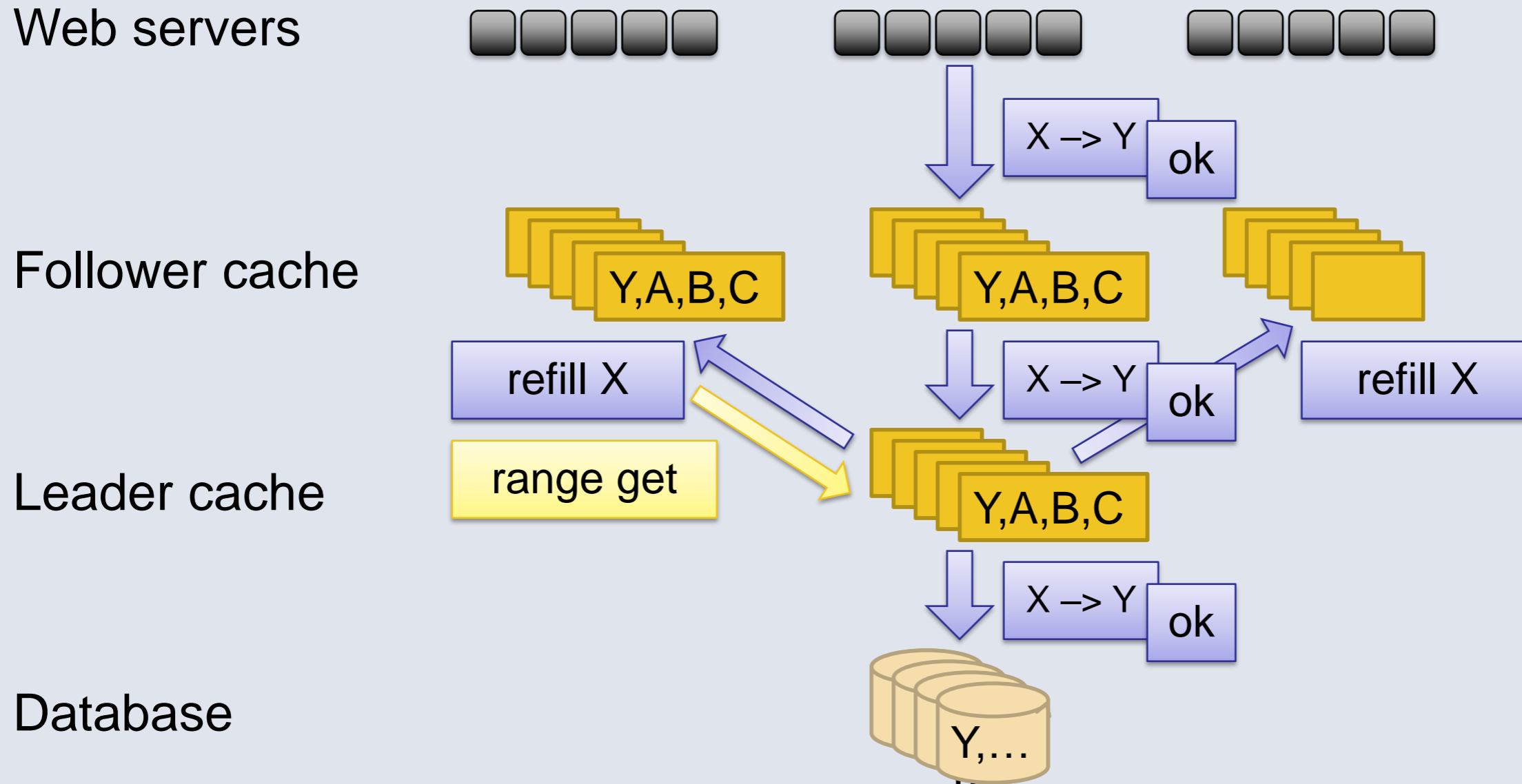
Inval and refill embedded in SQL

Delivery after DB replication done

# What Are TAO's Goals/Challenges?

- Efficiency at scale

- Low read latency

- Timeliness of writes

- High Read Availability

facebook

December 2010

# Improving Availability: Read Failover

Master data center
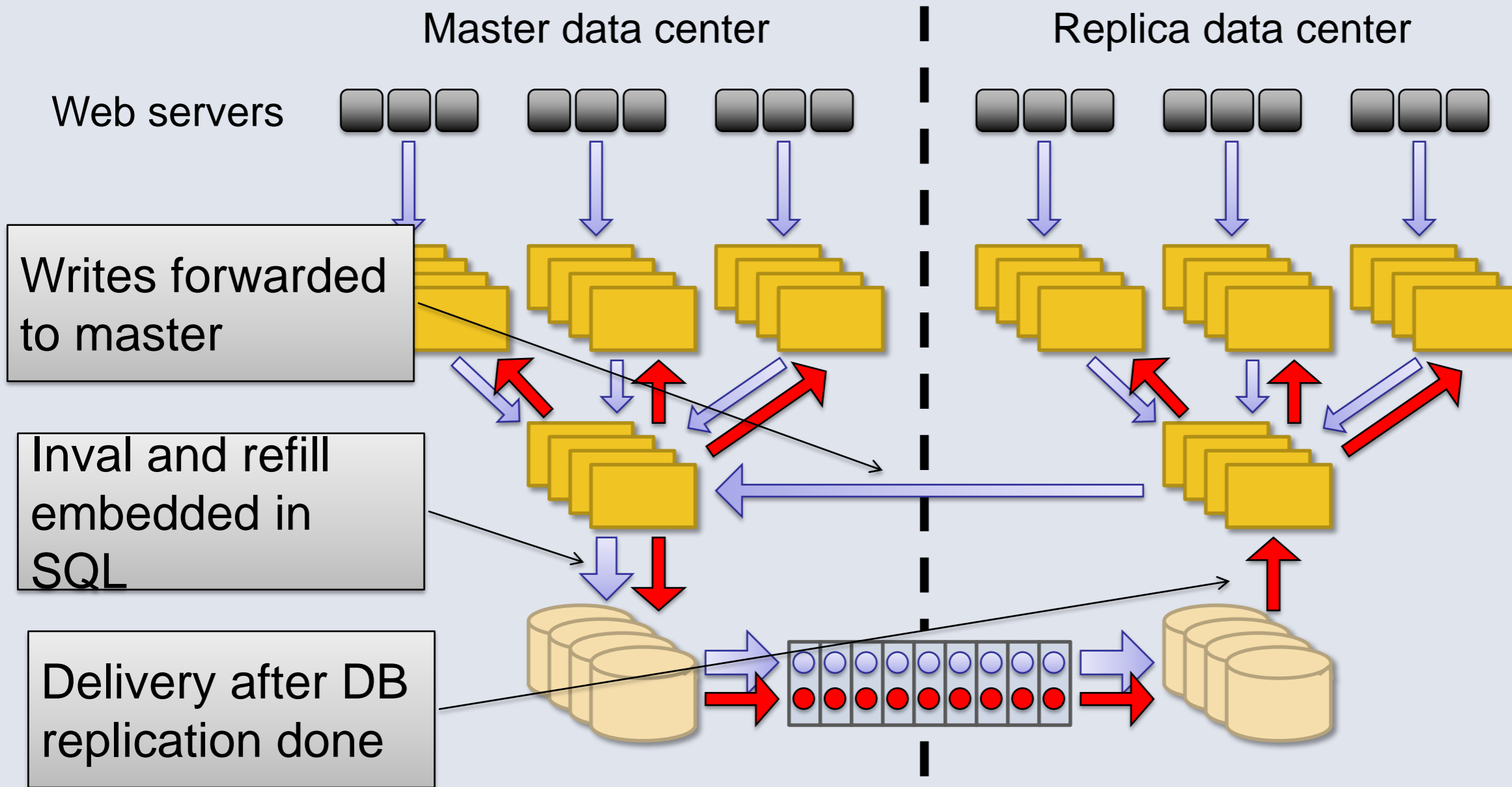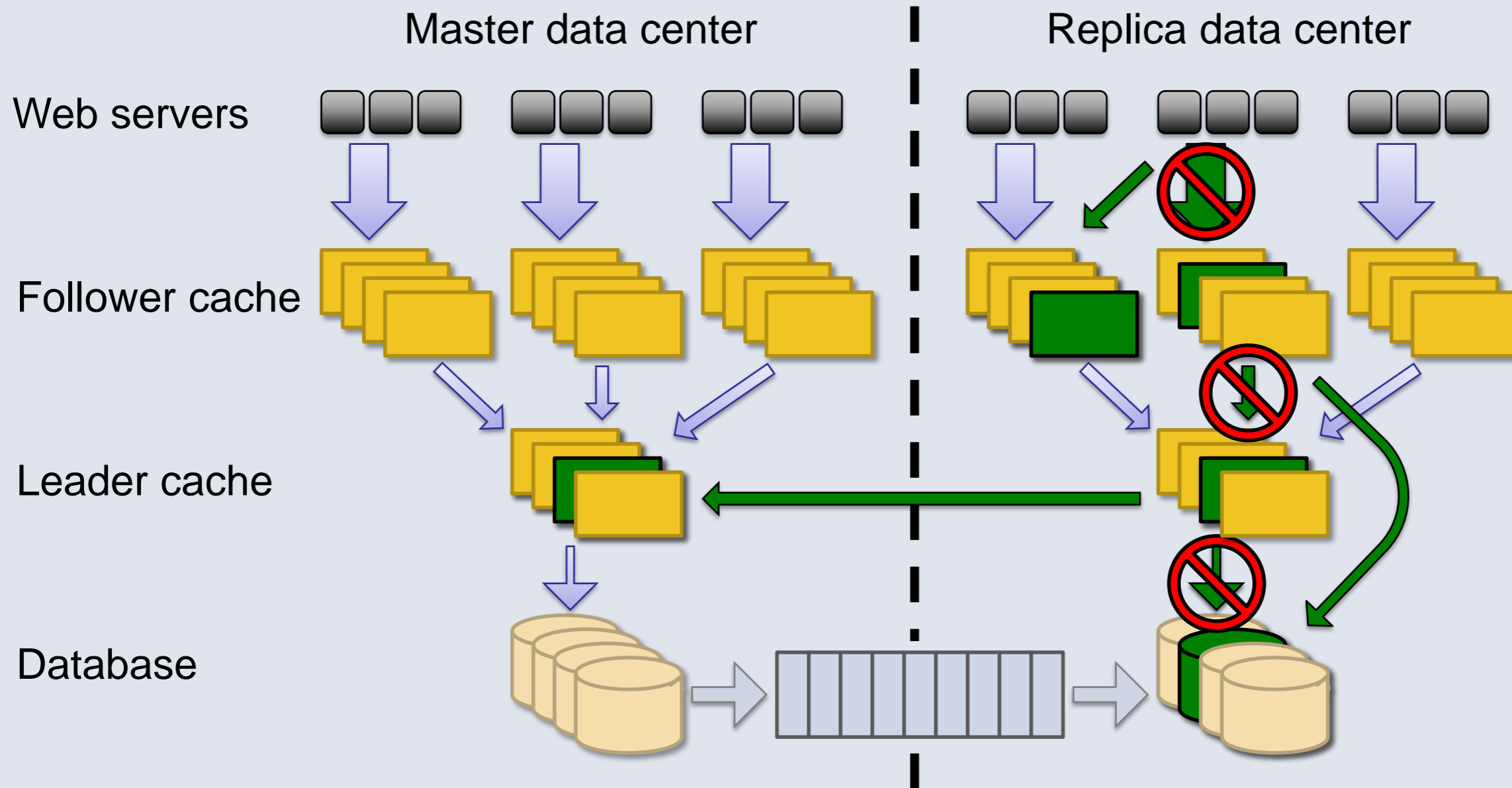
Replica data center

Web servers

Follower cache

Leader cache

Database

# TAO Summary

**Efficiency at scale**
**Read latency**

- Separate cache and DB
- Graph-specific caching
- Subdivide data centers

**Write timeliness**

- Write-through cache
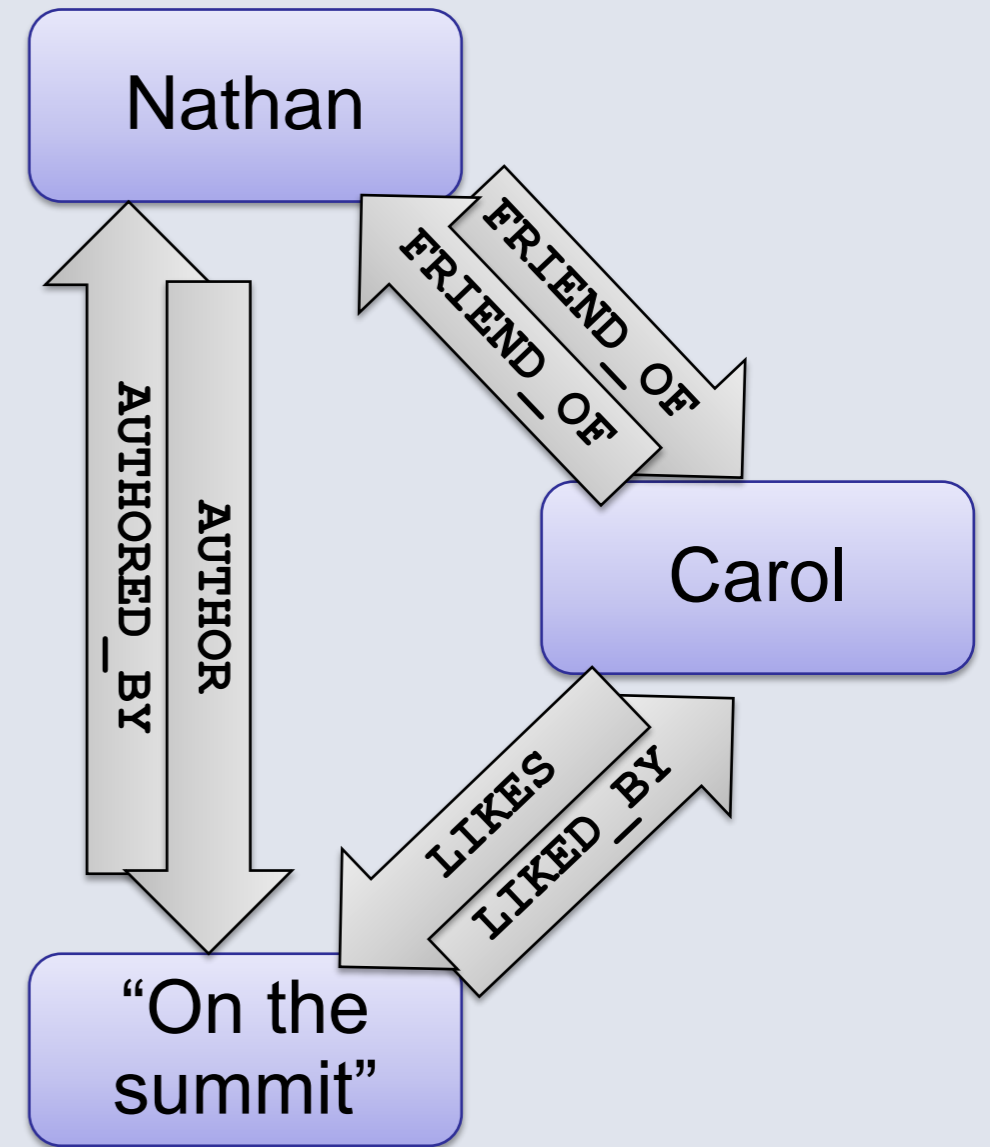- Asynchronous replication
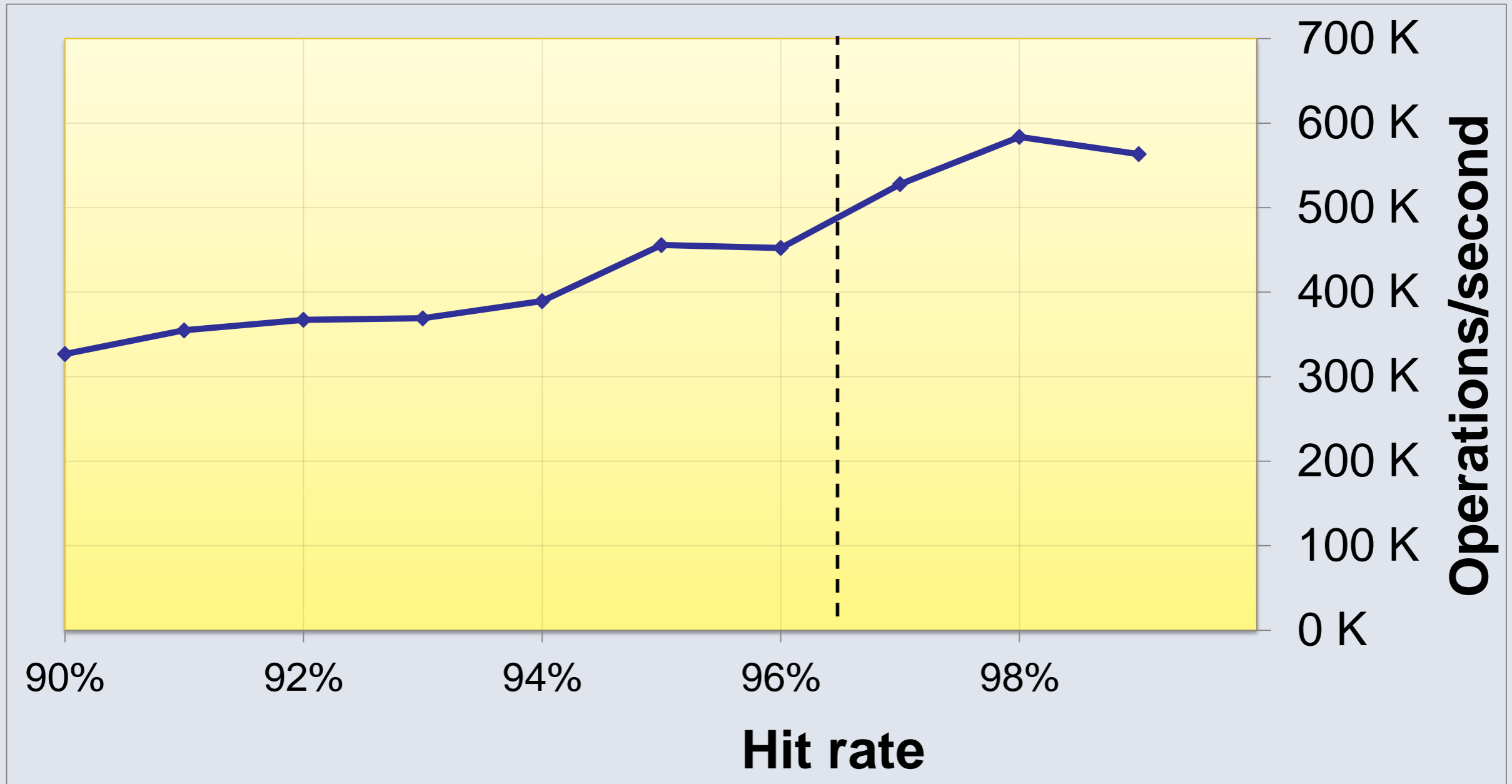
**Read availability**

- Alternate data sources

# facebook

# Inverse associations

- Bidirectional relationships have separate *a→b* and *b→a* edges

  - inv_type(`LIKES`) = `LIKED_BY`

  - inv_type(`FRIEND_OF`) = `FRIEND_OF`

- Forward and inverse types linked only during write

  - TAO `assoc_add` will update both

  - Not atomic, but failures are logged and repaired

Single-server Peak Observed Capacity

# Write latency

# More In the Paper

- The role of association time in optimizing cache hit rates

- Optimized graph-specific data structures

- Write failover

- Failure recovery

- Workload characterization