

facebook

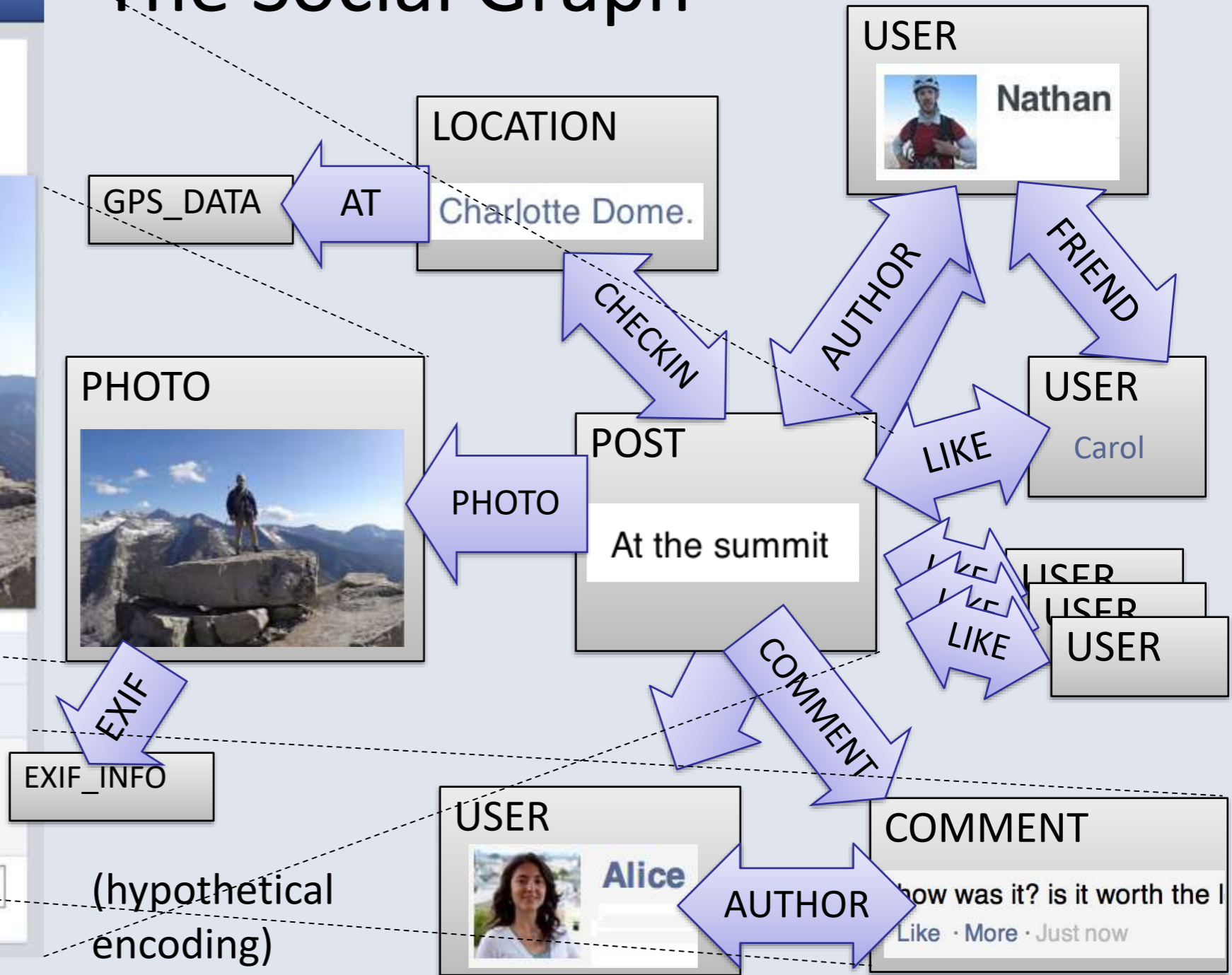
TAO

Facebook's Distributed Data Store for the Social Graph

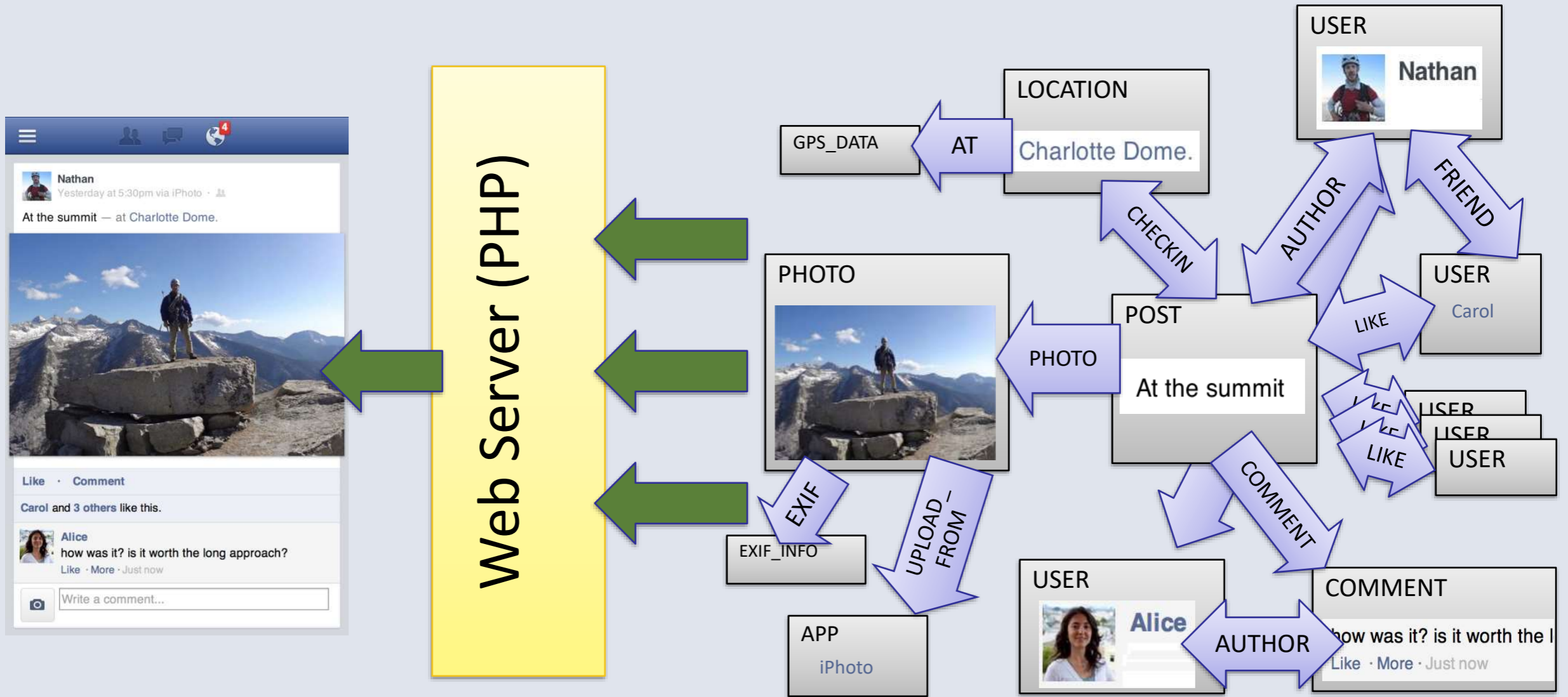
Nathan Bronson, Zach Amsden, George Cabrera, Prasad Chakka, Peter Dimov, Hui Ding, Jack Ferris, Anthony Giardullo, Sachin Kulkarni, Harry Li, Mark Marchukov, Dmitri Petrov, Lovro Puzar, Yee Jiun Song, Venkat Venkataramani

Presented at USENIX ATC – June 26, 2013

The Social Graph

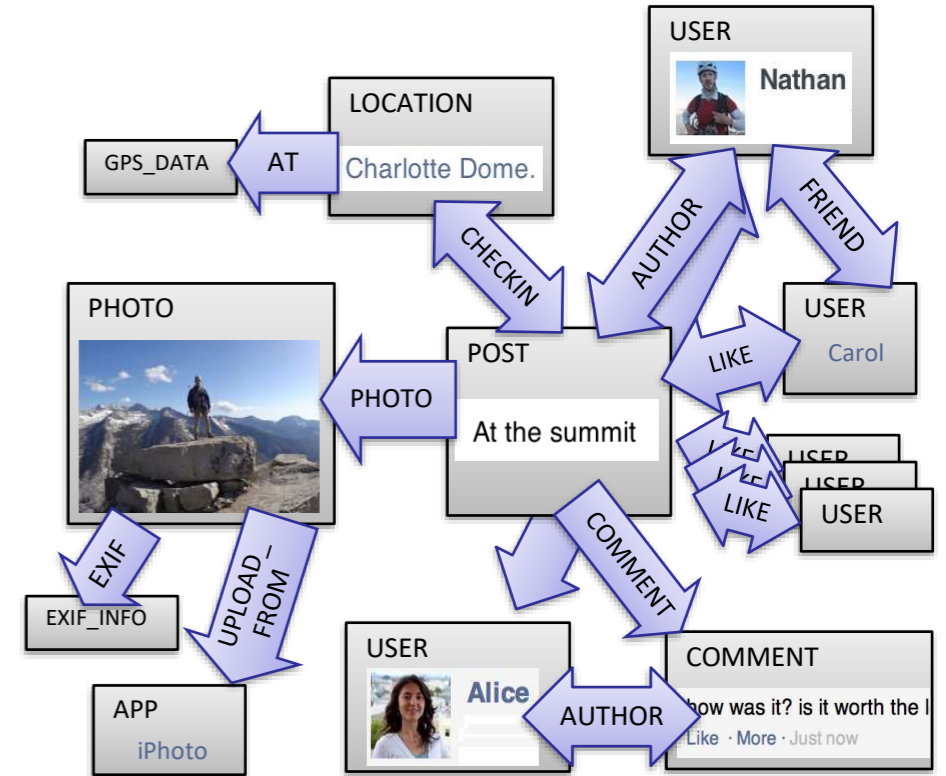
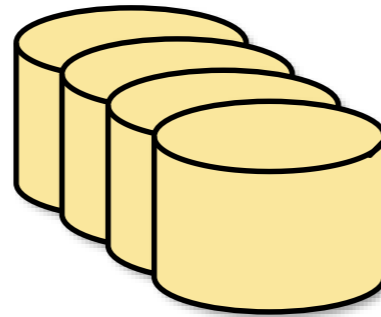
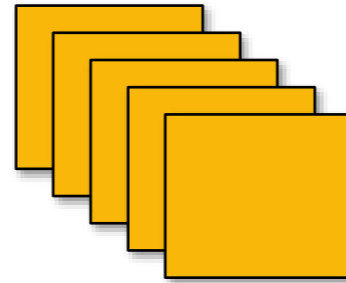


Dynamically Rendering the Graph



Dynamically Rendering the Graph

TAO



Web Server (PHP)

- 1 billion queries/second
- many petabytes of data



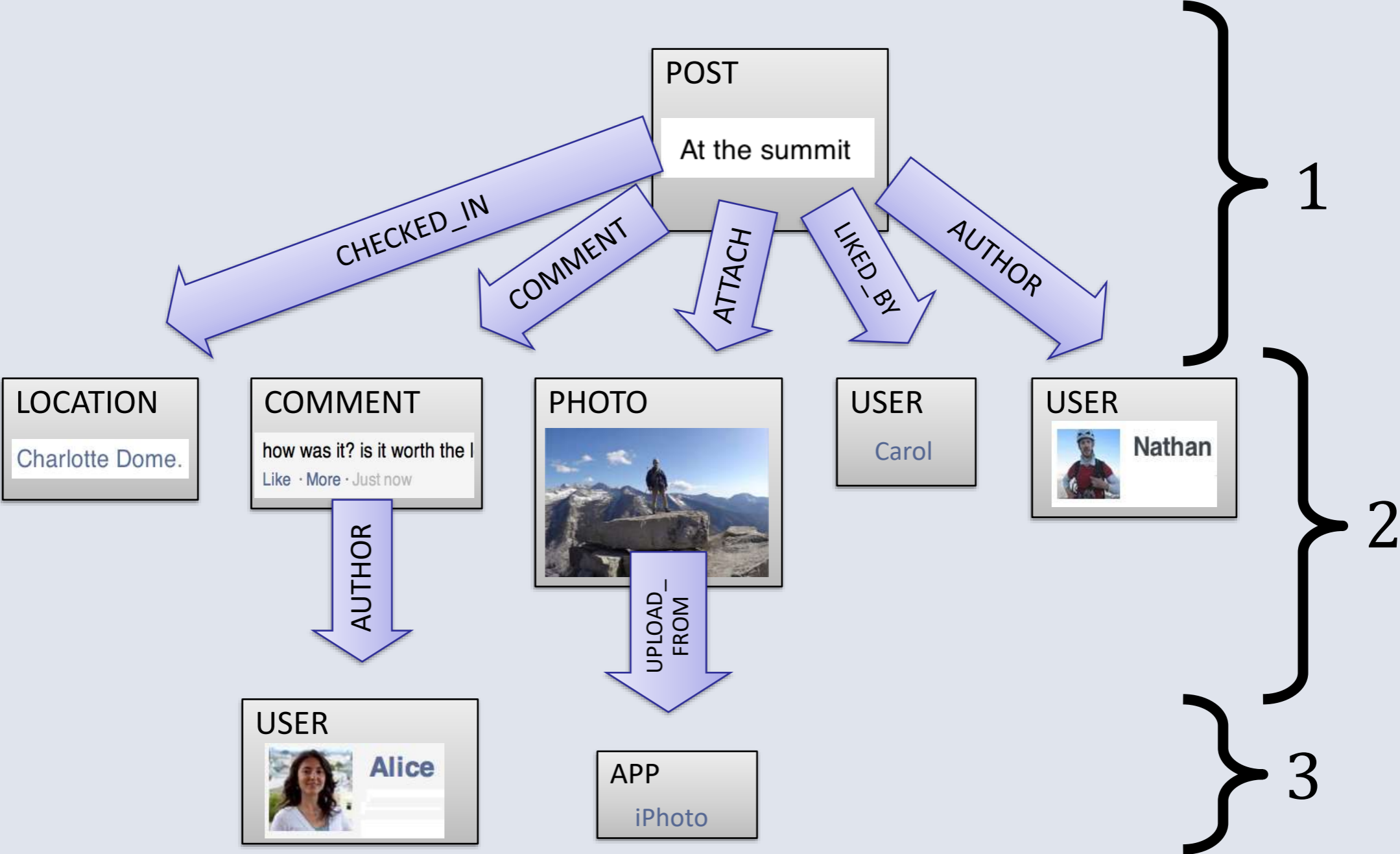
facebook data centers



What Are TAO's Goals/Challenges?

- Efficiency at scale

Dynamic Resolution of Data Dependencies



What Are TAO's Goals/Challenges?

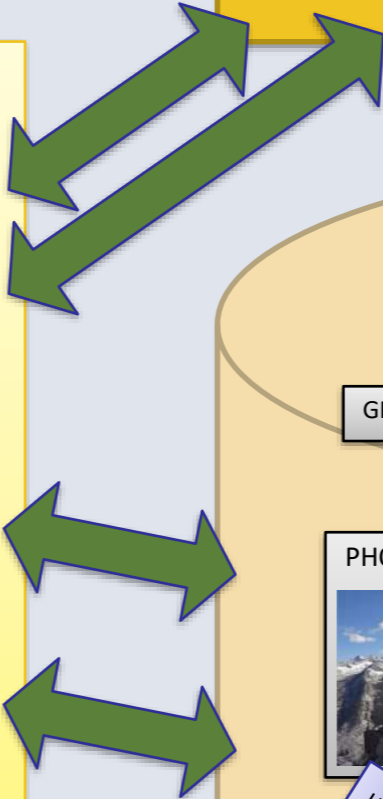
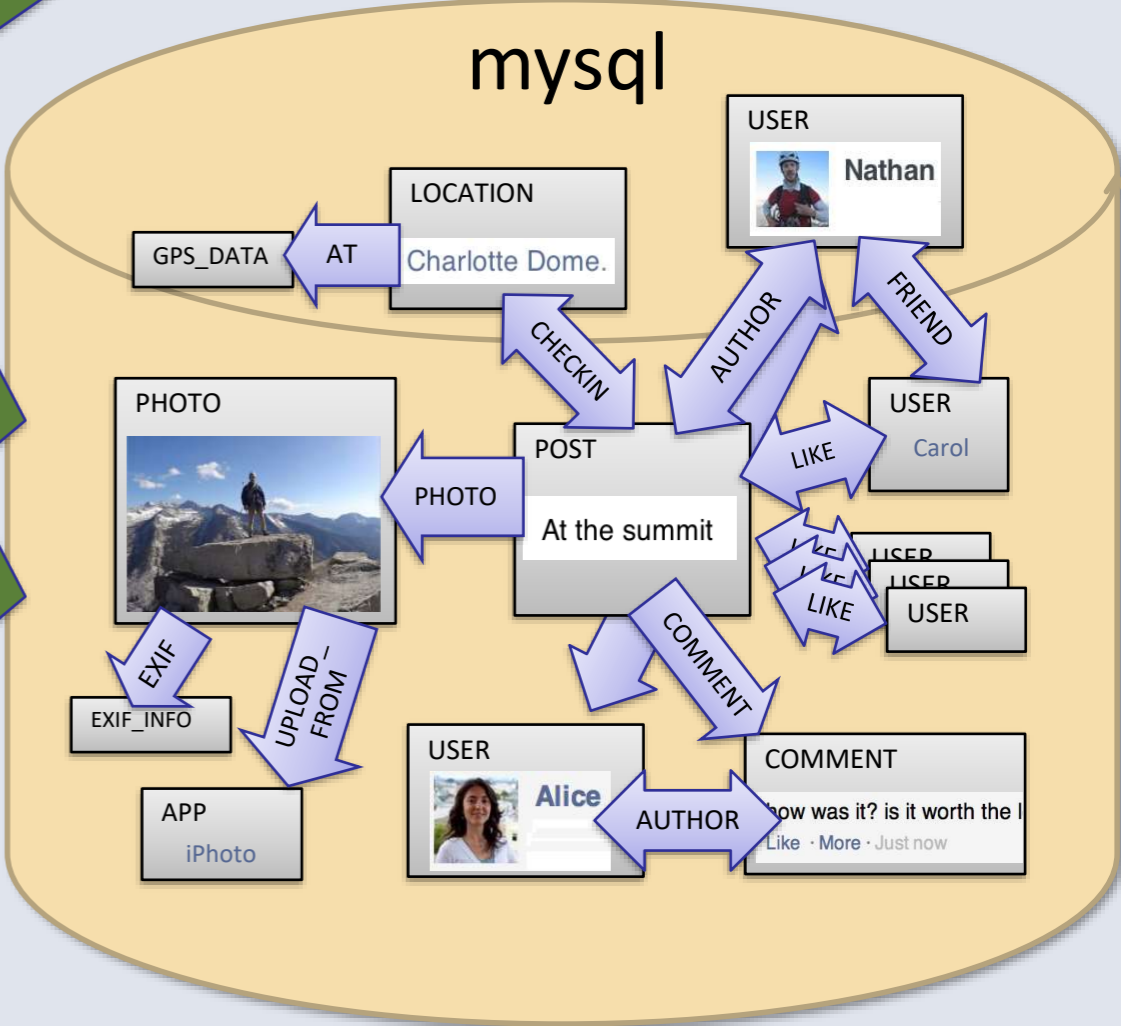
- Efficiency at scale
- Low read latency
- Timeliness of writes
- High Read Availability

Graph in Memcache

memcache
(nodes, edges, edge lists)



Web Server (PHP)
Obj & Assoc API

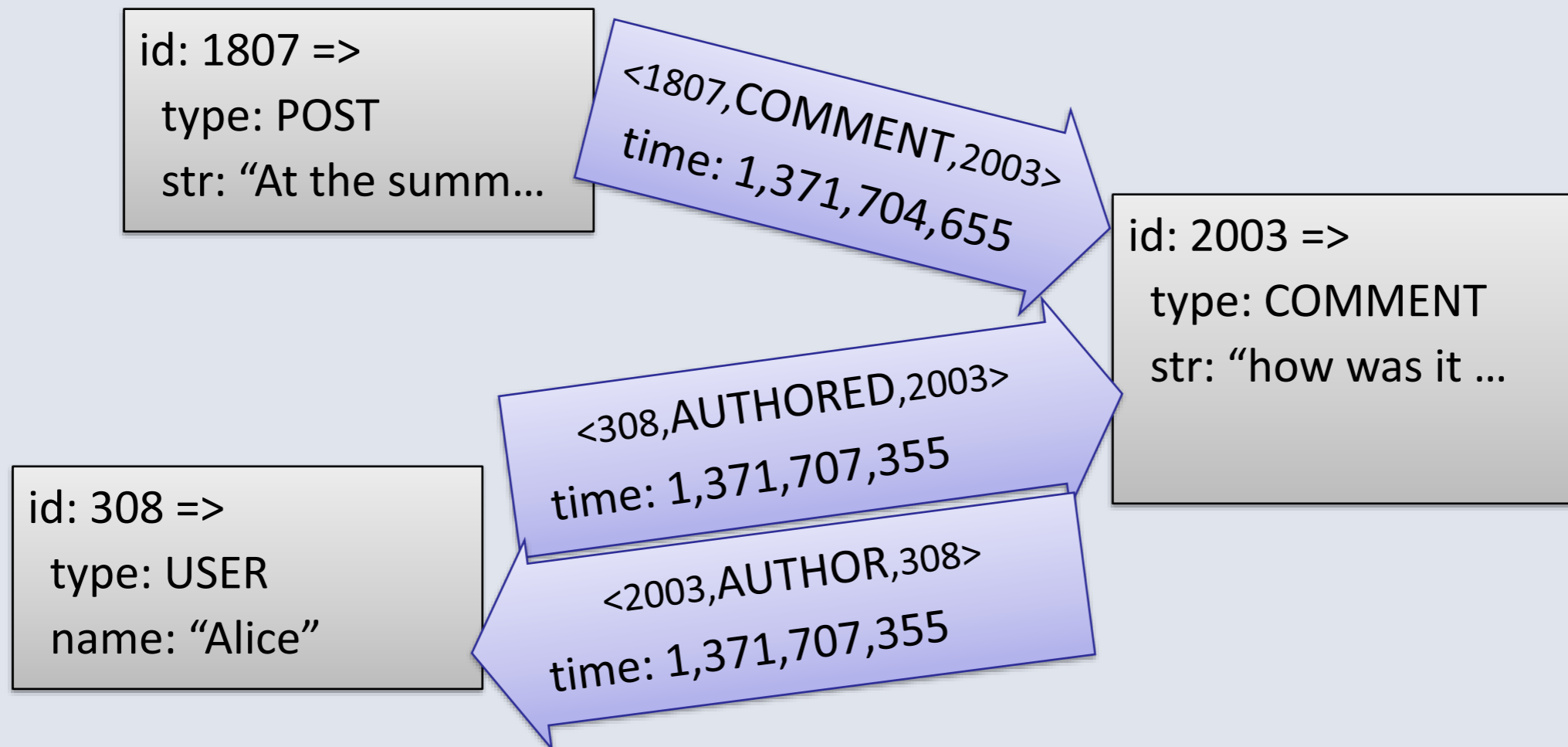


Objects = Nodes

- Identified by unique 64-bit IDs
- Typed, with a schema for fields

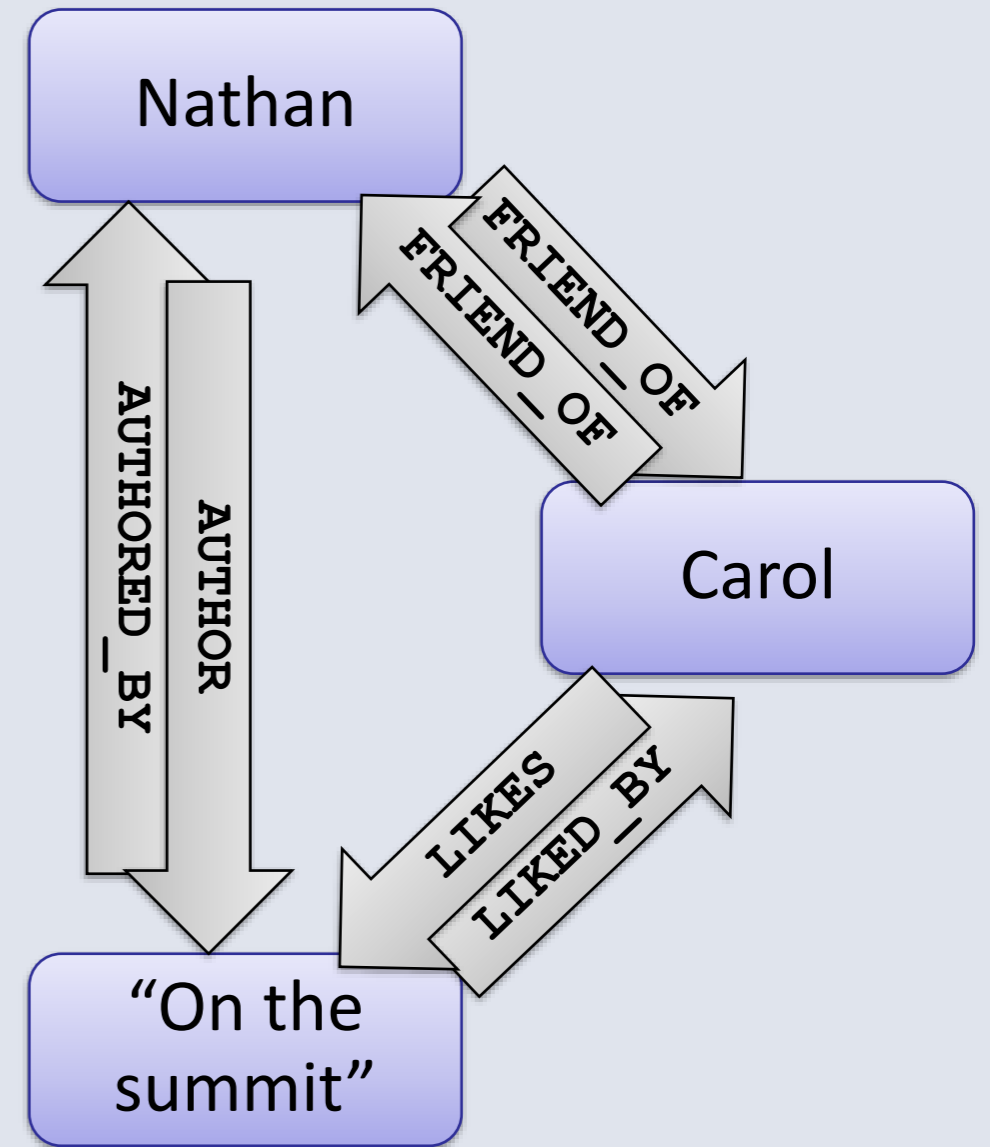
Associations = Edges

- Identified by $\langle \text{id1}, \text{type}, \text{id2} \rangle$
- Bidirectional associations are two edges, same or different type



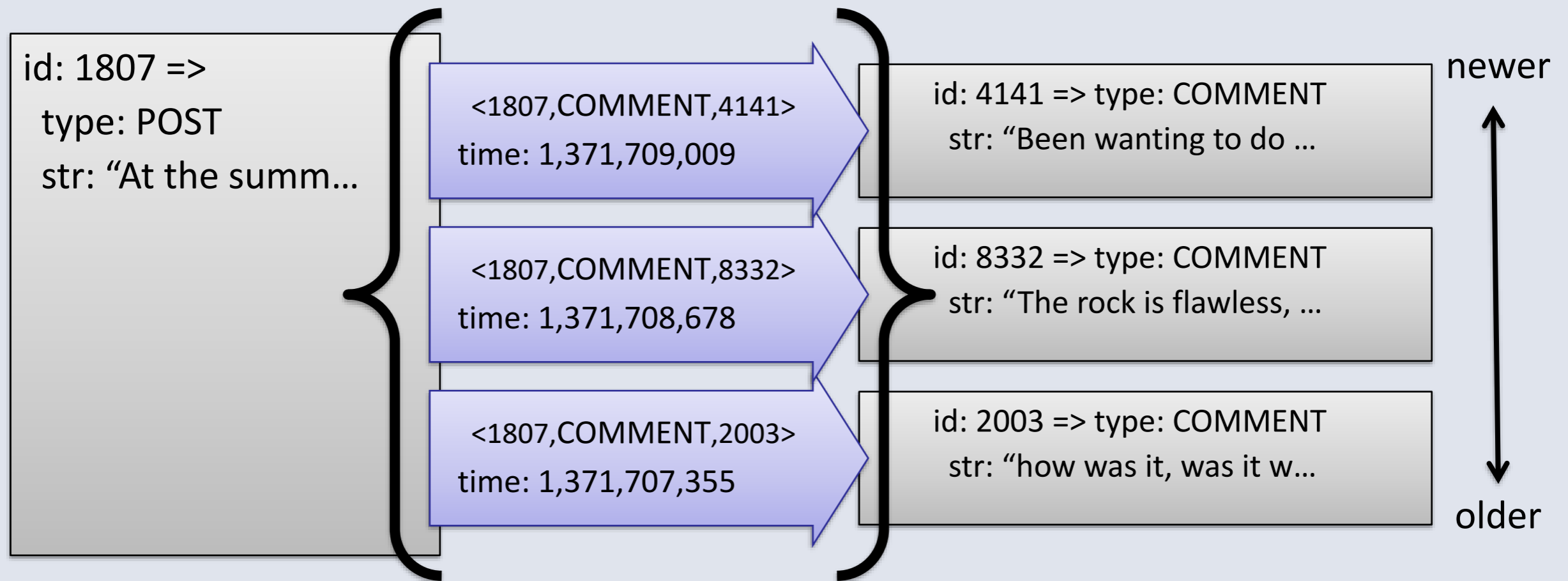
Inverse associations

- Bidirectional relationships have separate $a \rightarrow b$ and $b \rightarrow a$ edges
 - $\text{inv_type}(\text{LIKES}) = \text{LIKED_BY}$
 - $\text{inv_type}(\text{FRIEND_OF}) = \text{FRIEND_OF}$
- Forward and inverse types linked only during write
 - TAO `assoc_add` will update both
 - Not atomic, but failures are logged and repaired



Association Lists

- `<id1, type, *>`
- Descending order by time
- Query sublist by position or time
- Query size of entire list



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%

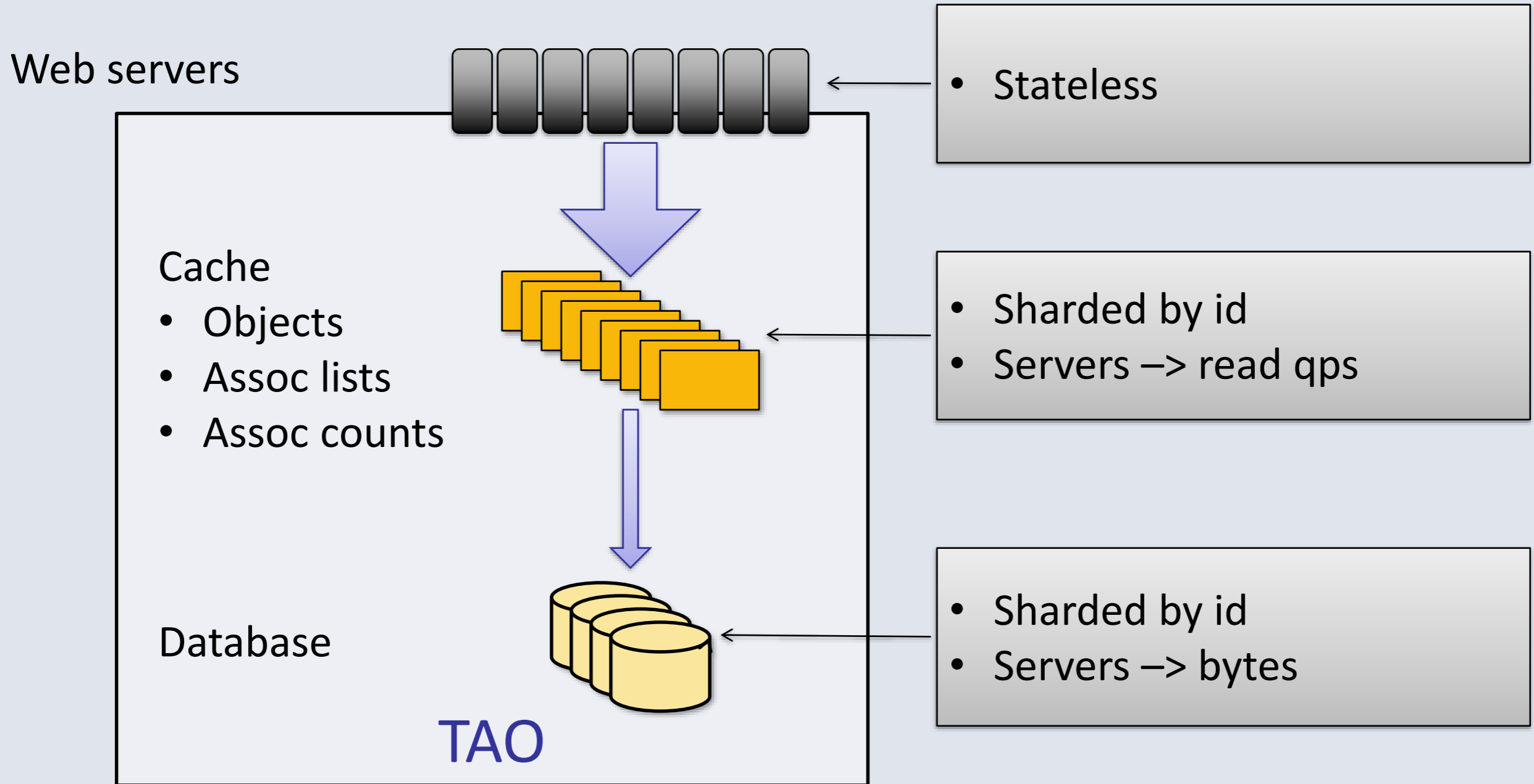
Writes – 0.2%

- 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%

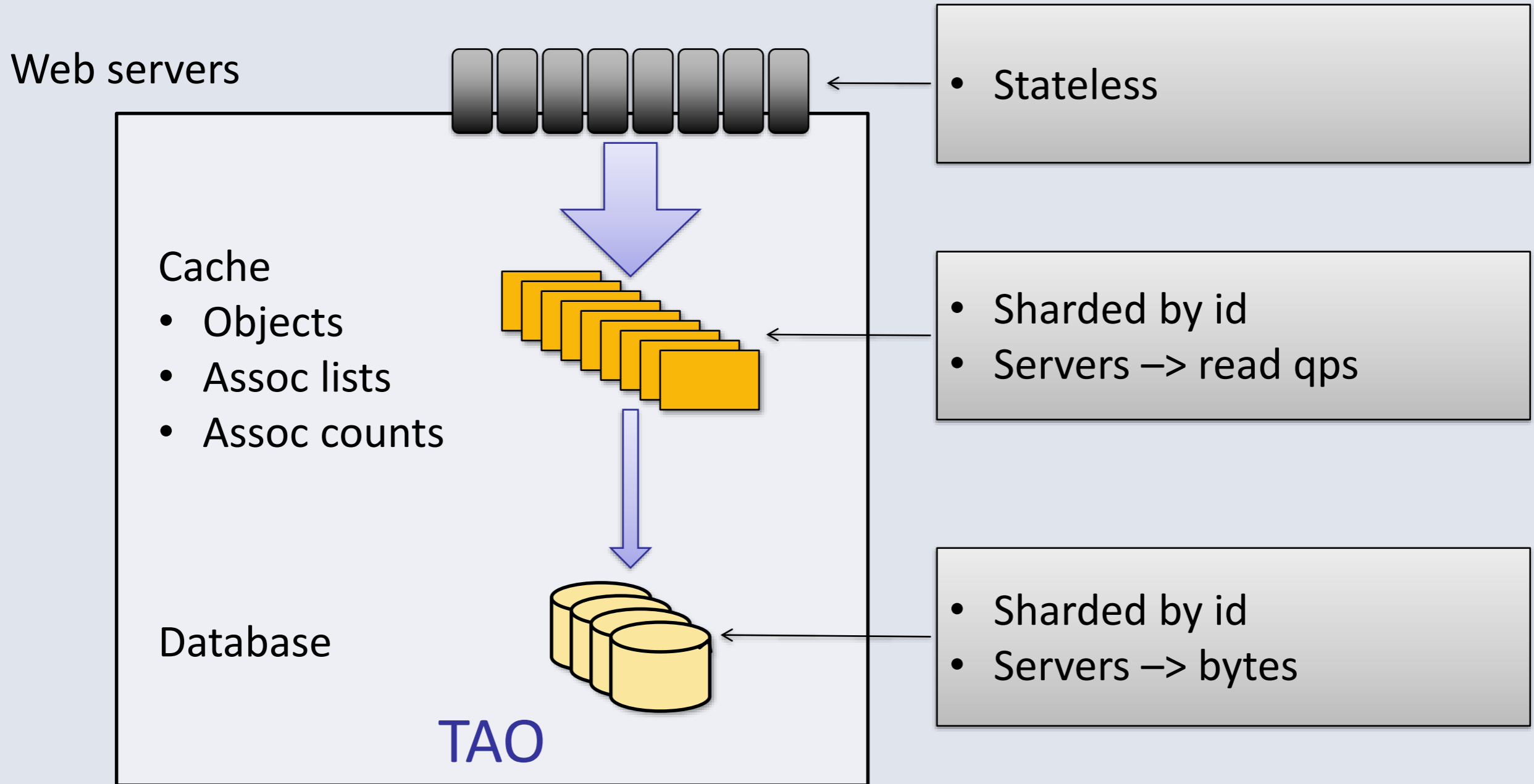
What Are TAO's Goals/Challenges?

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

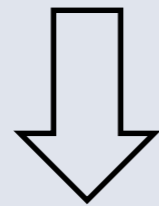
Independent Scaling by Separating Roles



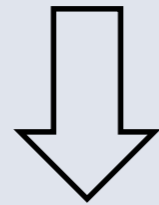
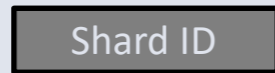
Independent Scaling by Separating Roles



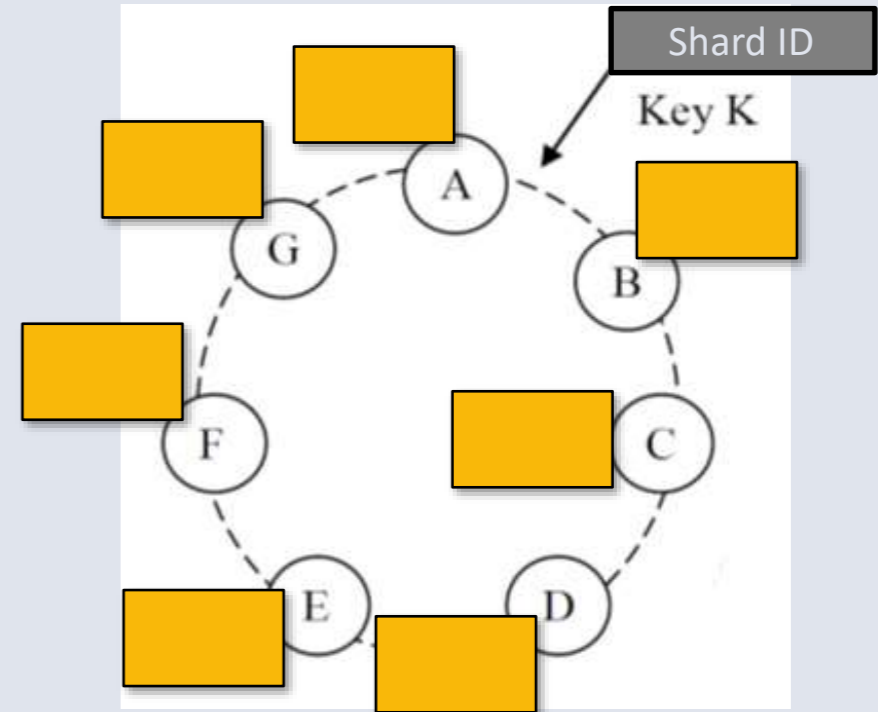
Object ID -> Shard -> Cache node



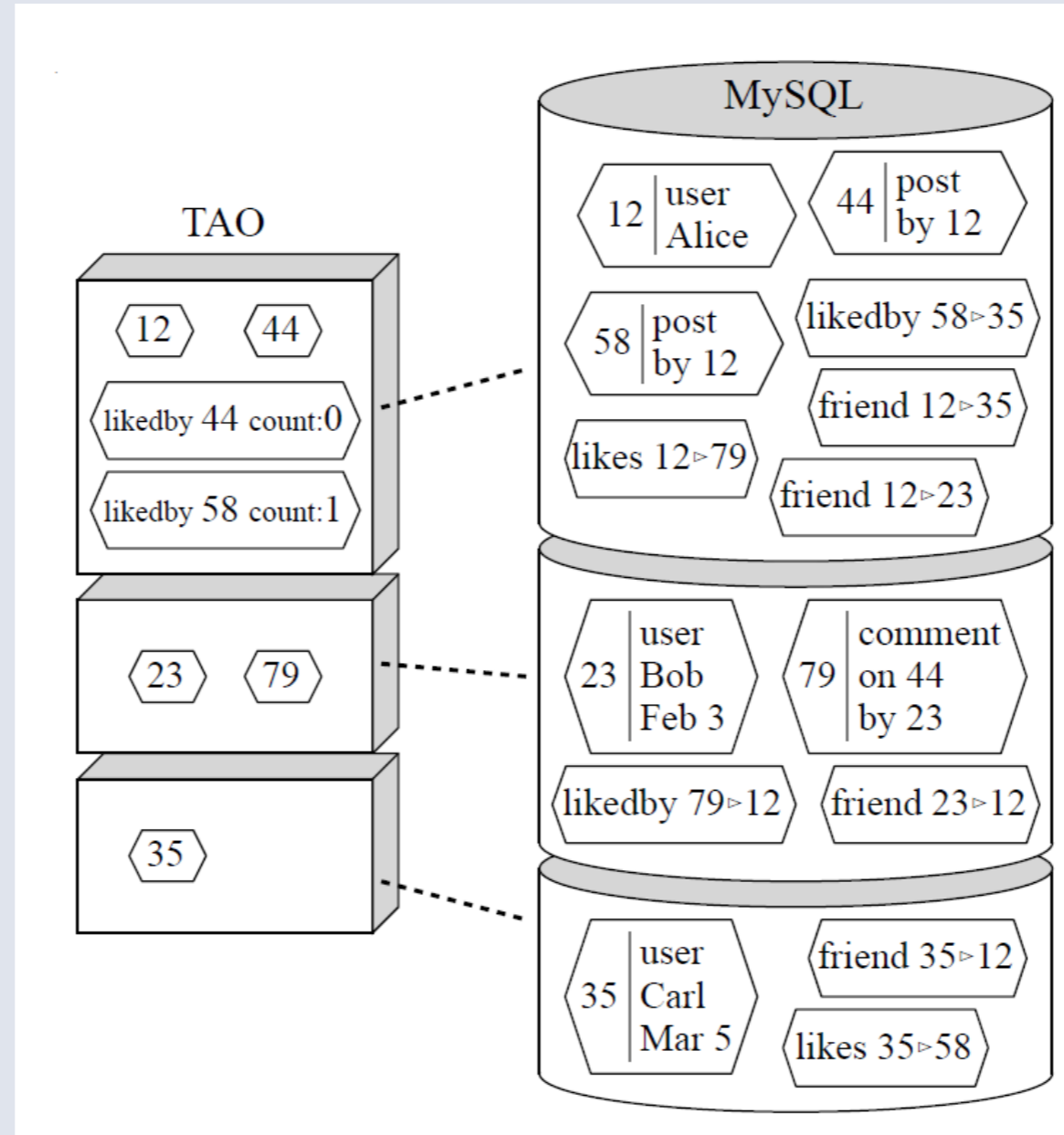
Static assignment, at creation time

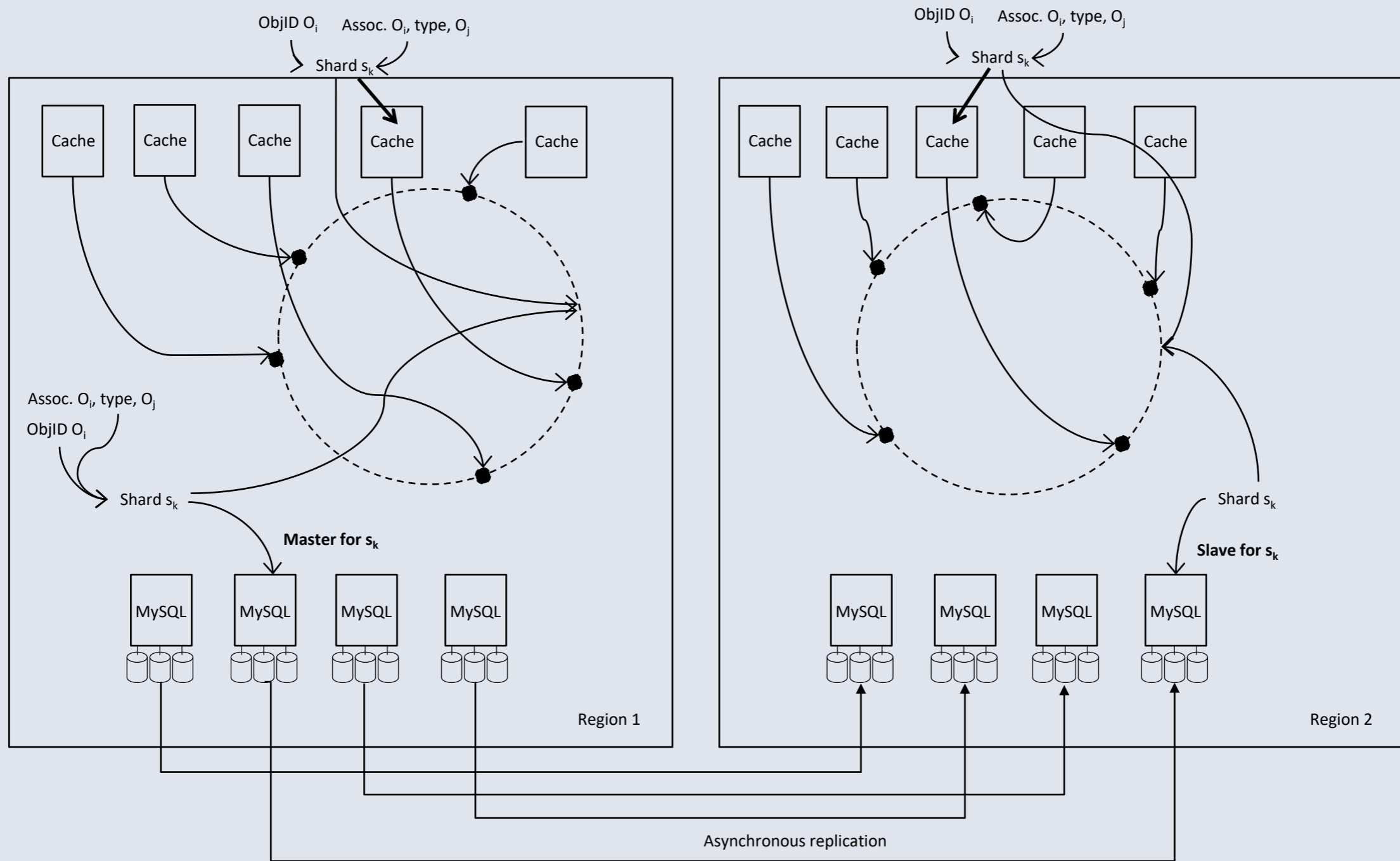


Consistent hashing



Sharding in Tao





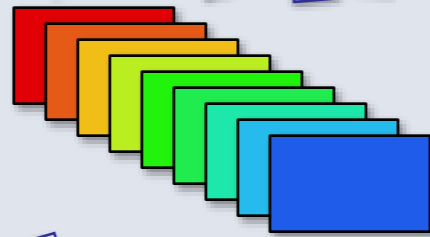
Subdividing the Data Center

Web servers



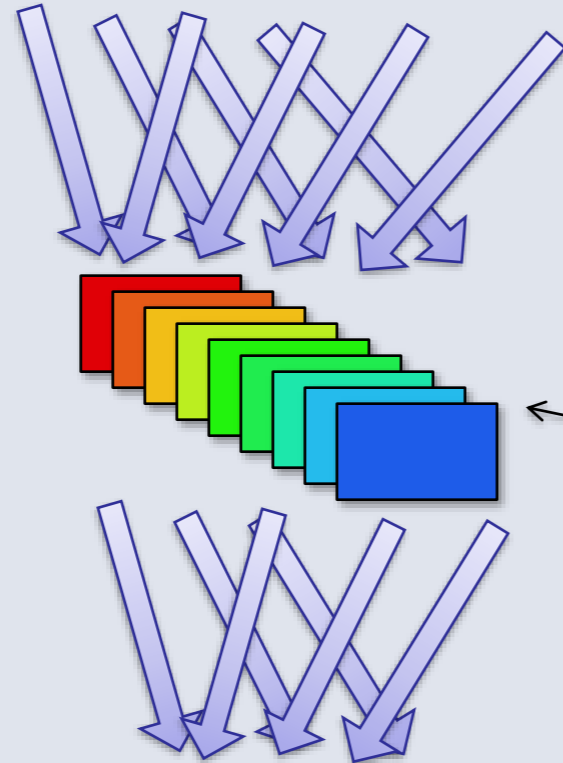
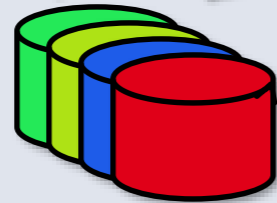
- Inefficient failure detection
- Many switch traversals

Cache

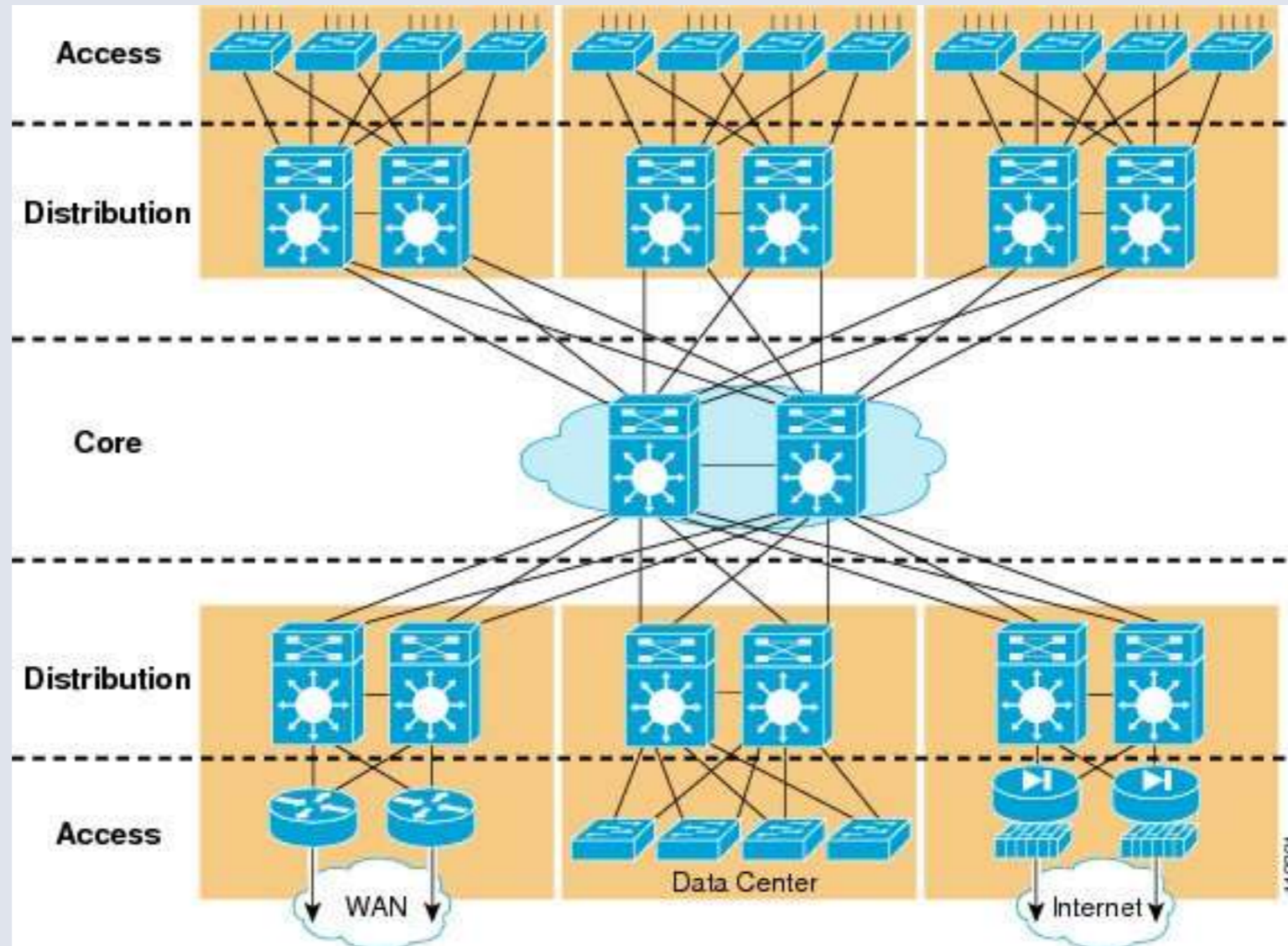


- Many open sockets
- Lots of hot spots

Database

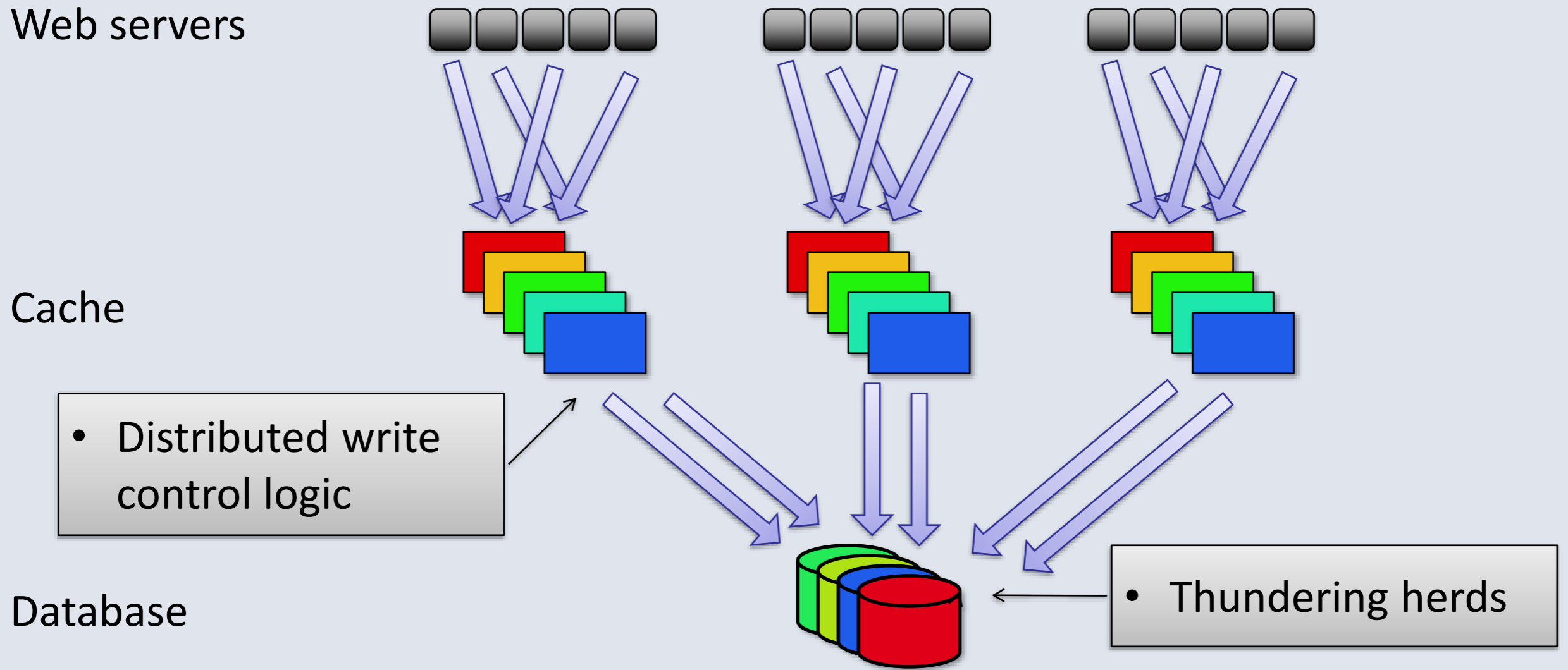


Network design

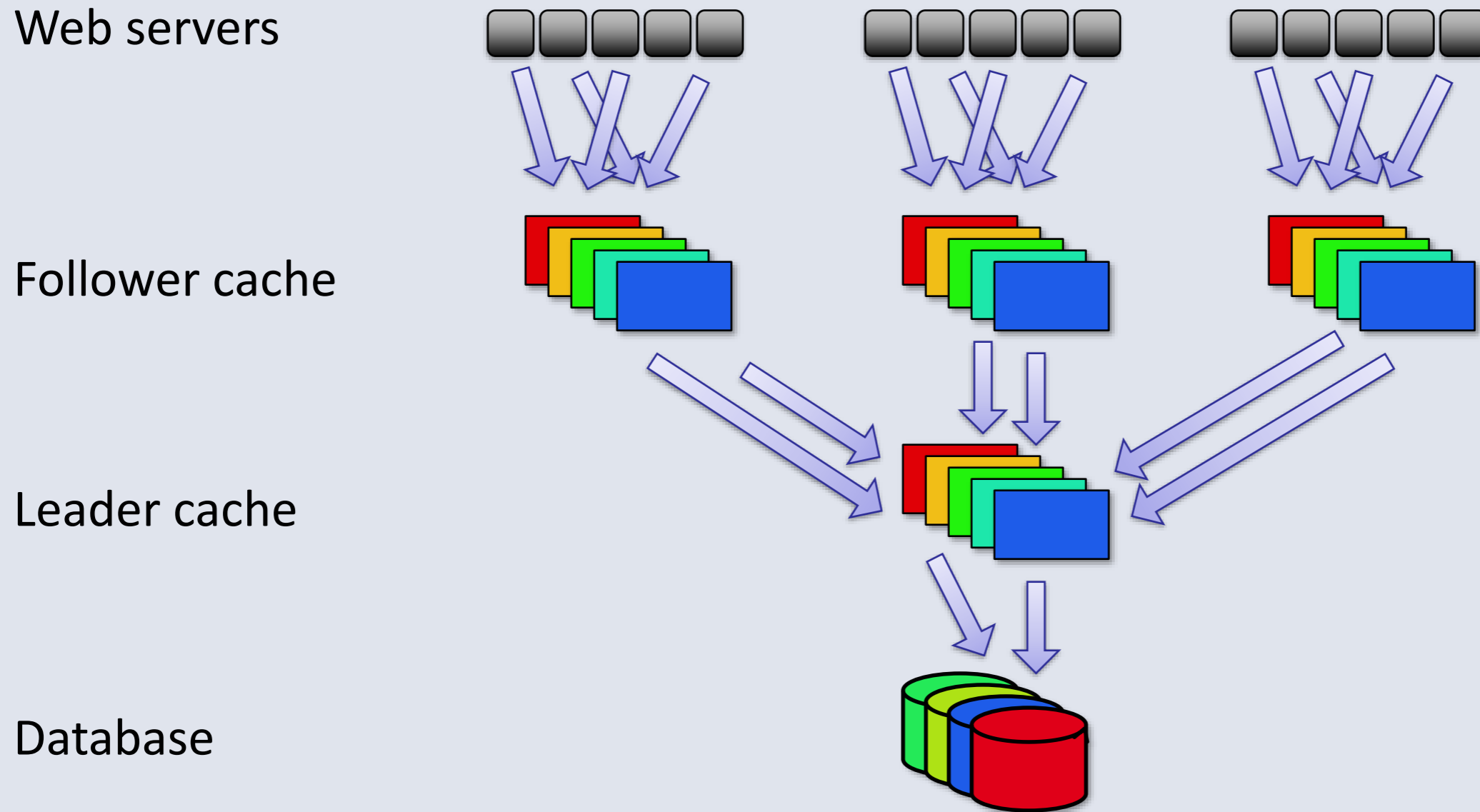


Source: http://www.cisco.com/c/en/us/td/docs/solutions/Enterprise/Campus/HA_campus_DG/hacampusdg.html

Subdividing the Data Center



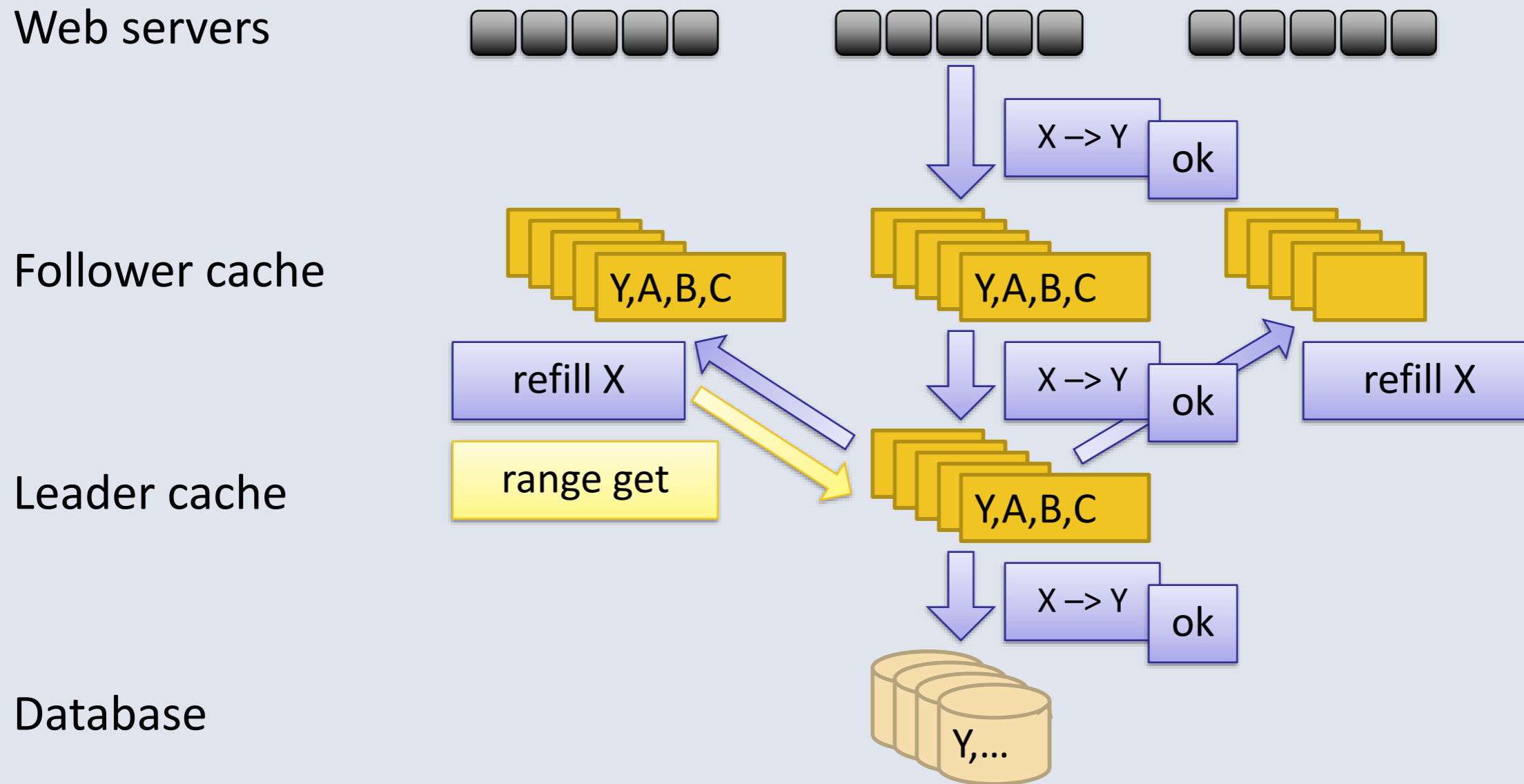
Follower and Leader Caches



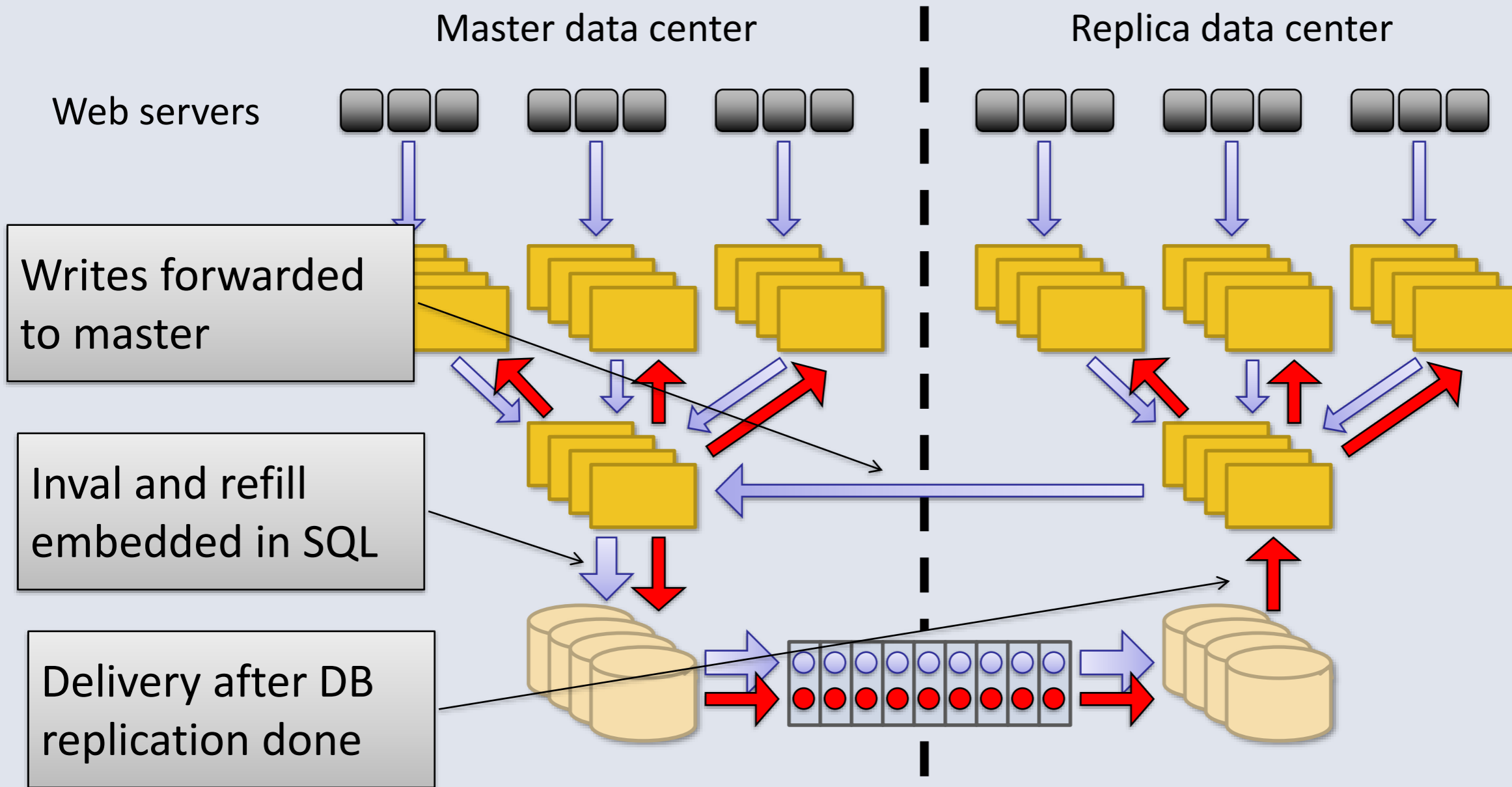
What Are TAO's Goals/Challenges?

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

Write-through Caching – Association Lists



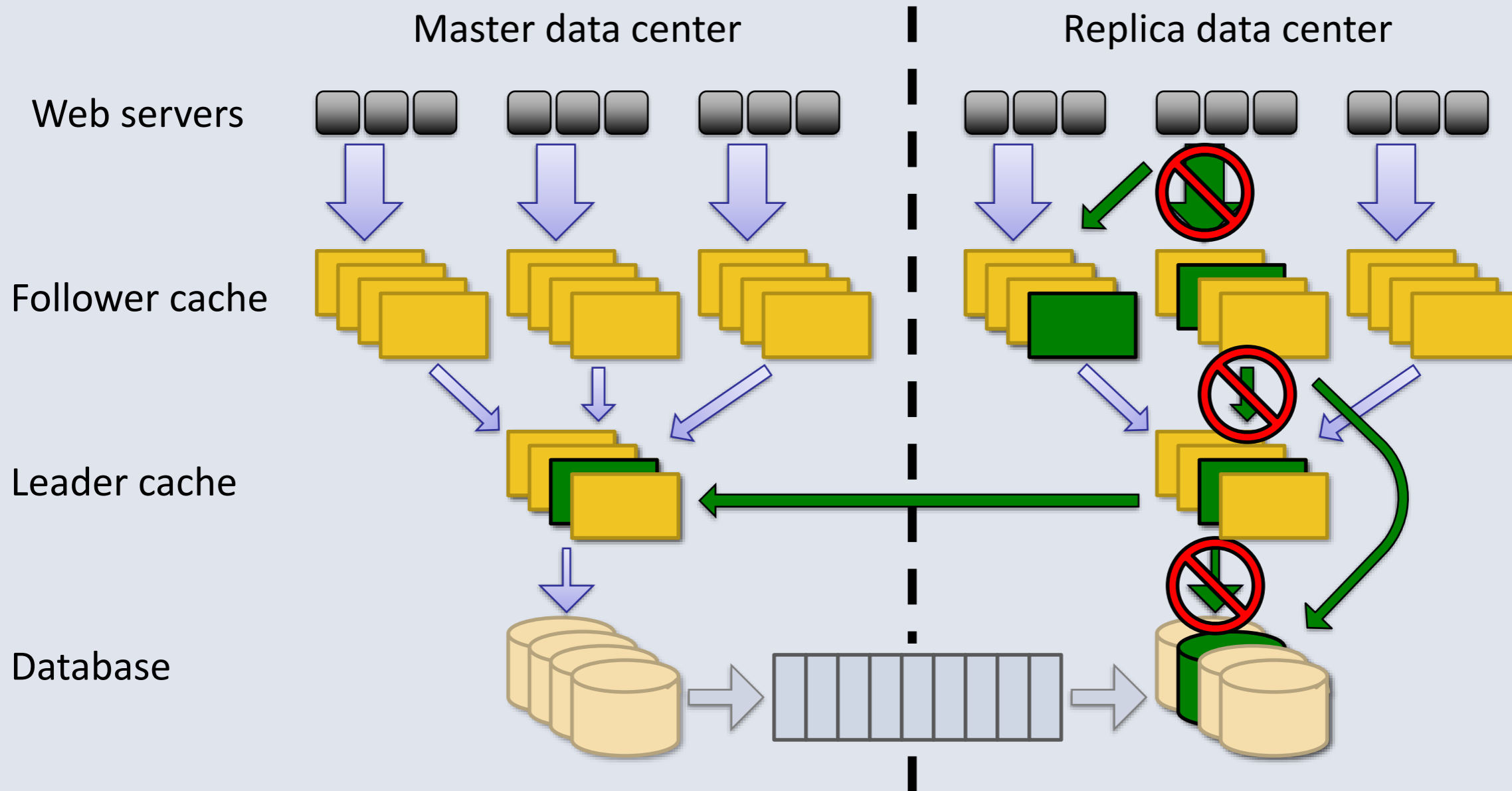
Asynchronous DB Replication



What Are TAO's Goals/Challenges?

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

Improving Availability: Read Failover



TAO Summary

Efficiency at scale
Read latency

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

Write timeliness

- Write-through cache
- Asynchronous replication

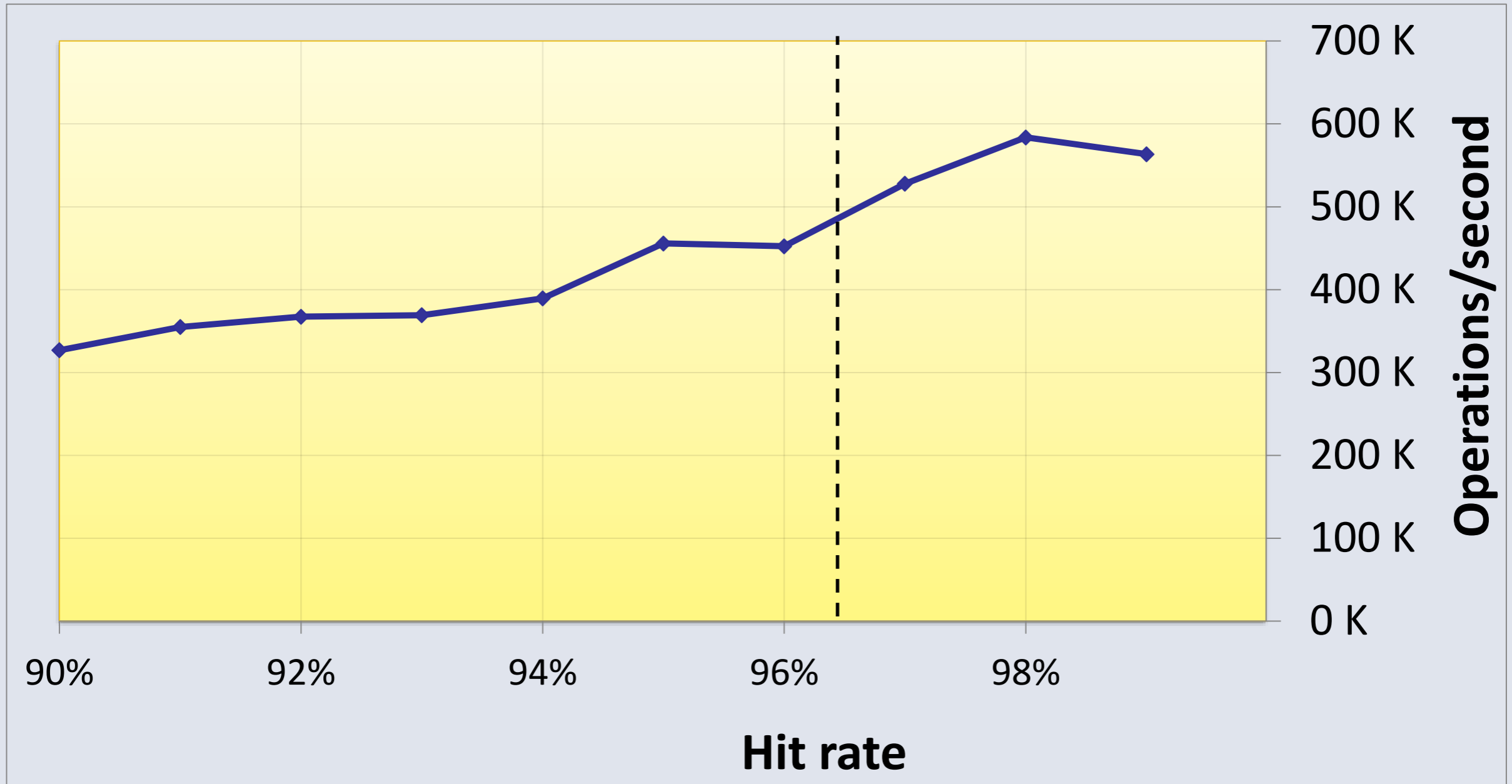
Read availability

- Alternate data sources

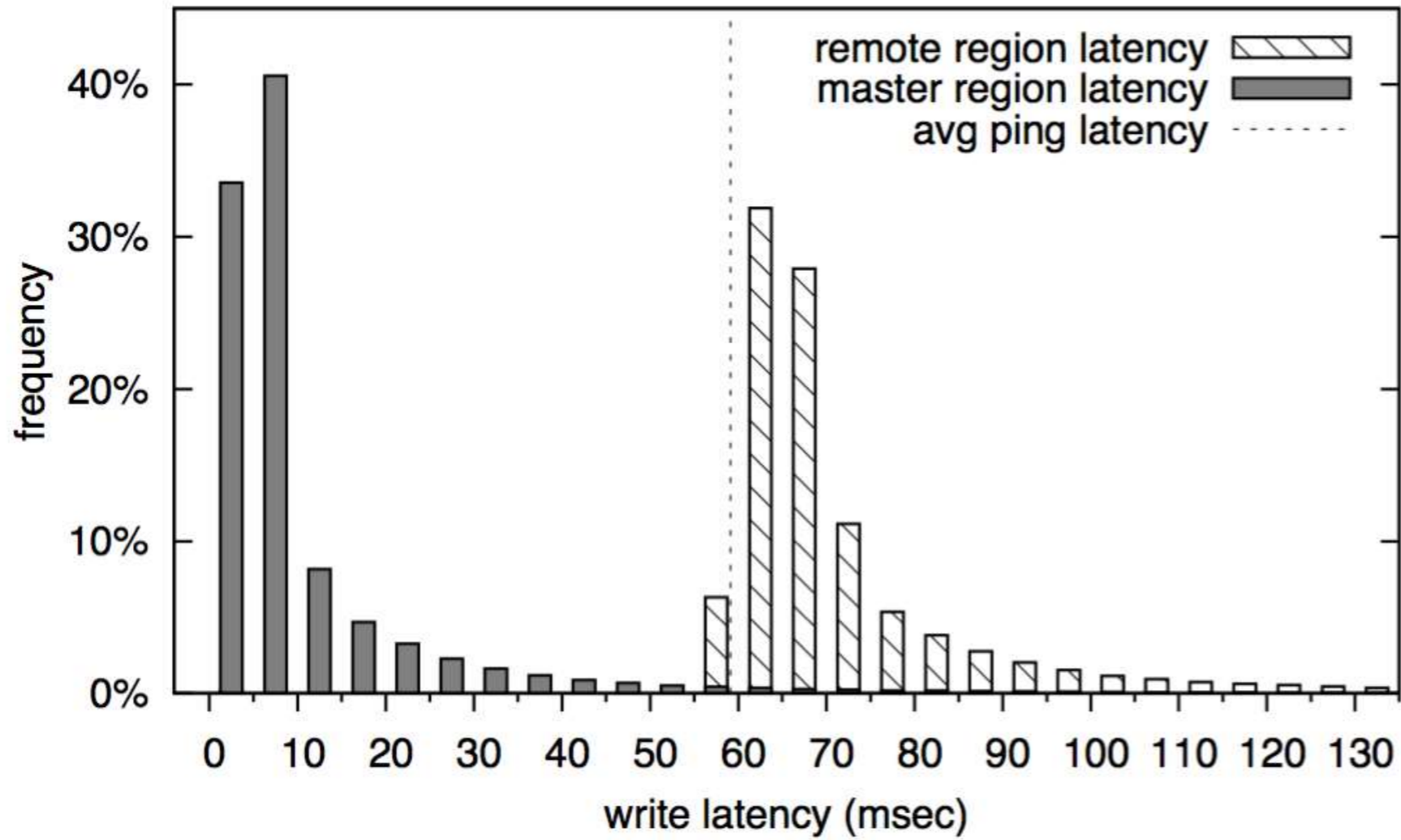
facebook

(c) 2009 Facebook, Inc. or its licensors. "Facebook" is a registered trademark of Facebook, Inc.. All rights reserved. 1.0

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

Alice was at the Golden Gate Bridge with Bob
 Cathy : Wish we were there! David likes this

