

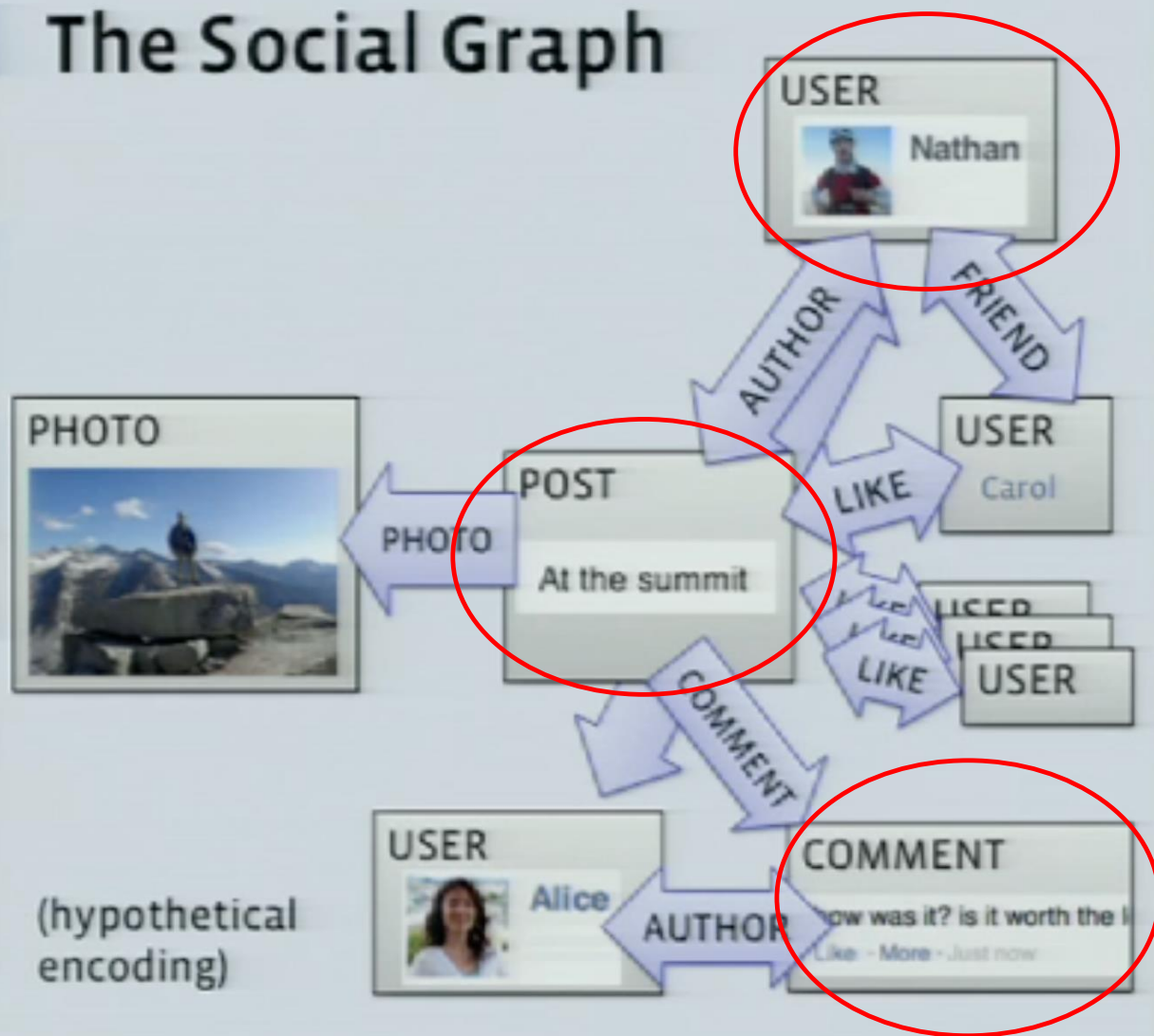
TAO: Facebook's Distributed Data Store for the Social Graph

Authors: Nathan Bronson, *et. al*

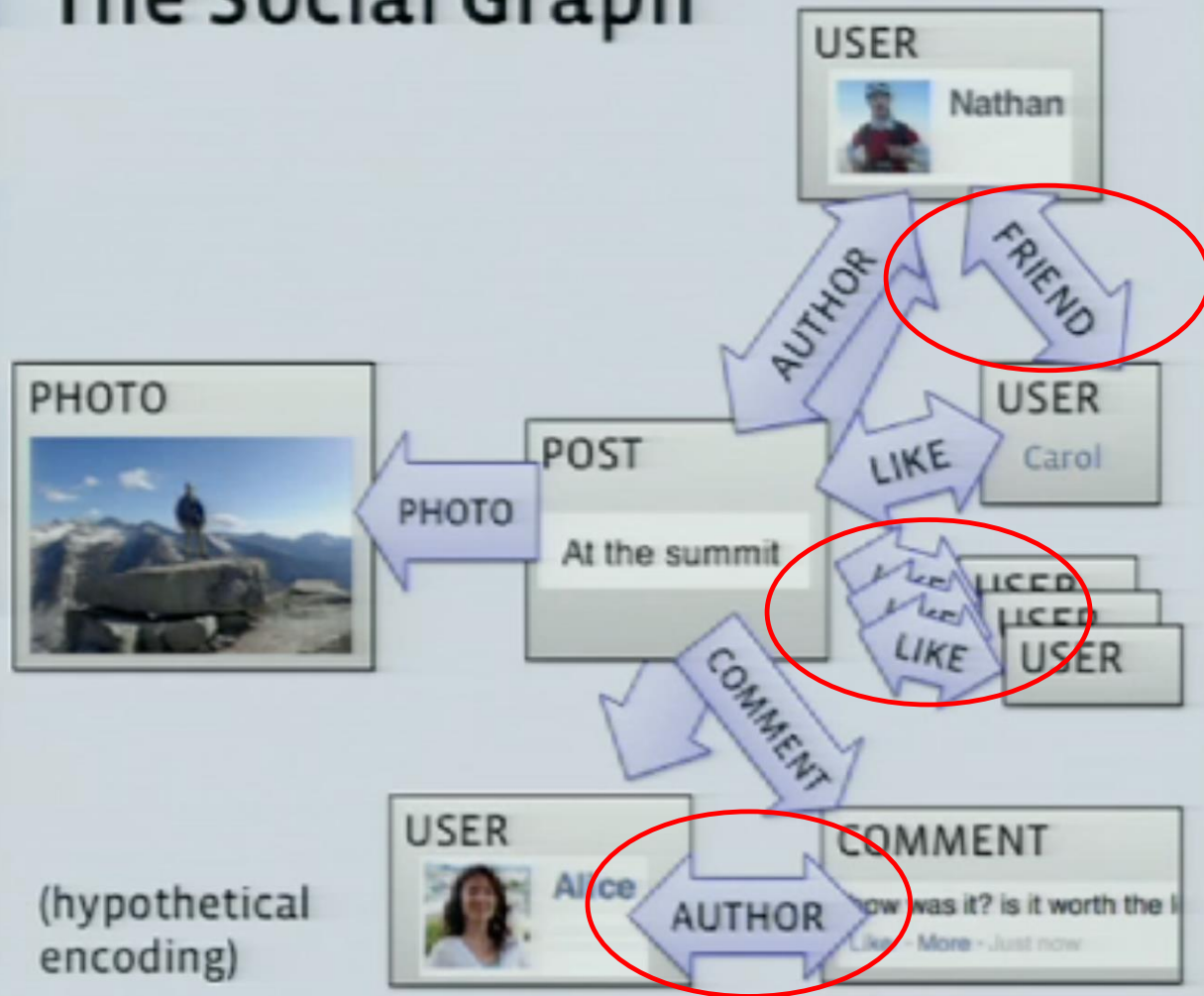
Presentors: Zhichao Cao, Fenggang Wu

10/01/2018

The Social Graph



The Social Graph

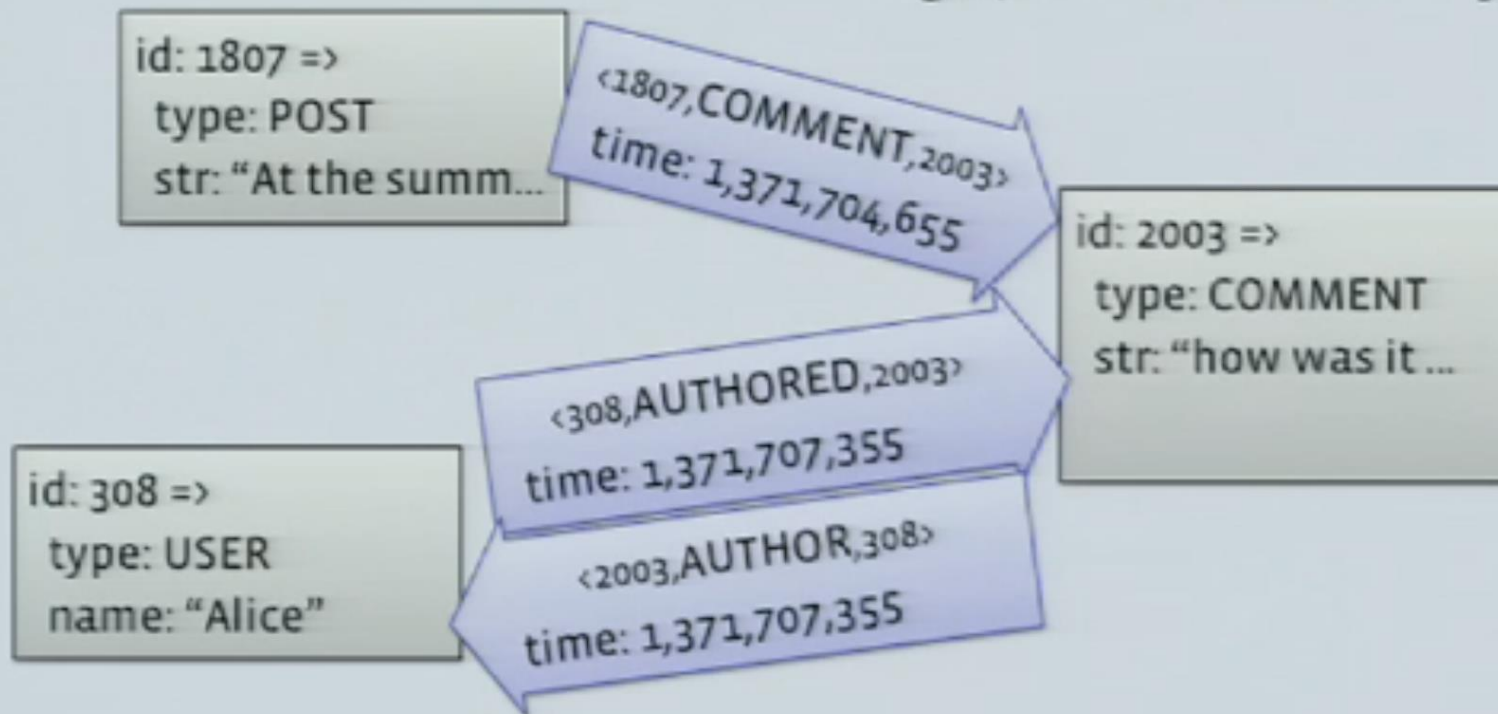


Objects = Nodes

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

Associations = Edges

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



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%

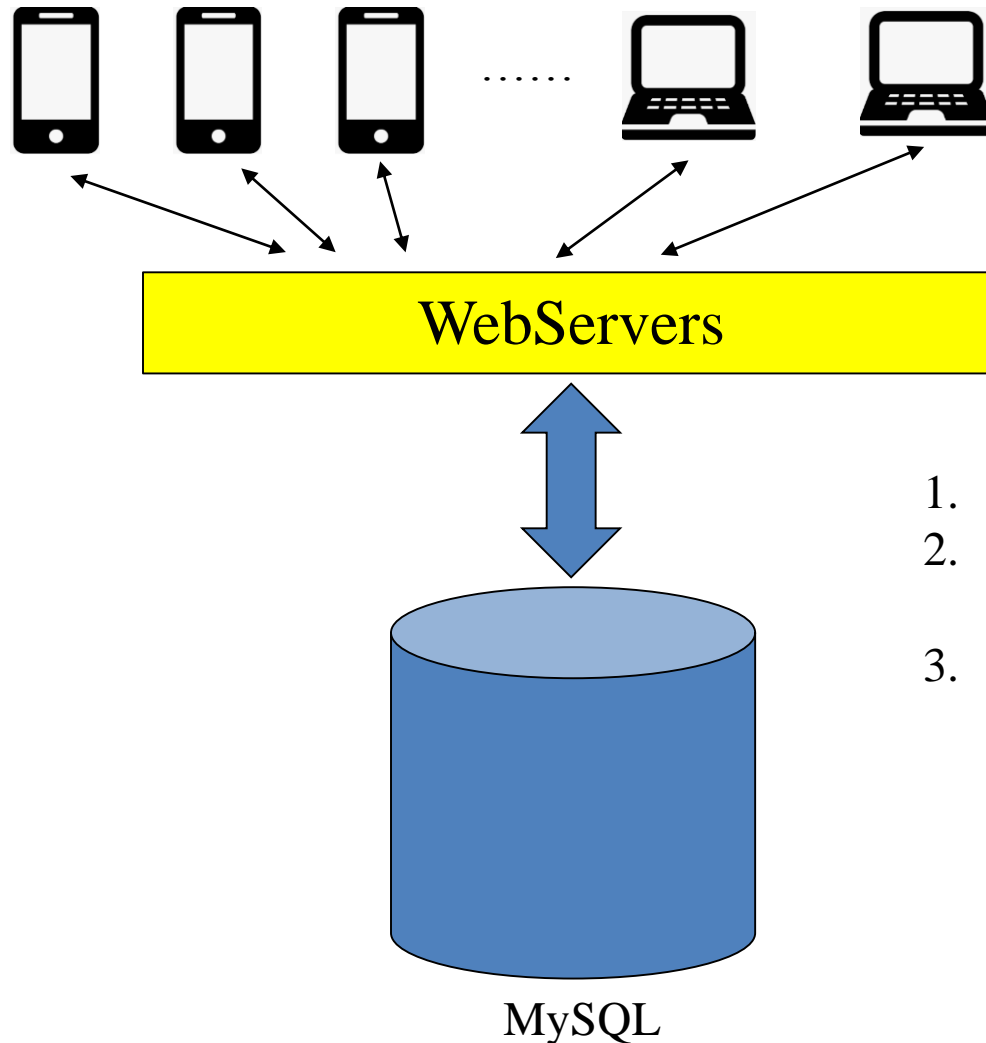
- Create, update, delete for objects

- `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%

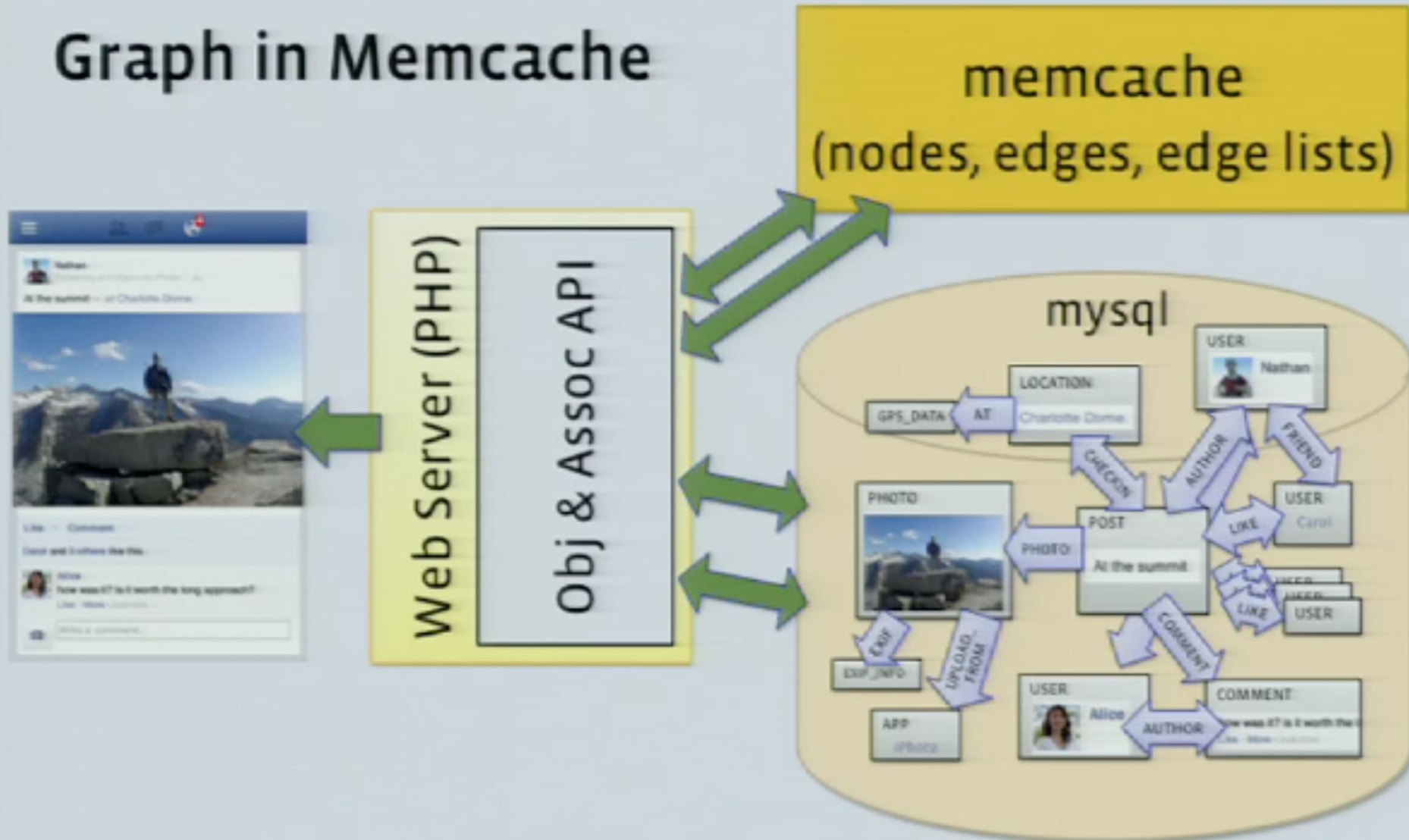
Stage 1: Relational Database Supported



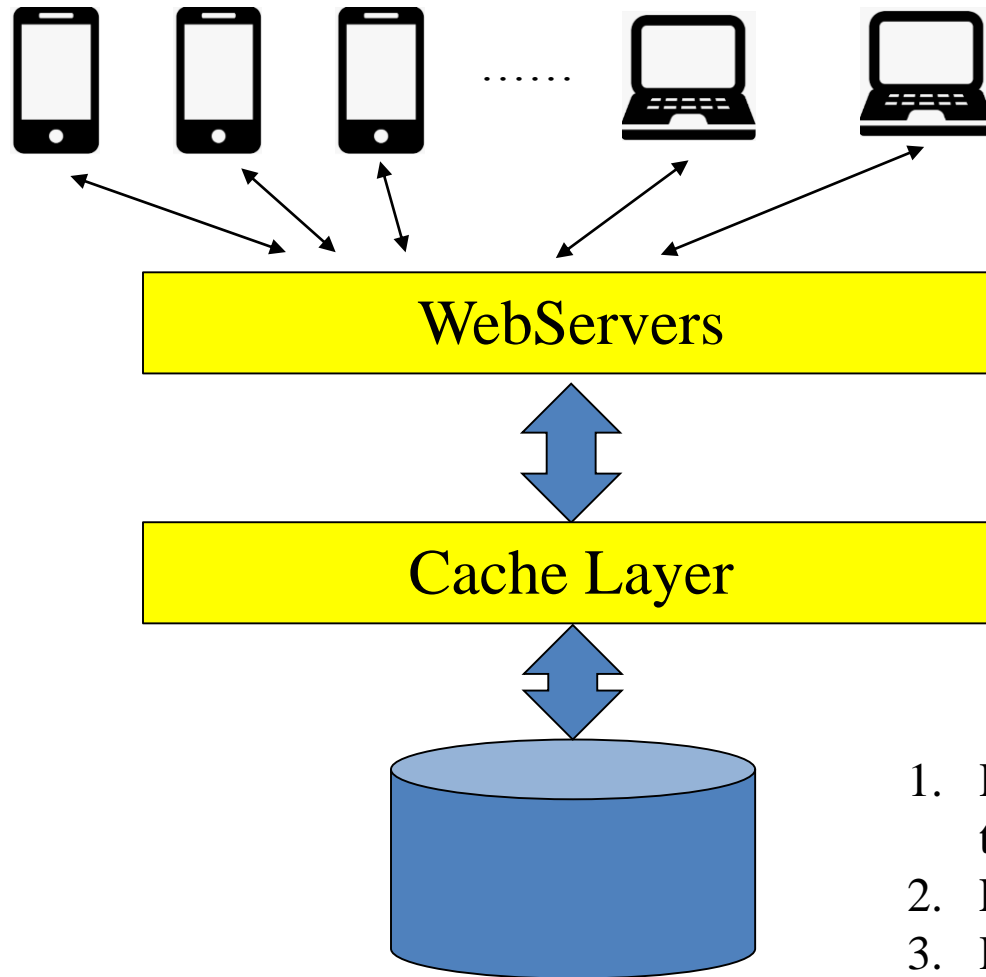
1. Low write latency
2. High read latency, difficult to achieve load balancing
3. Low scalability

Stage 2: Cache Supported (memcache)

Graph in Memcache



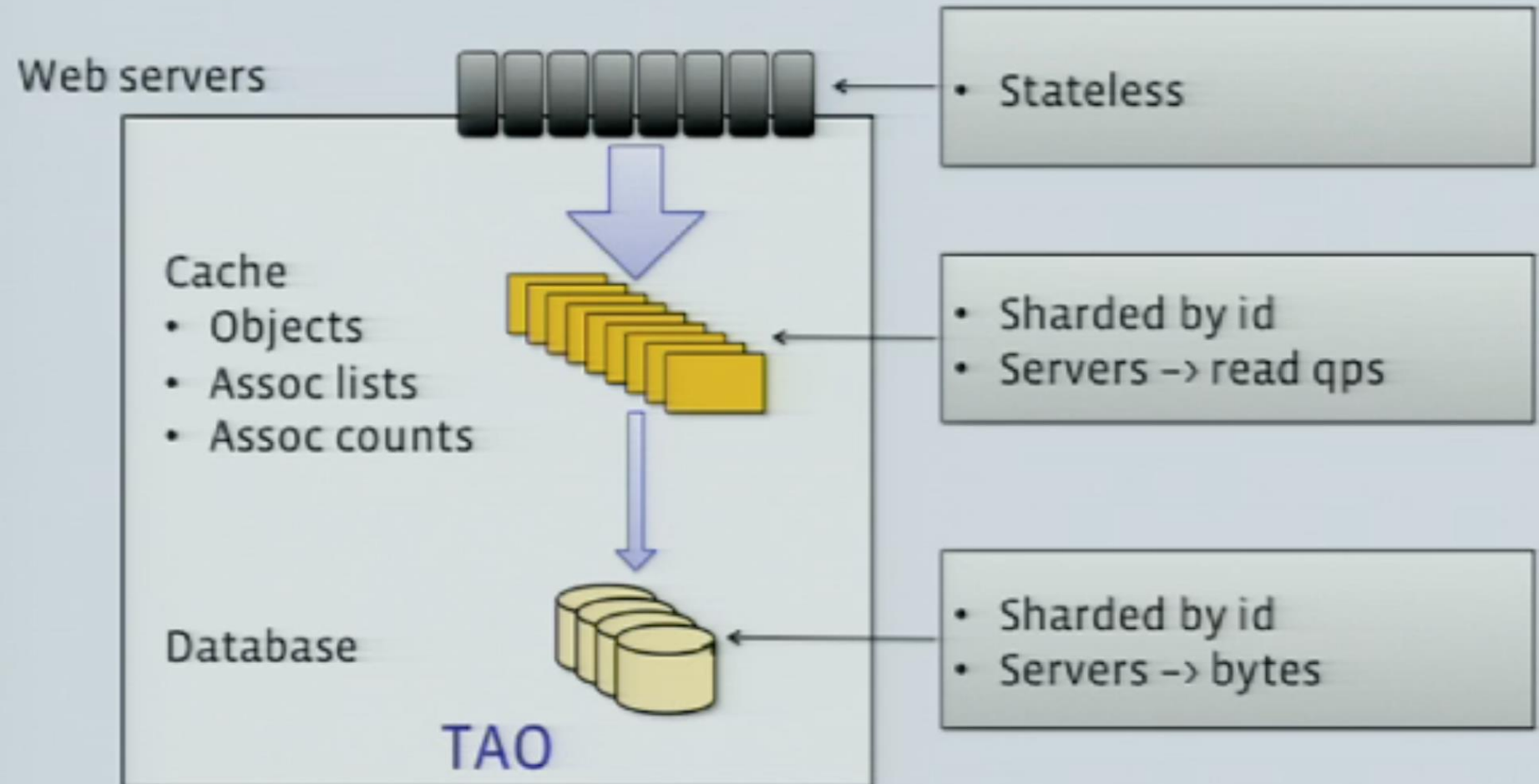
Stage 2: Cache Supported (cache in between)



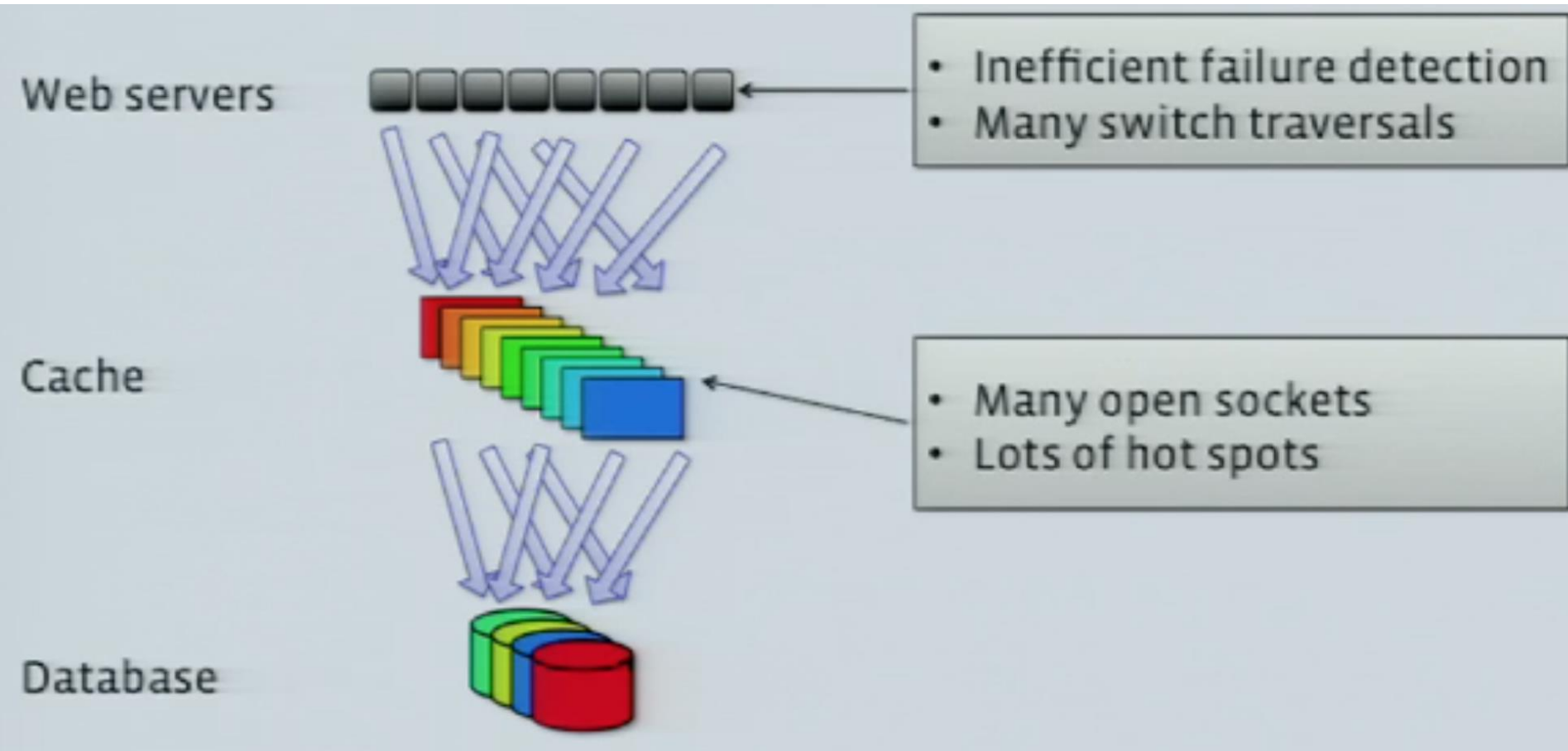
1. Higher write latency (clean the cache)
2. Lower read latency
3. Higher scalability

Stage 3: TAO (The Associations and Objects)

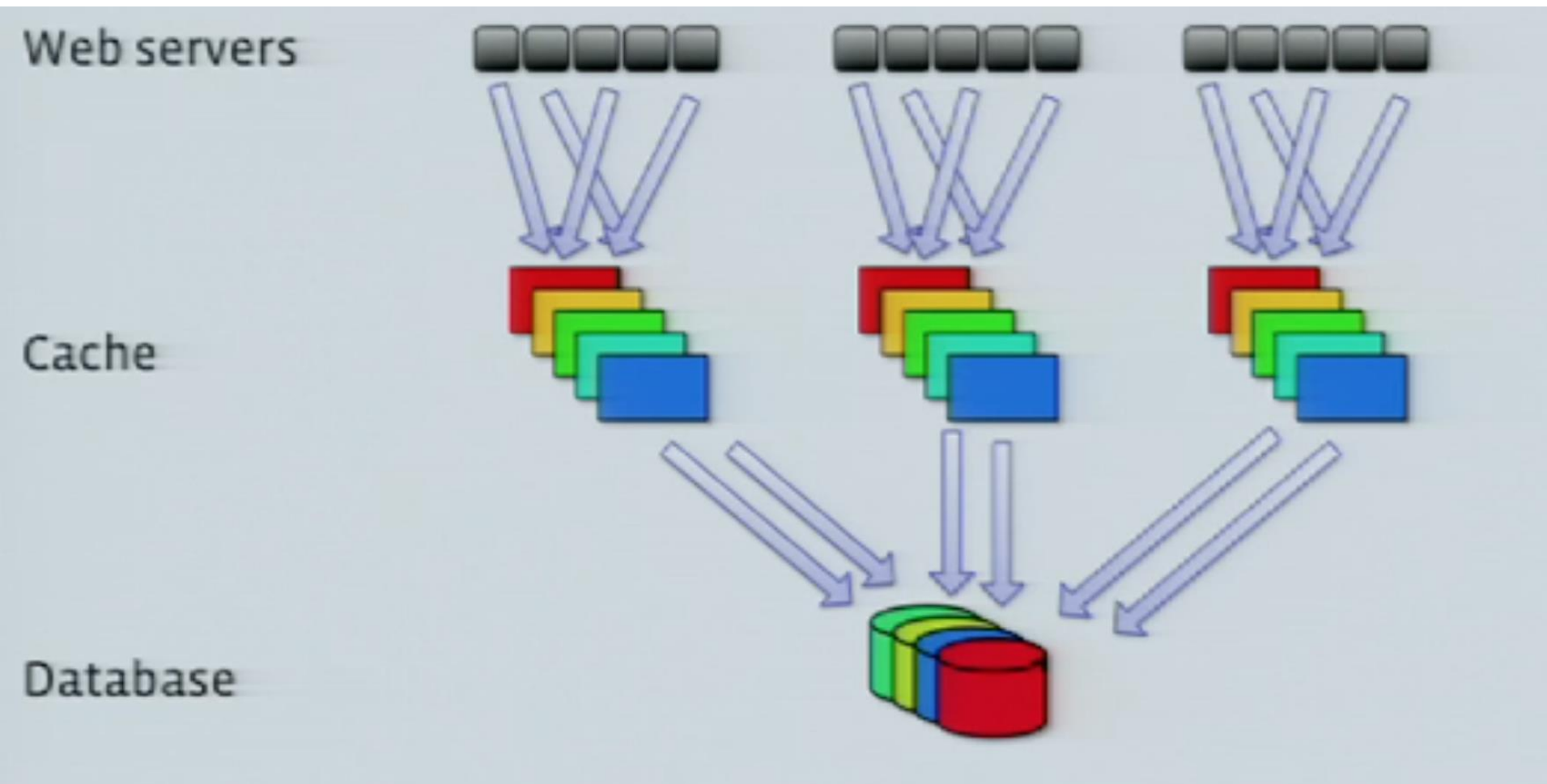
Independent Scaling by Separating Roles

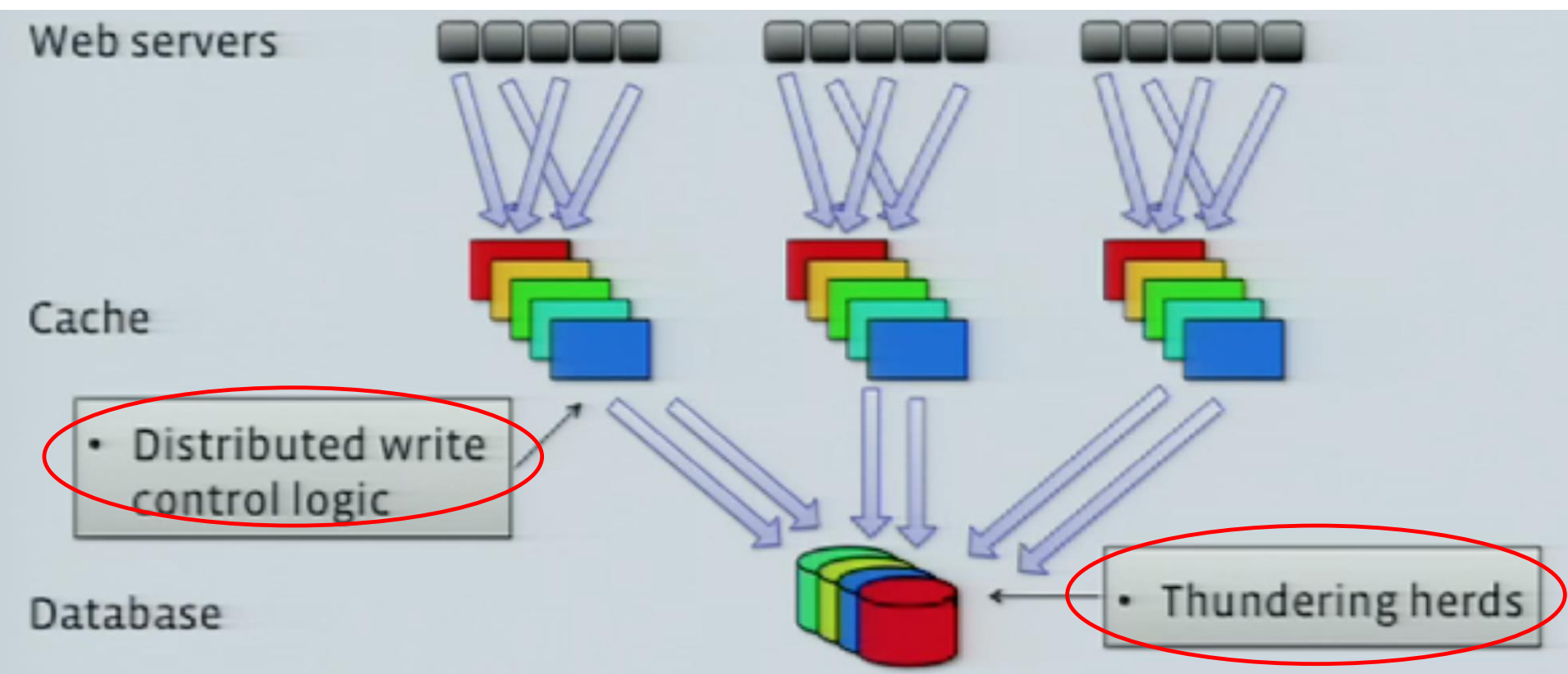


TAO Design 1: One Tier with Sharding



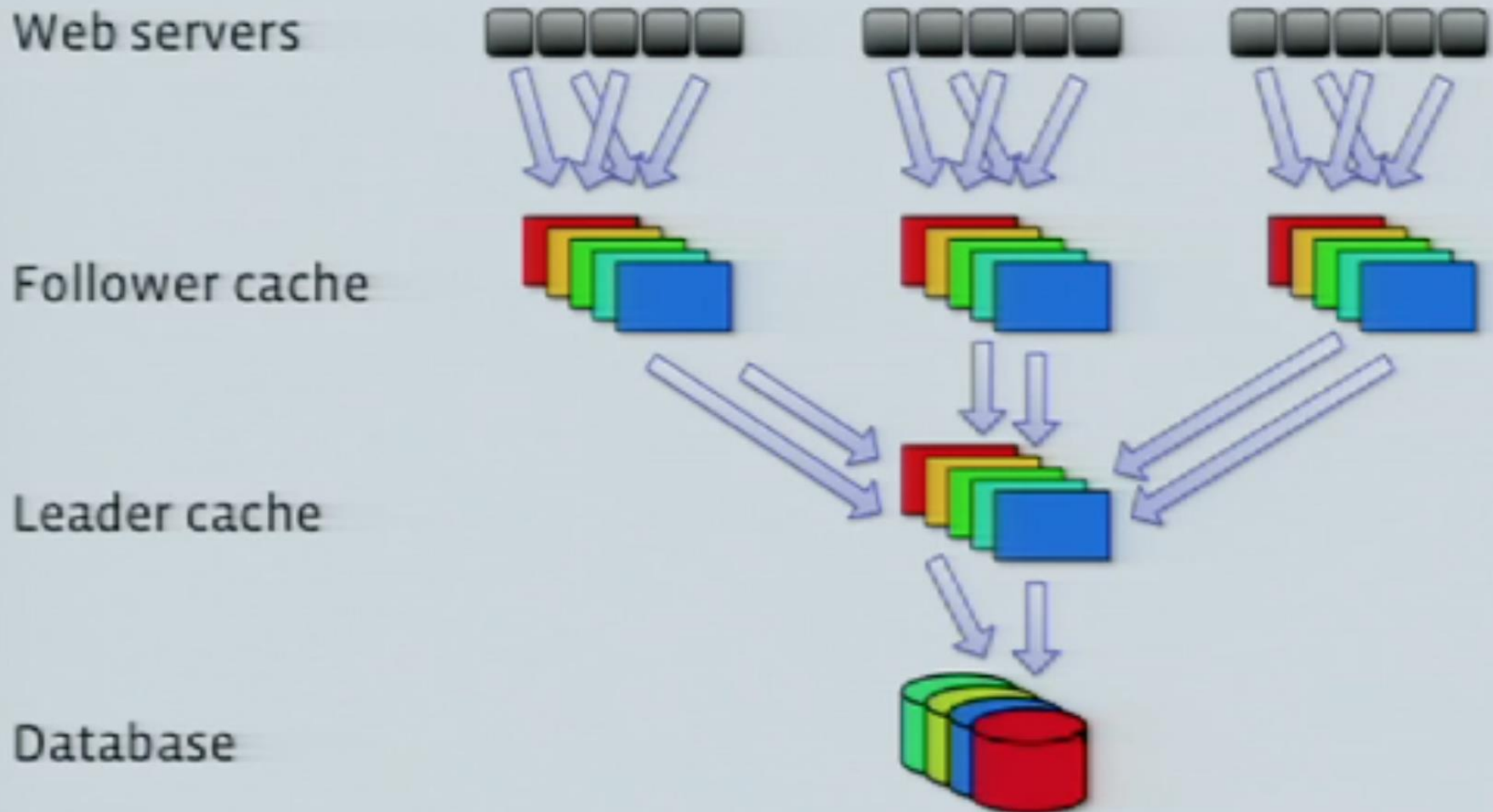
TAO Design 2: Multiple Caching Tiers



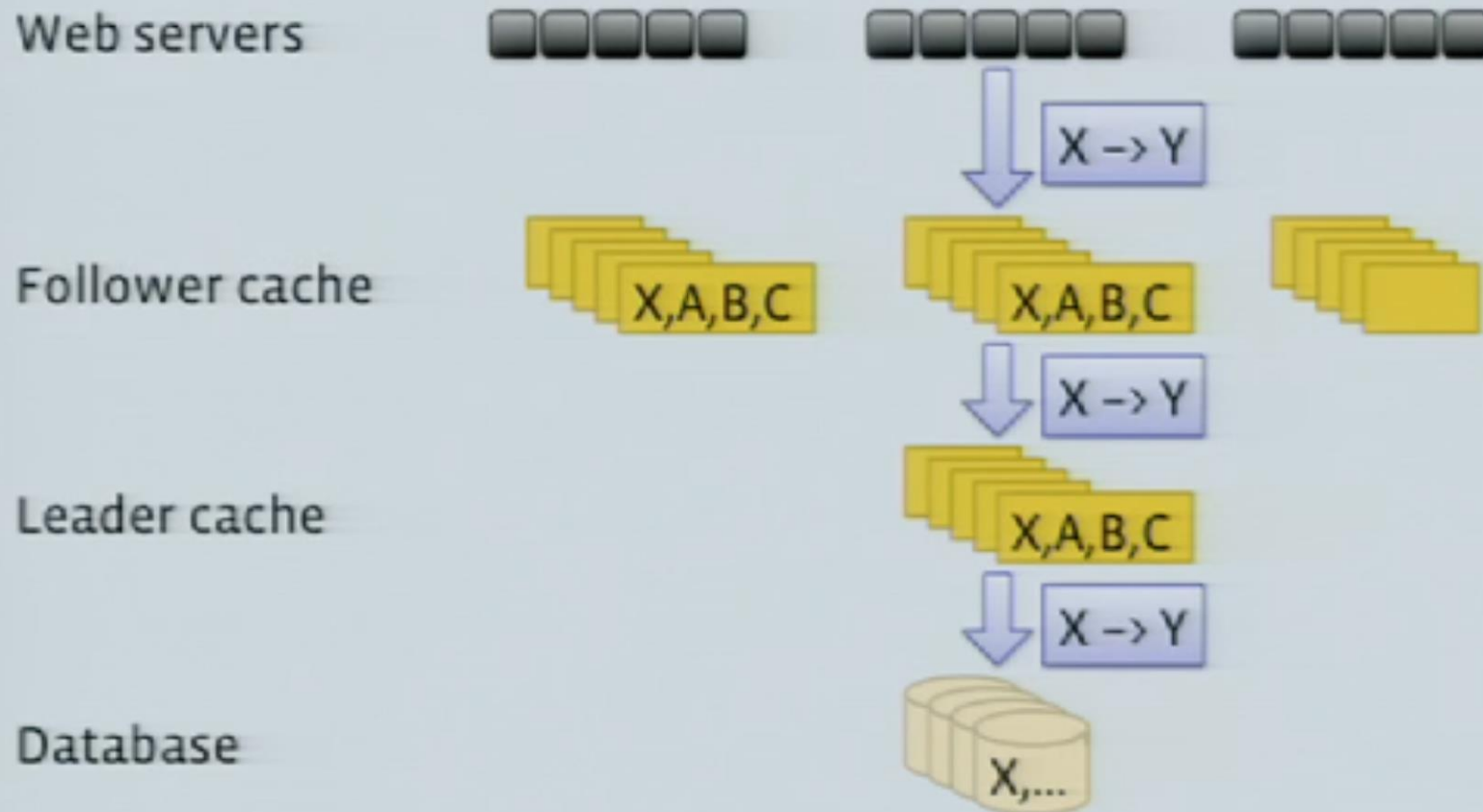


1. Complex distributed write logic to handle to ensure data correctness, consistency
2. A lot of concurrent read or write go to the MySQL server, for the same data
3. Make tradeoffs between limiting the maximum tier size and scaling the cache

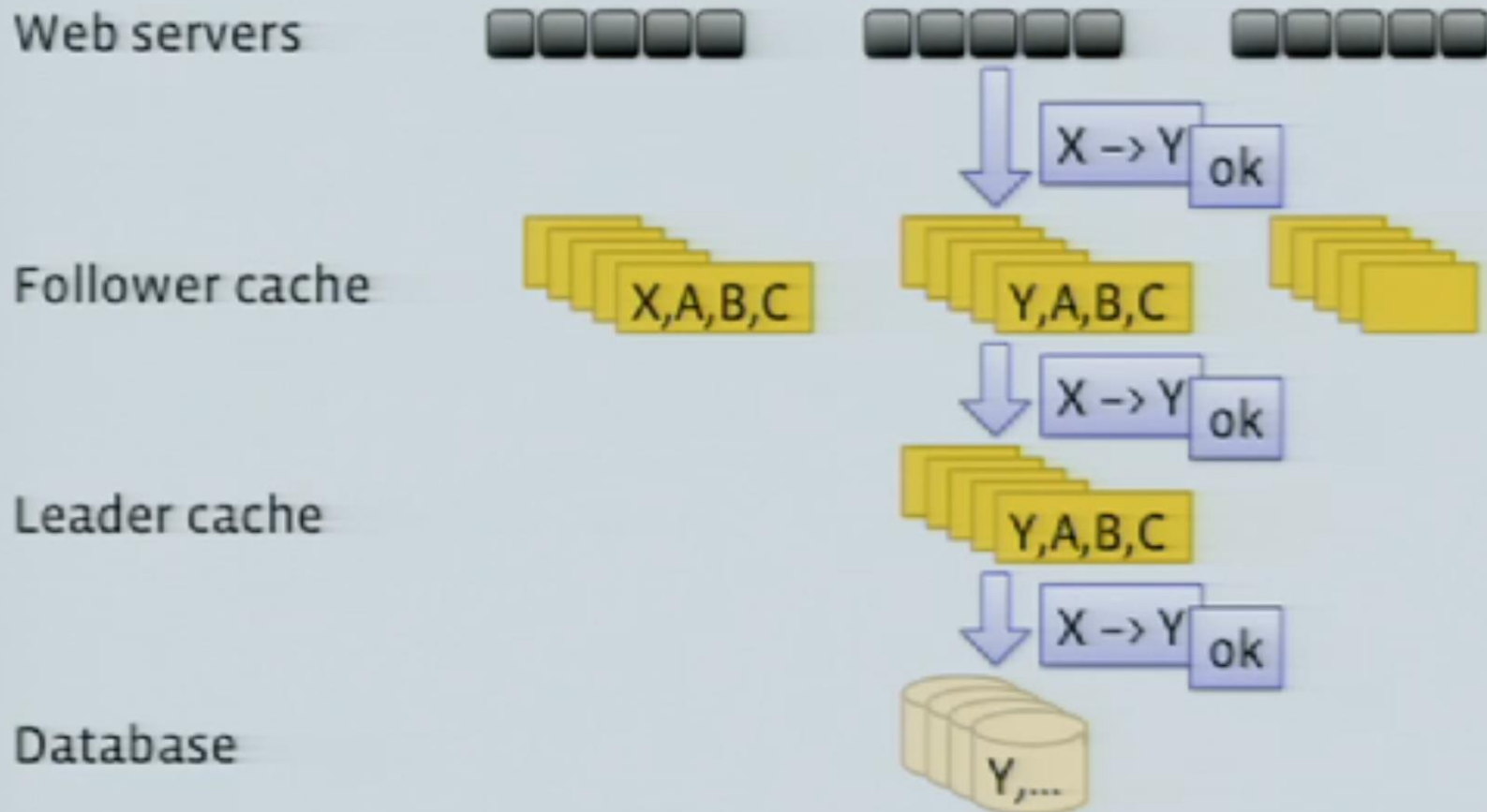
Follower and Leader Caches



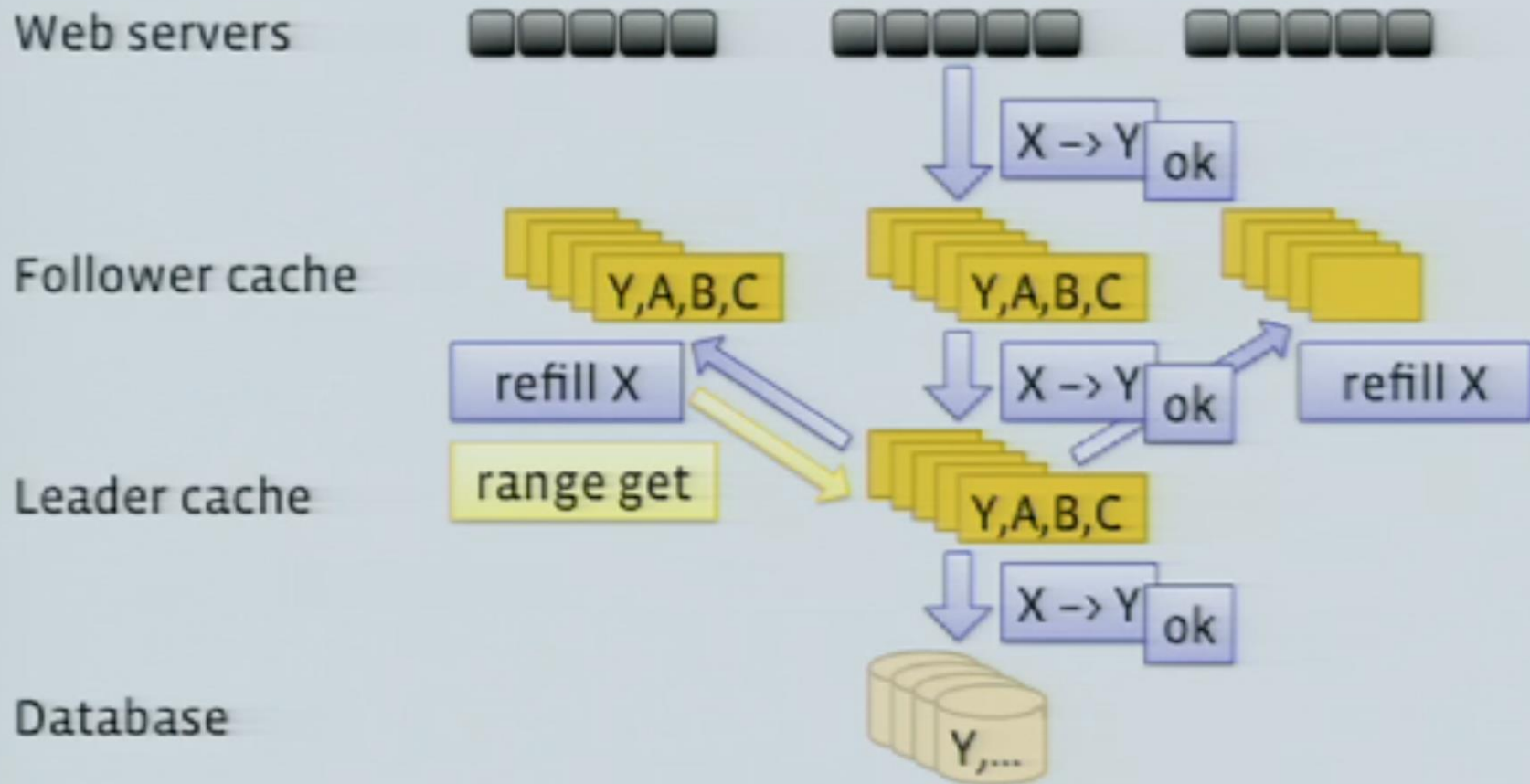
Write-through Caching – Association Lists



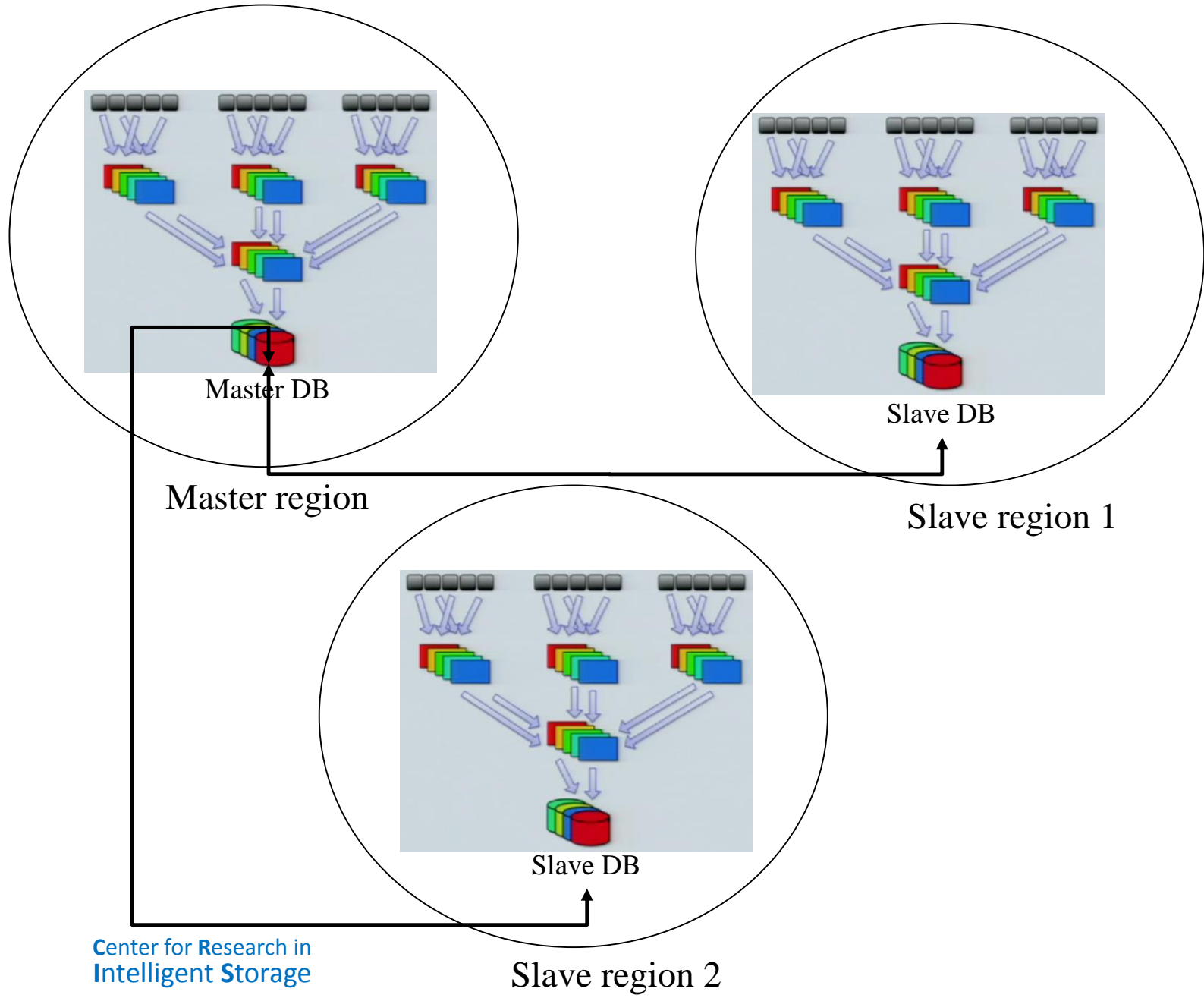
Write-through Caching – Association Lists



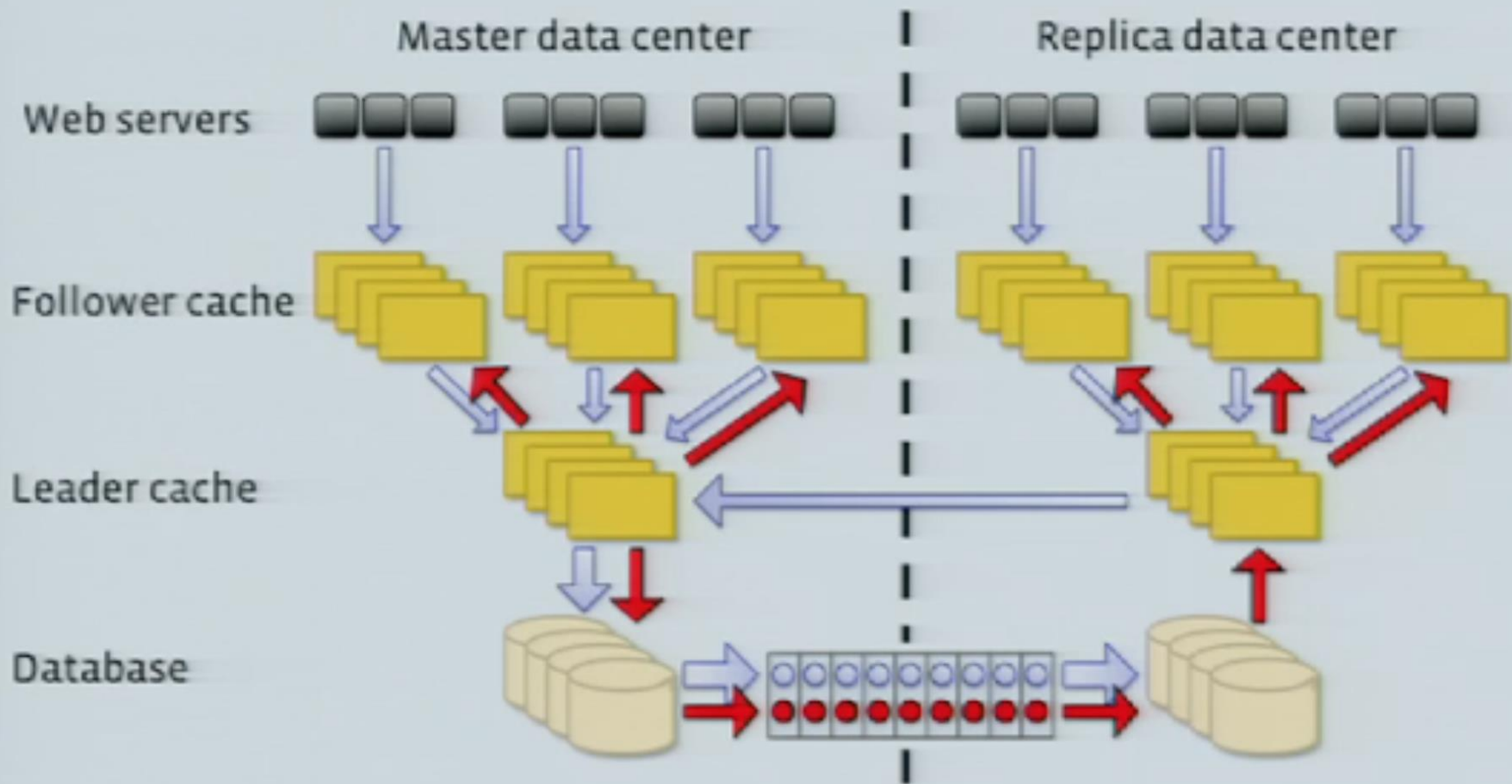
Write-through Caching – Association Lists



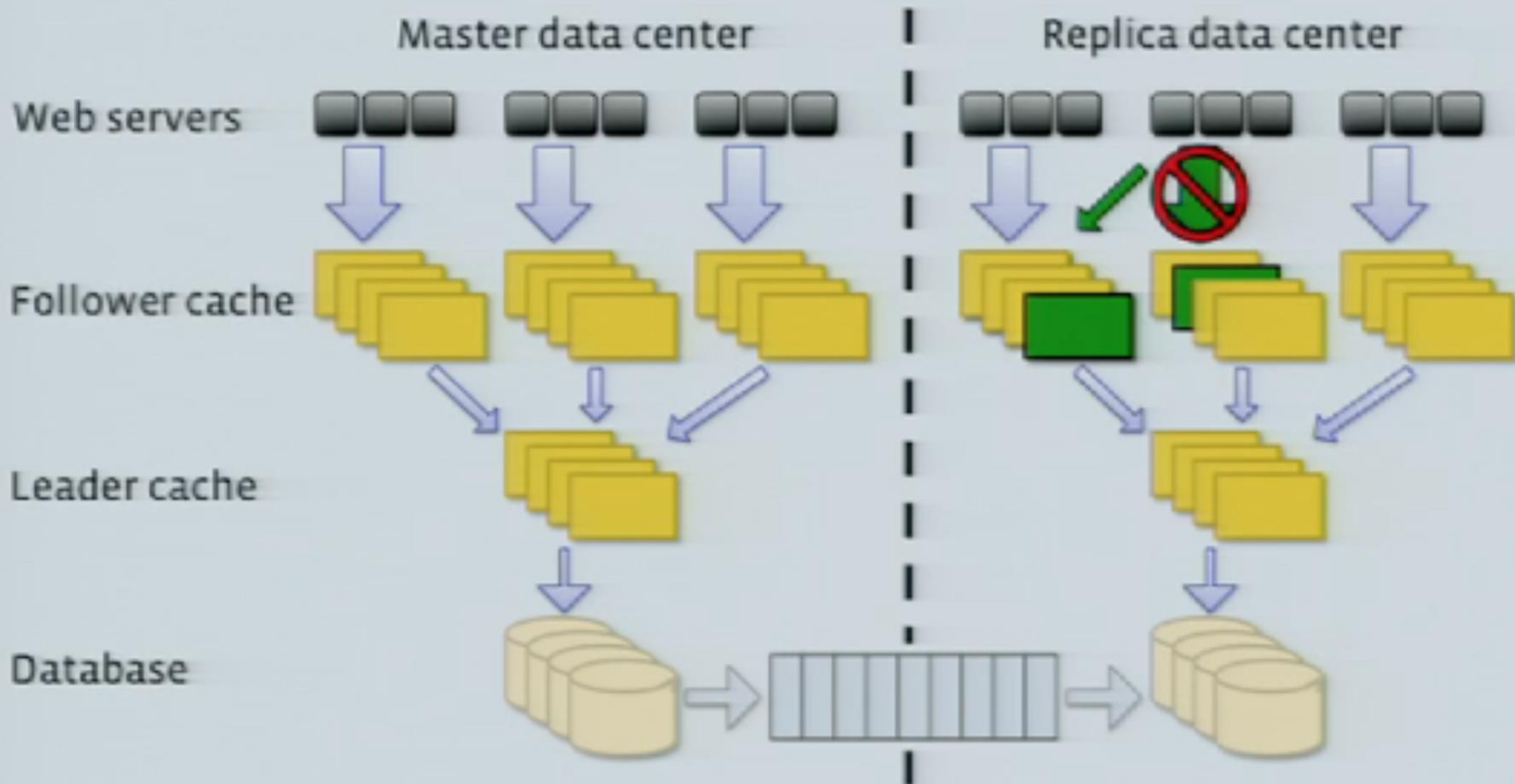
How to Scale out



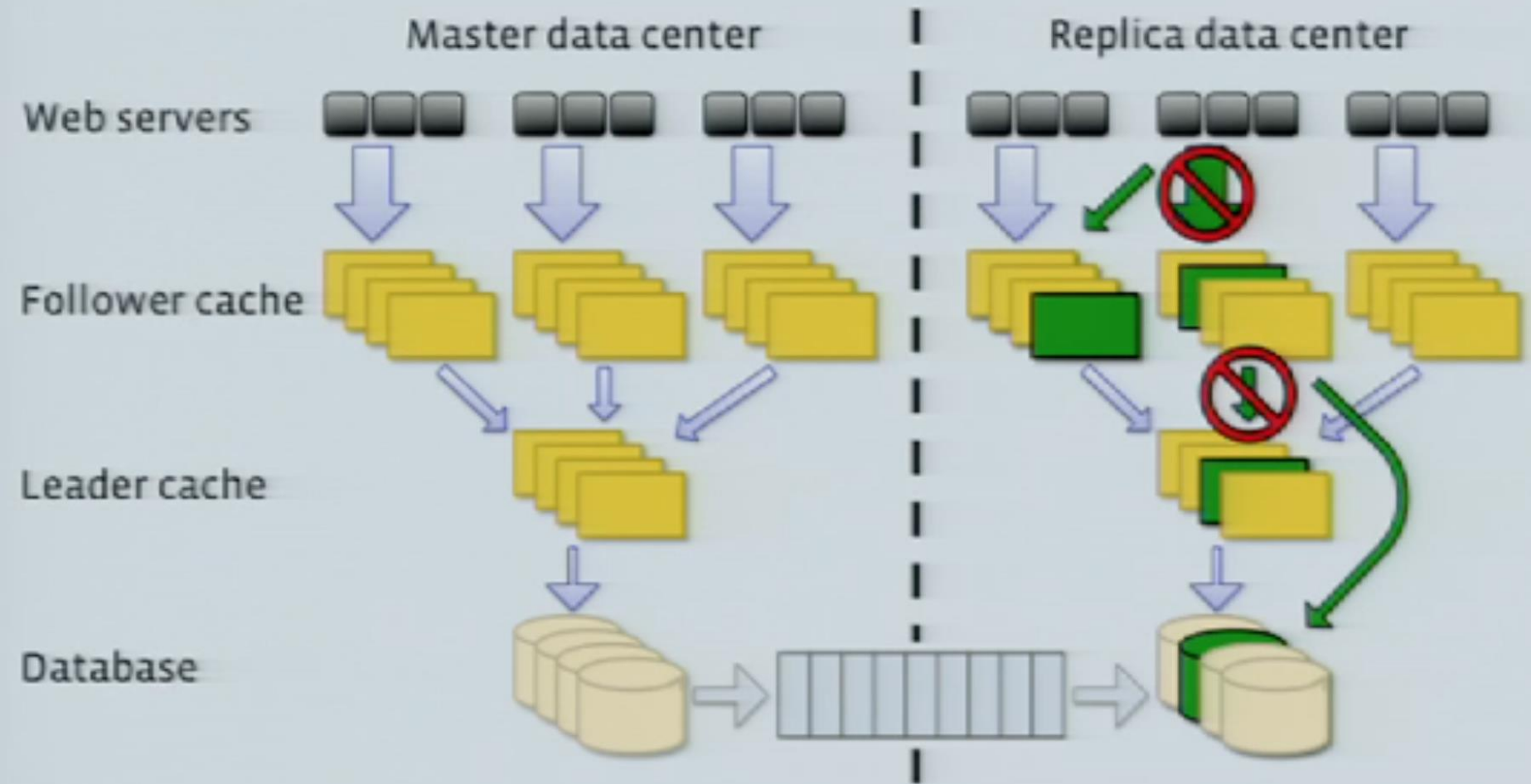
Asynchronous DB Replication



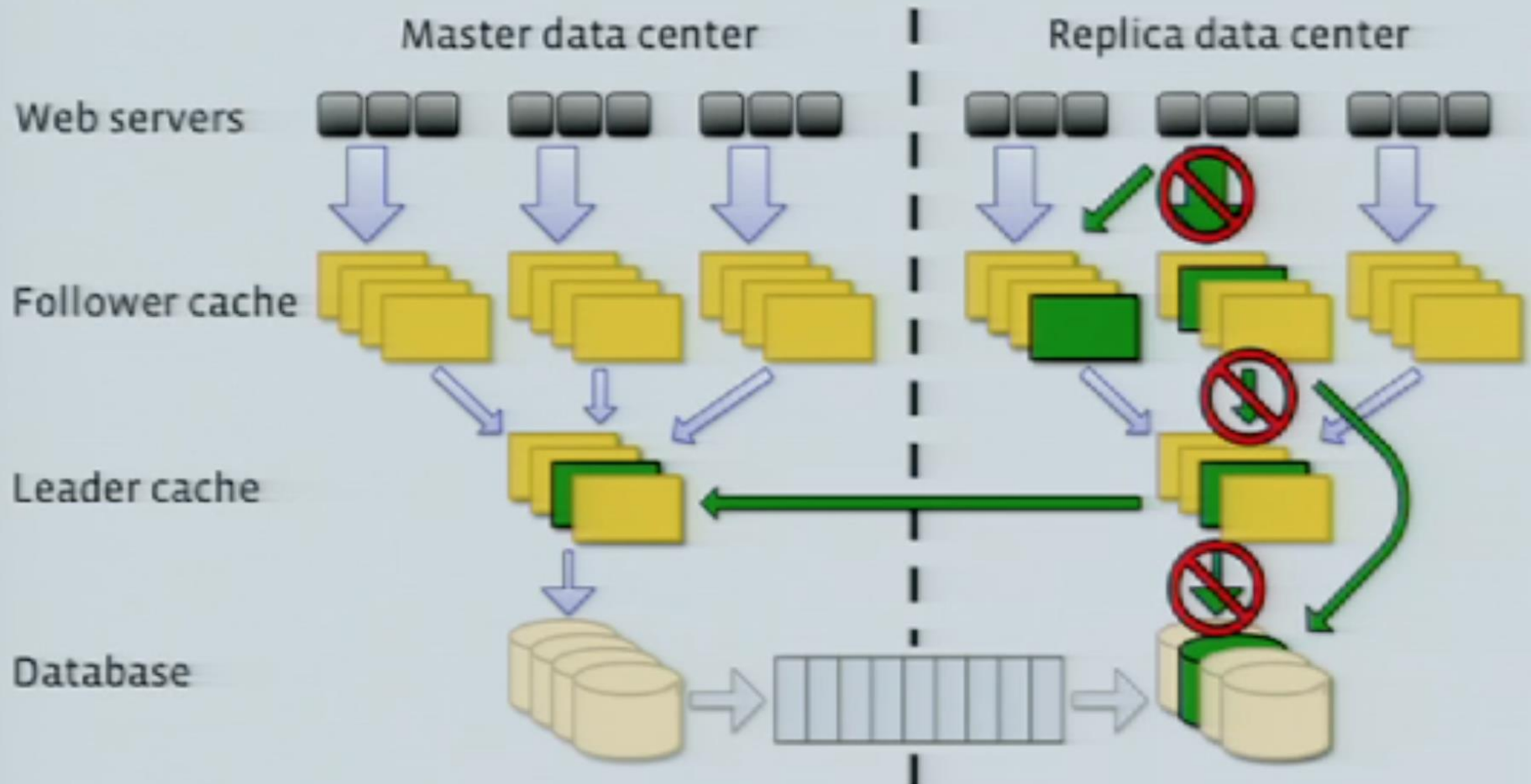
Improving Availability: Read Failover



Improving Availability: Read Failover



Improving Availability: Read Failover



MySQL Architecture

Server's functionality like connection management/authentication... is done in this layer

Connection Management/security

SQL Parsing, execution and caching...

Responsible for storage and retrieval of all available information

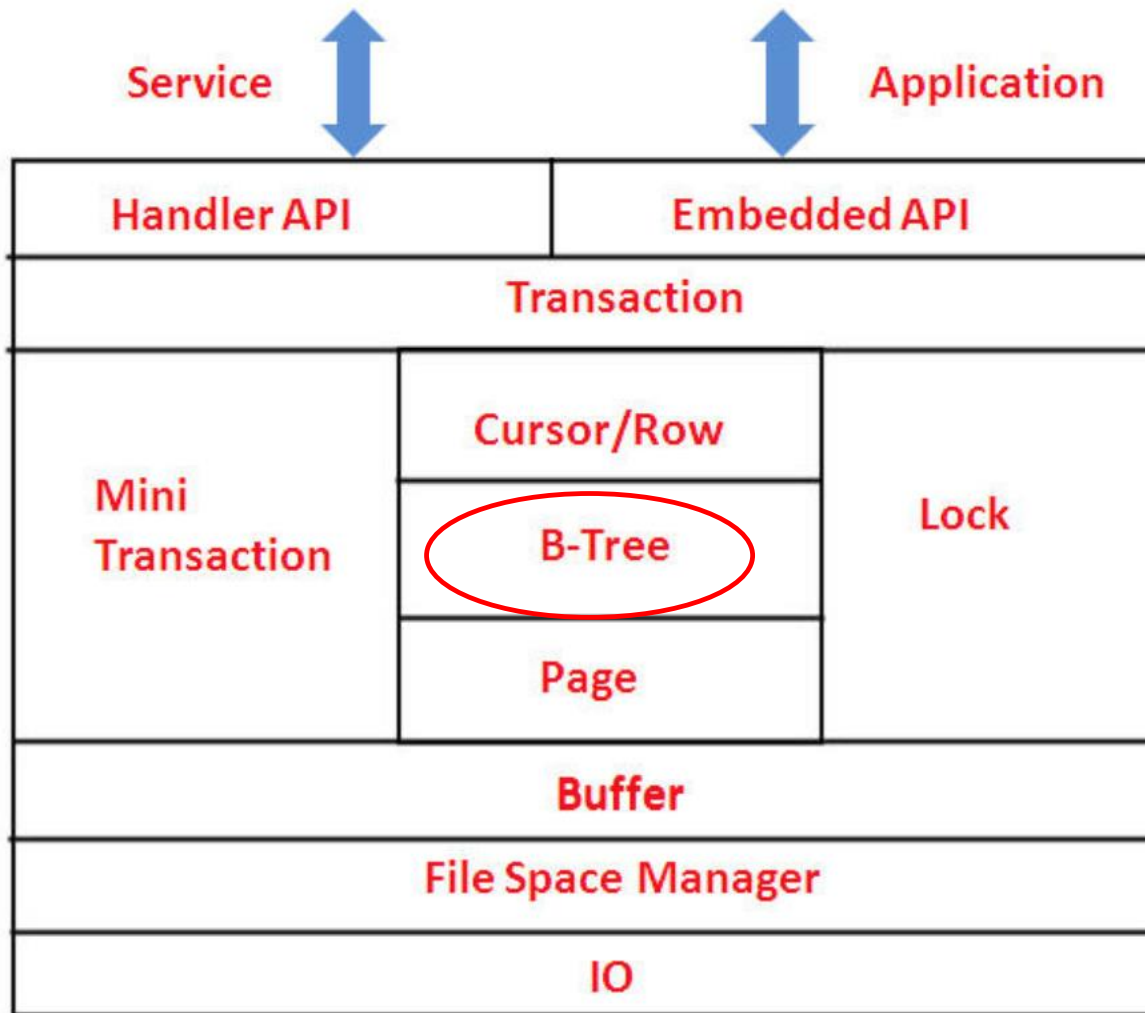
MyISAM

Innodb

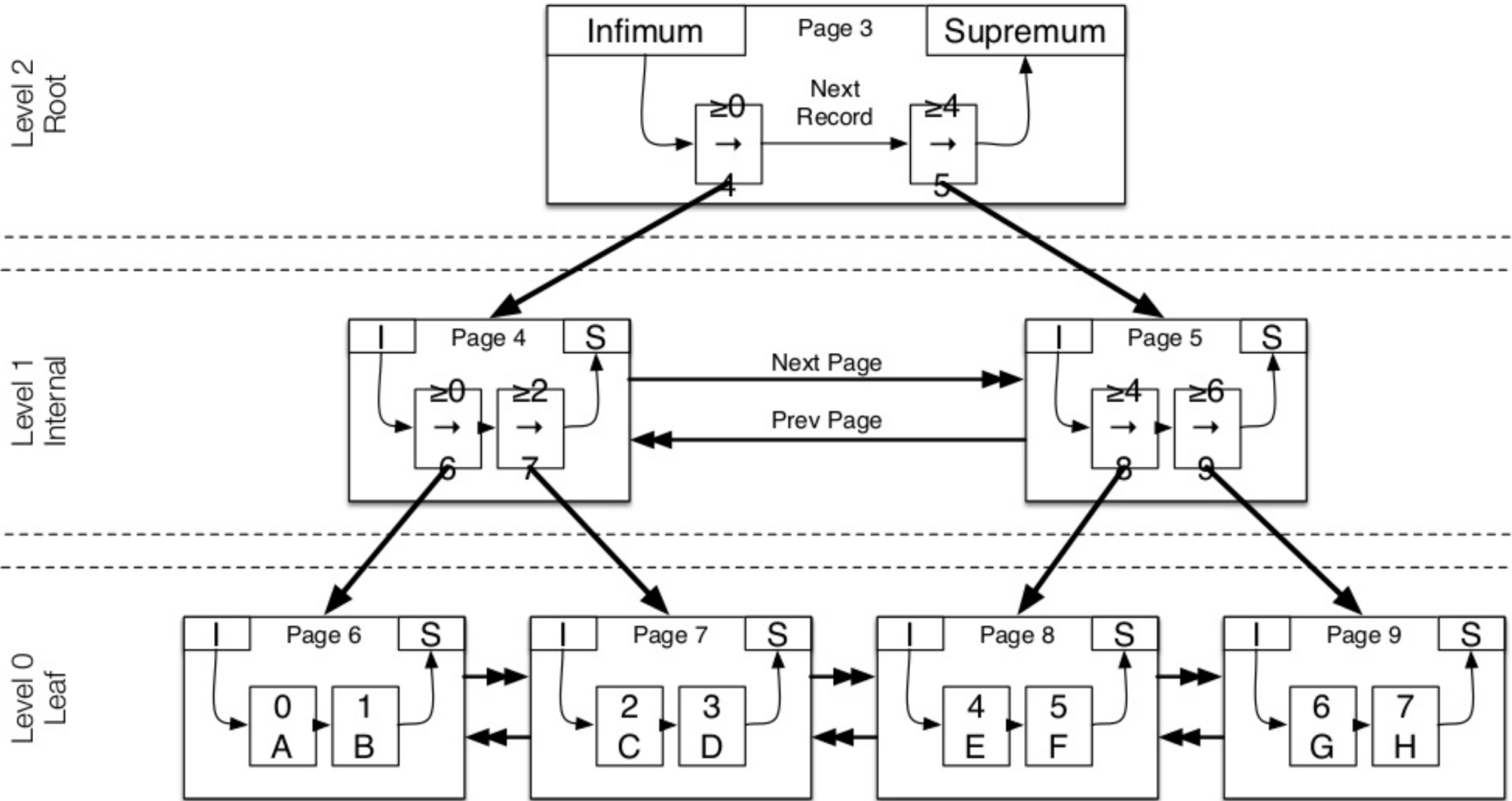
Heap
(In-Memory)

NDB
Network DB

InnoDB Architecture



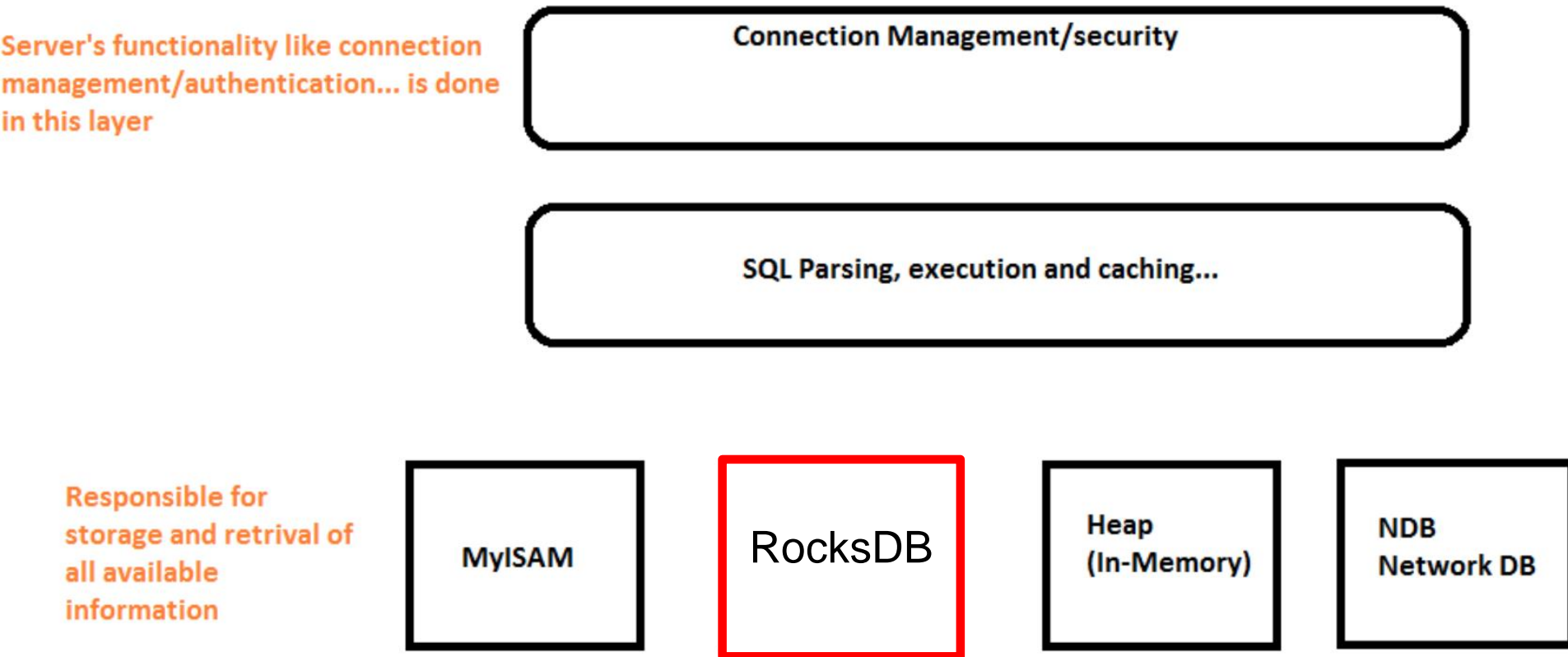
Index Structure



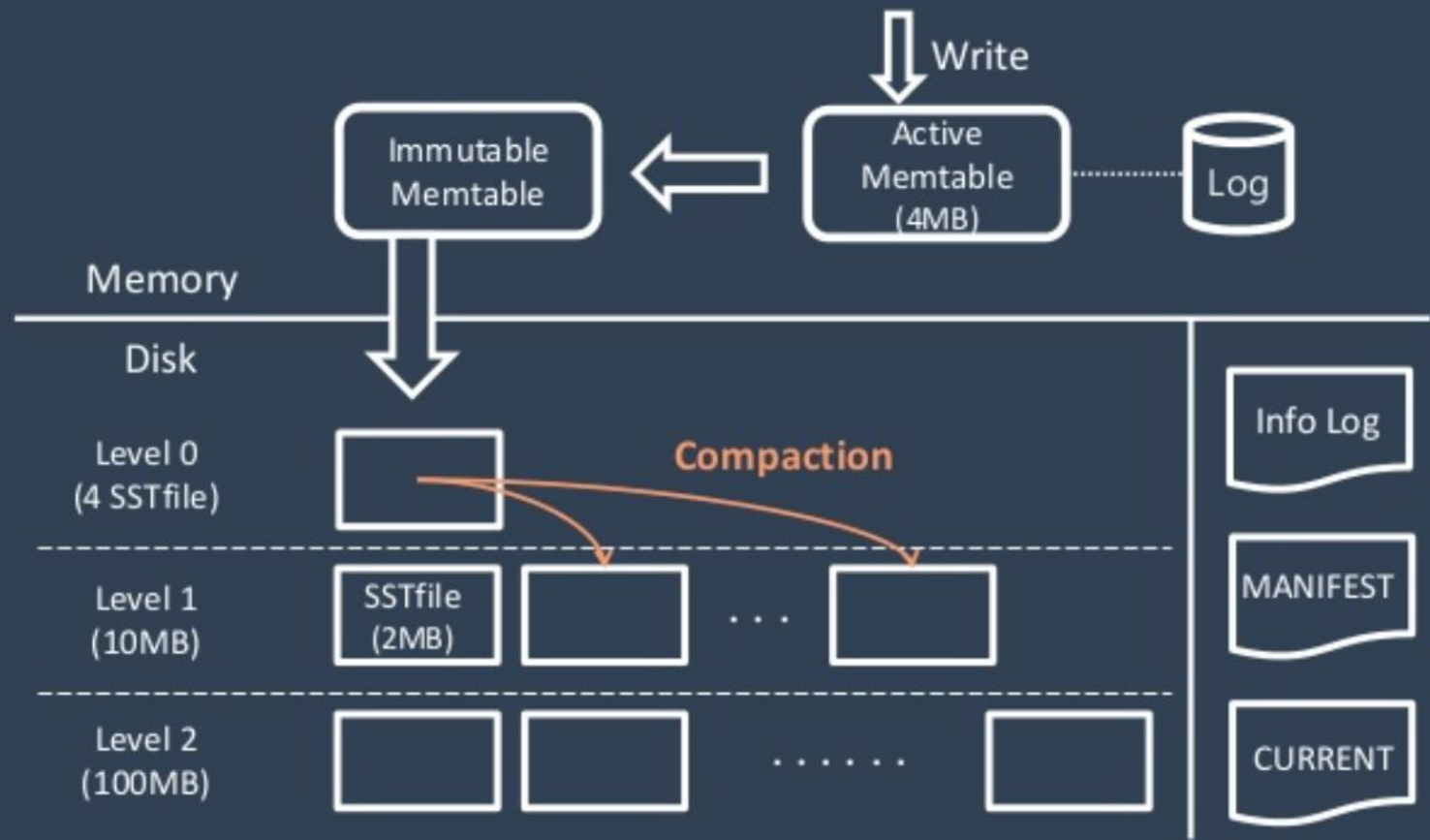
Pros?

Cons?

MySQL Architecture in FaceBook



RocksDB Architecture



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

Thanks
Q/A