

facebook

Scaling Memcache_d at Facebook

Presenter: Rajesh Nishtala (rajesh.nishtala@fb.com)

Co-authors: Hans Fugal, Steven Grimm, Marc Kwiatkowski, Herman Lee, Harry C. Li, Ryan McElroy, Mike Paleczny, Daniel Peek, Paul Saab, David Stafford, Tony Tung, Venkateshwaran Venkataramani

facebook

The Facebook logo is positioned in the bottom right corner of the slide. It consists of the word "facebook" in a white, lowercase, sans-serif font. The background of the slide features a faint, light blue network graph pattern of interconnected nodes and lines, which is more prominent on the right side.

Infrastructure Requirements for Facebook

1. Near real-time communication *social*
2. Aggregate content on-the-fly from multiple sources *skewed*
3. Be able to access and update very popular shared content
4. Scale to process millions of user requests per second

Design Requirements

Support a very heavy read load

- Over 1 billion reads / second
- Insulate backend services from high read rates ACID,

Geographically Distributed

Support a constantly evolving product

- System must be flexible enough to support a variety of use cases
- Support rapid deployment of new features

Persistence handled outside the system

- Support mechanisms to refill after updates

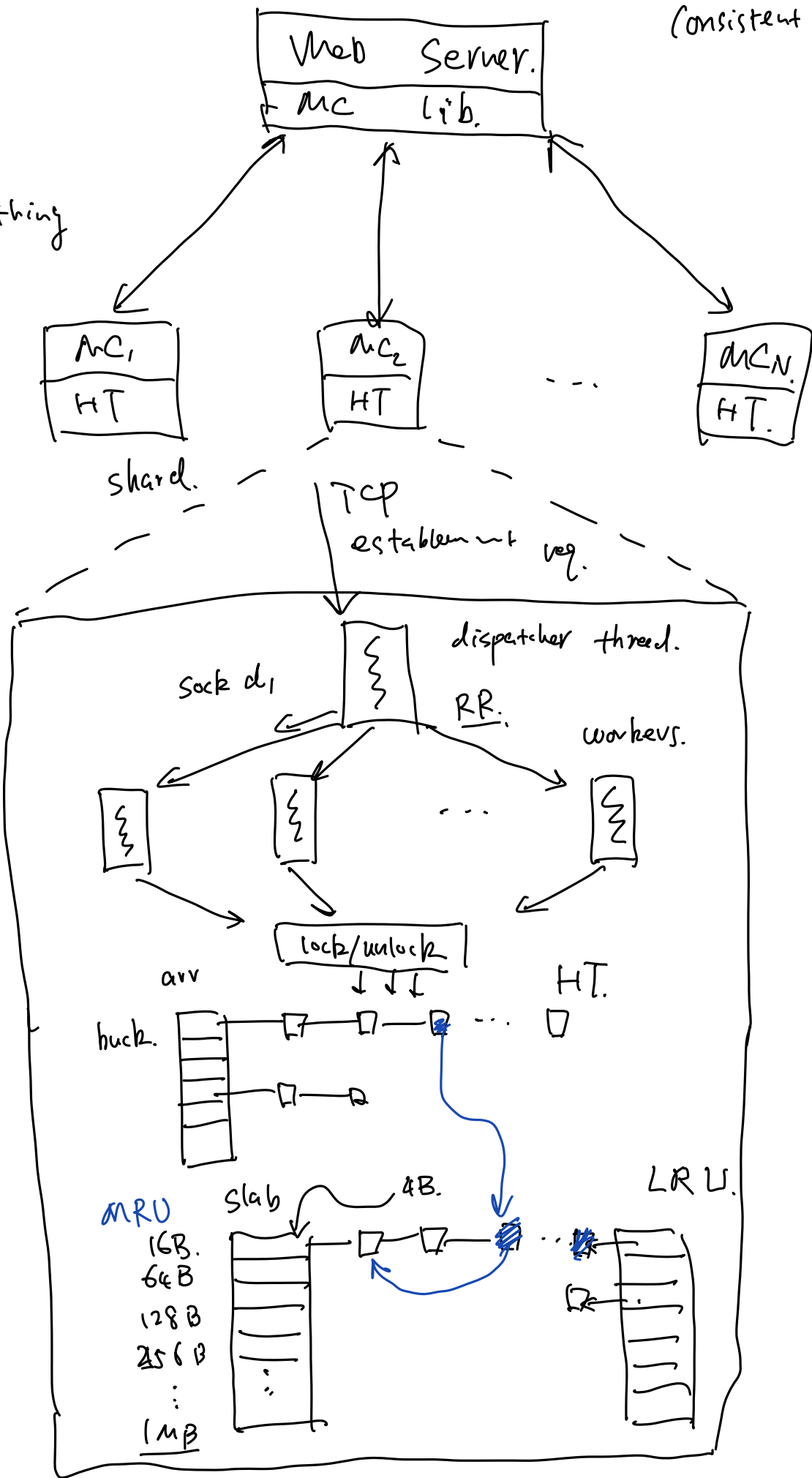
memcached

LRU.

- Basic building block for a distributed key-value store for Facebook
 - Trillions of items
 - Billions of requests / second
- Network attached in-memory hash table
 - Supports LRU based eviction

Consistent hashing.

Shared-nothing
Arch.



UserID → names

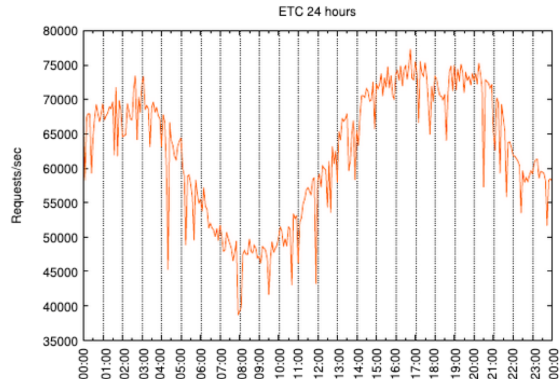
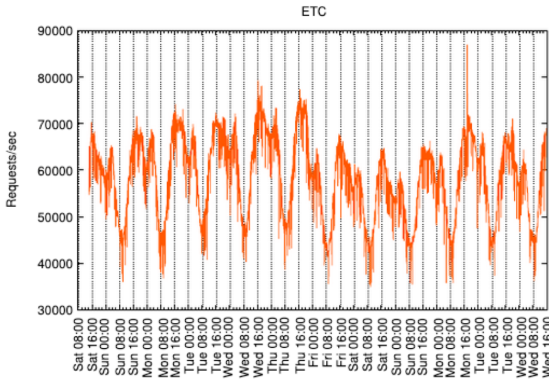
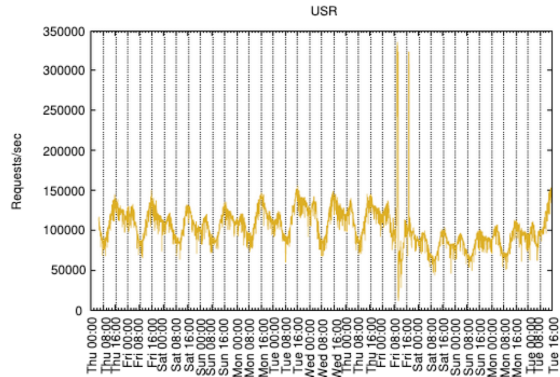
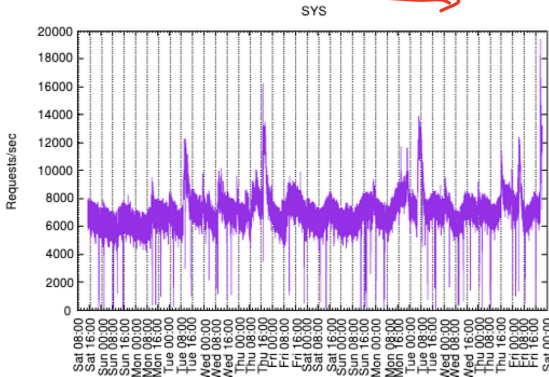
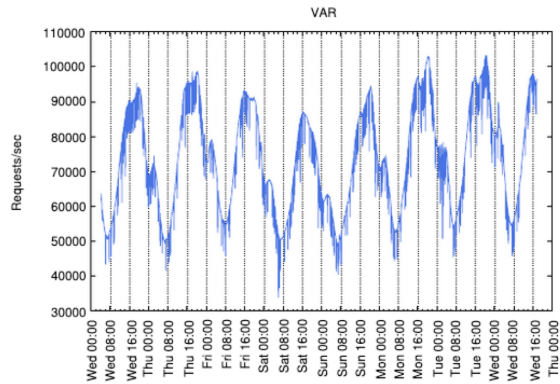
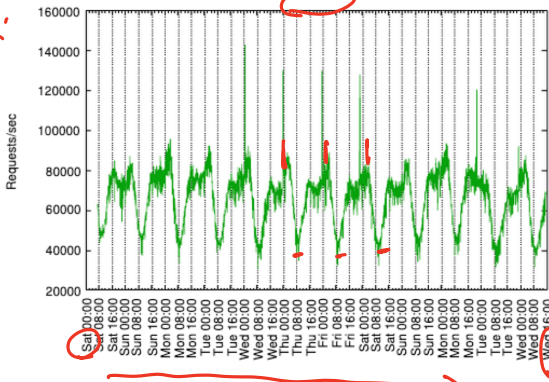
postID → text
news feed

Table 1: Memcached pools sampled (in one cluster). These pools do not match their UNIX namesakes, but are used for illustrative purposes here instead of their internal names.

Pool	Size	Description
<u>USR</u>	<u>few</u>	user-account status information
<u>APP</u>	dozens	object metadata of one application
ETC	hundreds	nonspecific, general-purpose
VAR	dozens	server-side browser information
<u>SYS</u>	few	<u>system data on service location</u>

copies of data by DBMS

GPS.



Roadmap

0. Pre-memcached. Era.

1. Single front-end cluster

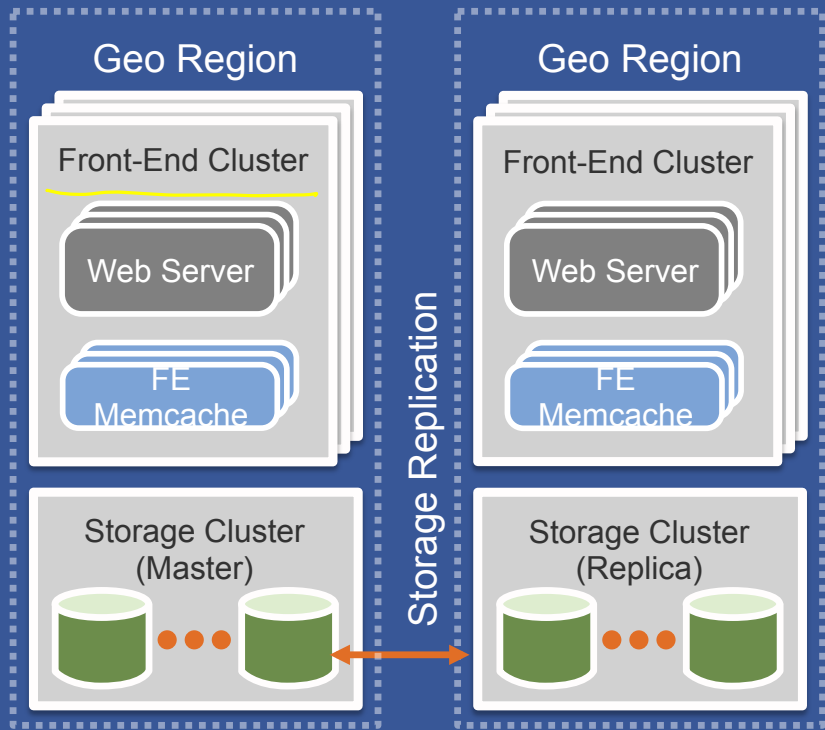
- Read heavy workload
- Wide fanout
- Handling failures

2. Multiple front-end clusters

- Controlling data replication
- Data consistency

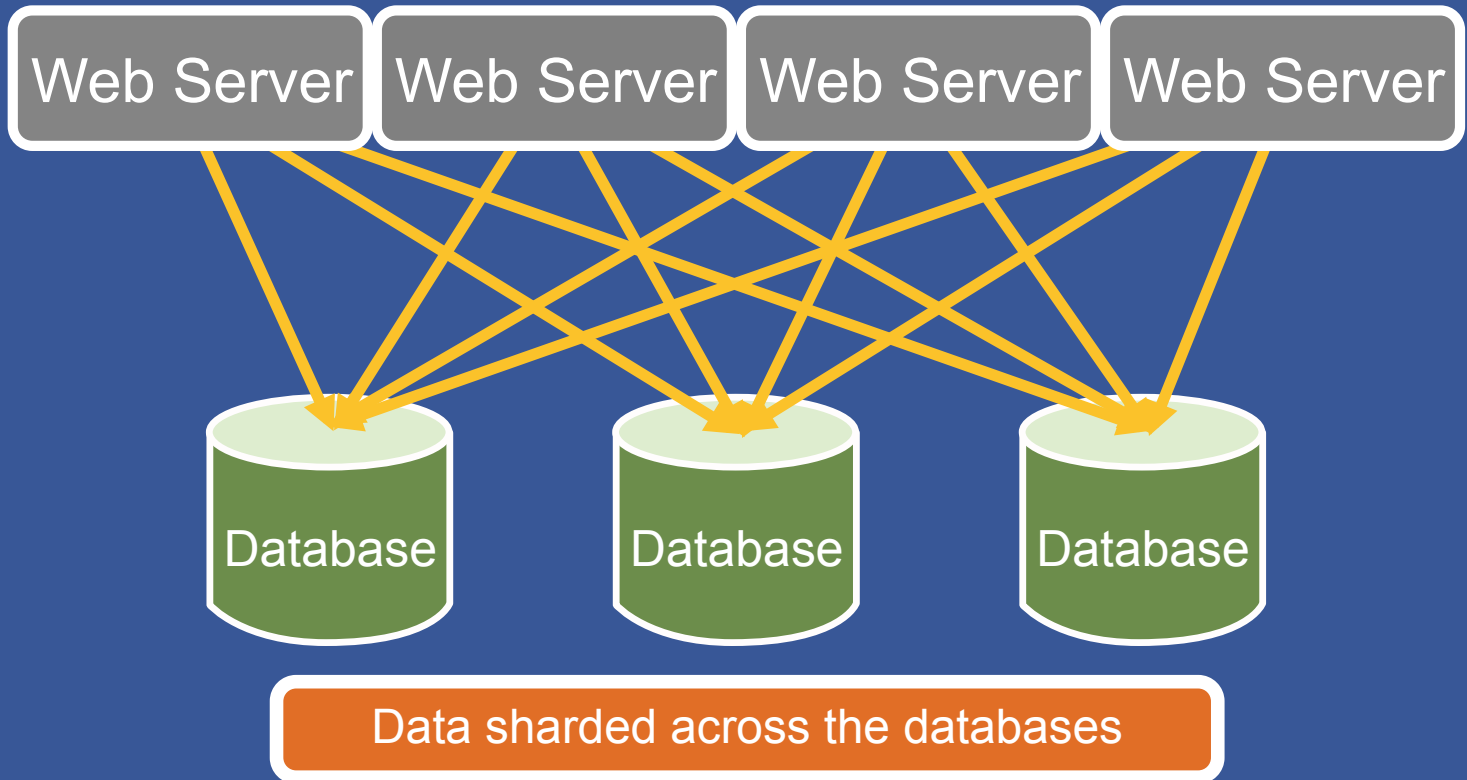
3. Multiple Regions

- Data consistency



Pre-memcache

Just a few databases are enough to support the load

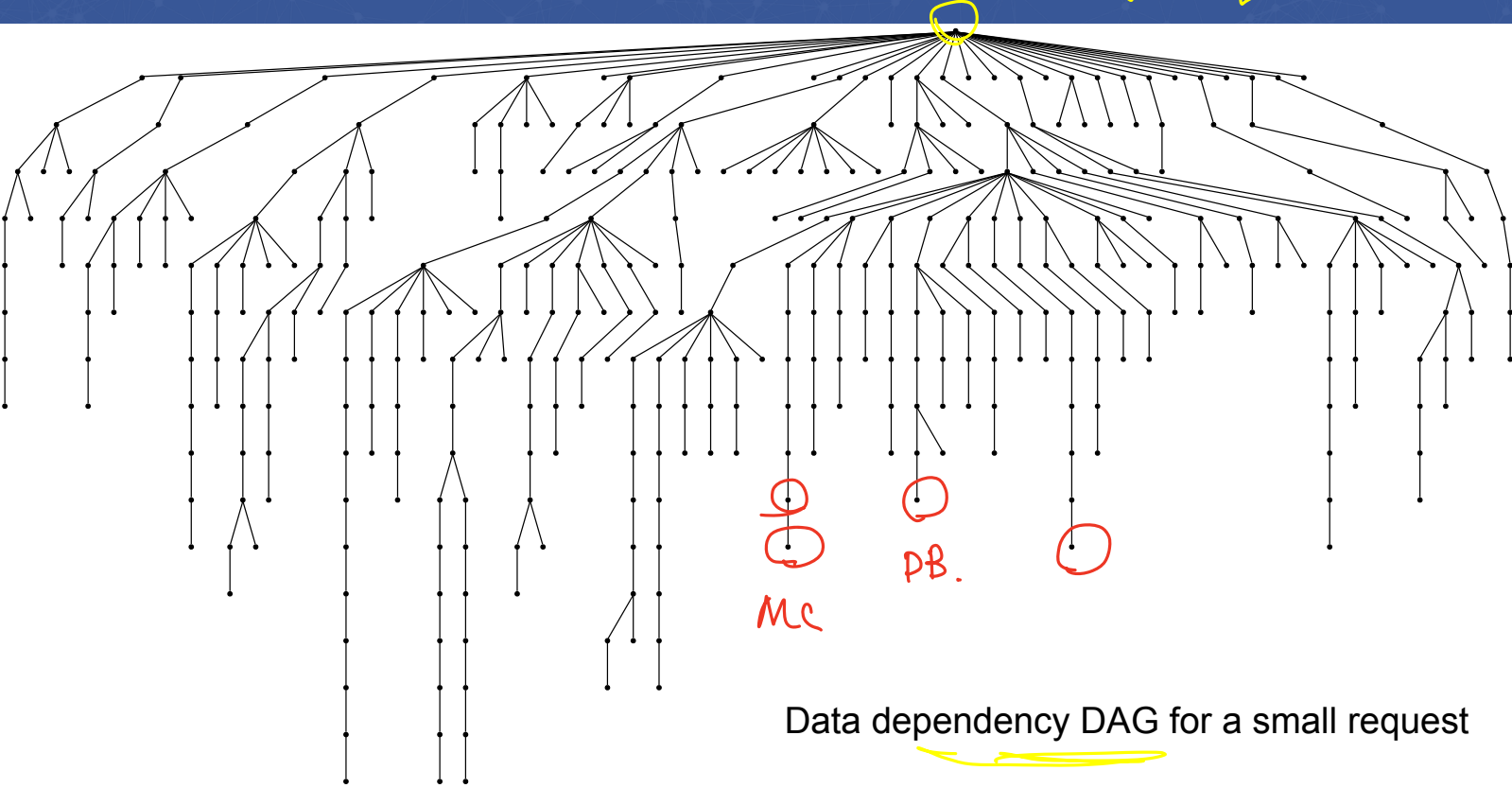


Why Separate Cache?

multi-get.

High fanout and multiple rounds of data fetching

Page Rendering Req



Data dependency DAG for a small request

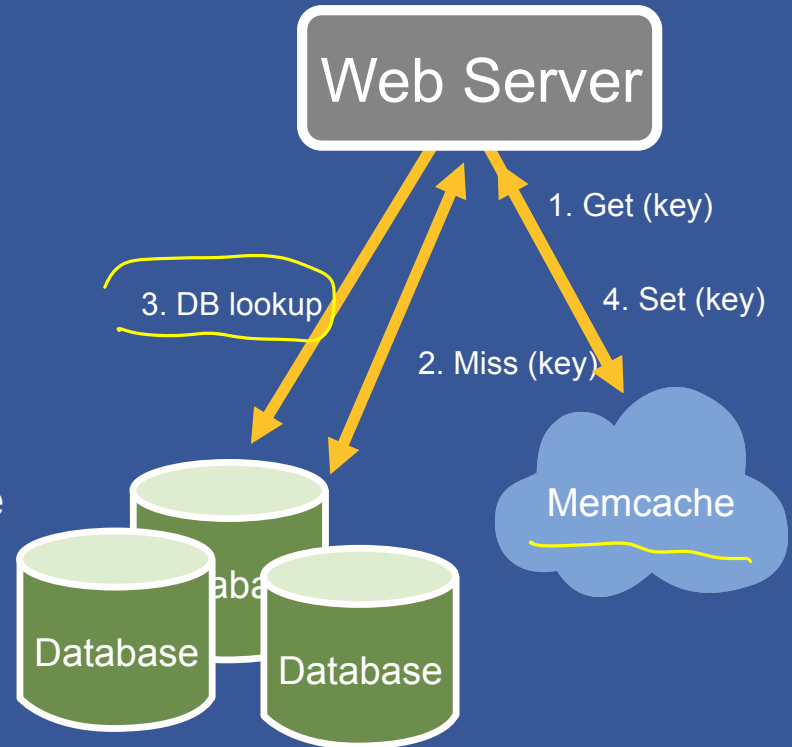
Scaling memcache in 4 “easy” steps

10s of servers & millions of operations per second

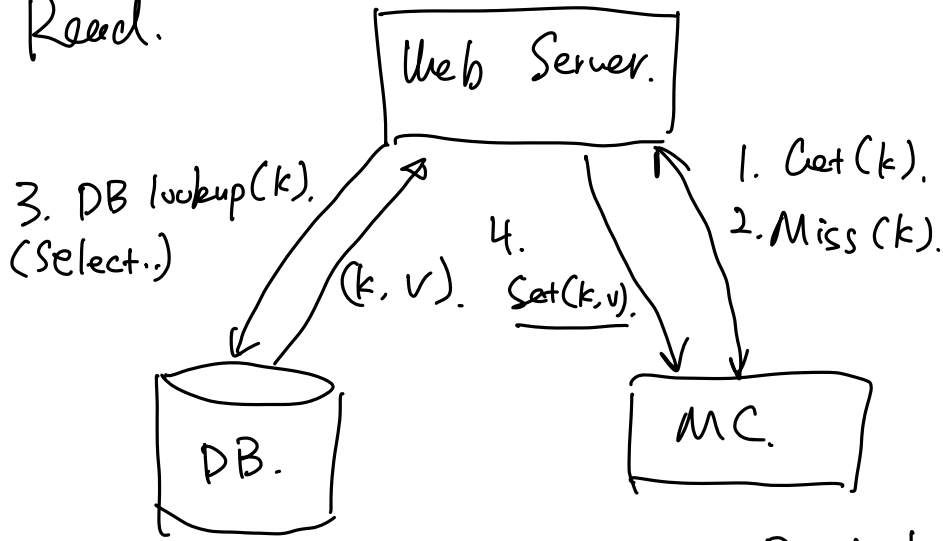
0	No memcache servers
1	A few memcache servers
2	Many memcache servers in one cluster
3	Many memcache servers in multiple clusters
4	Geographically distributed clusters

Need more read capacity

- Two orders of magnitude more reads than writes
- Solution: Deploy a few memcache hosts to handle the read capacity
- How do we store data?
 - Demand-filled look-aside cache
 - Common case is data is available in the cache



Read.

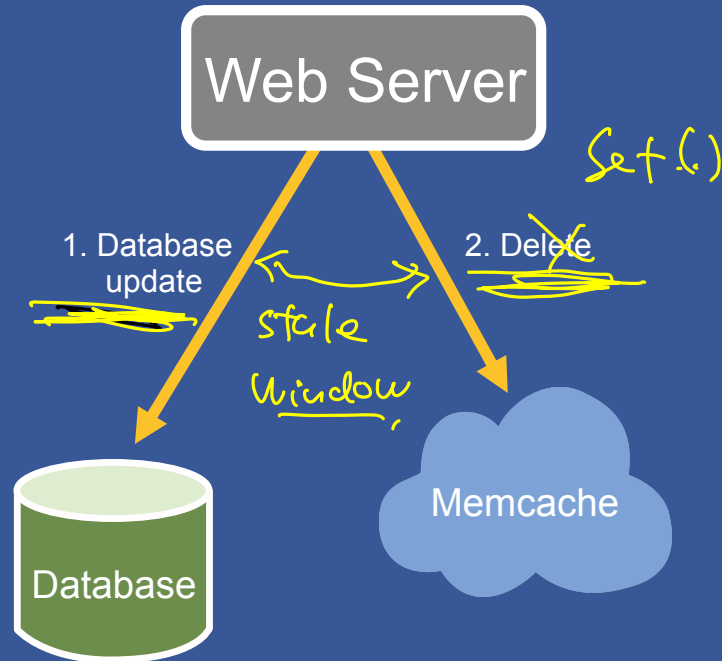


Read hits ratio $> 95\%$

misses are rare!

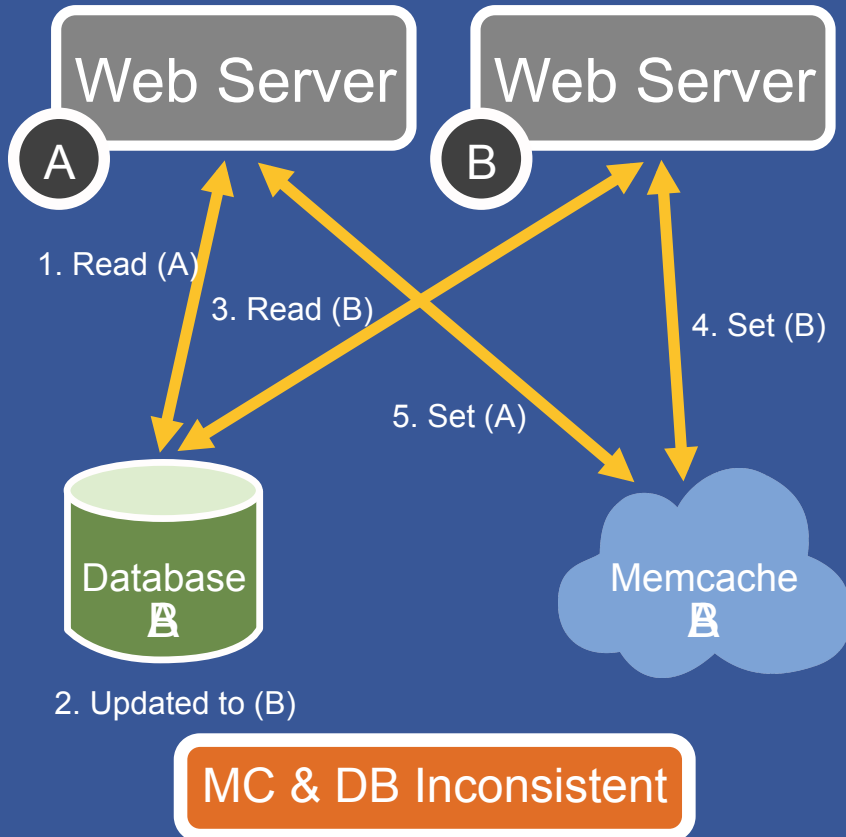
Handling updates

- Memcache needs to be invalidated after DB write
- Prefer deletes to sets
 - Idempotent
 - Demand filled
- Up to web application to specify which keys to invalidate after database update



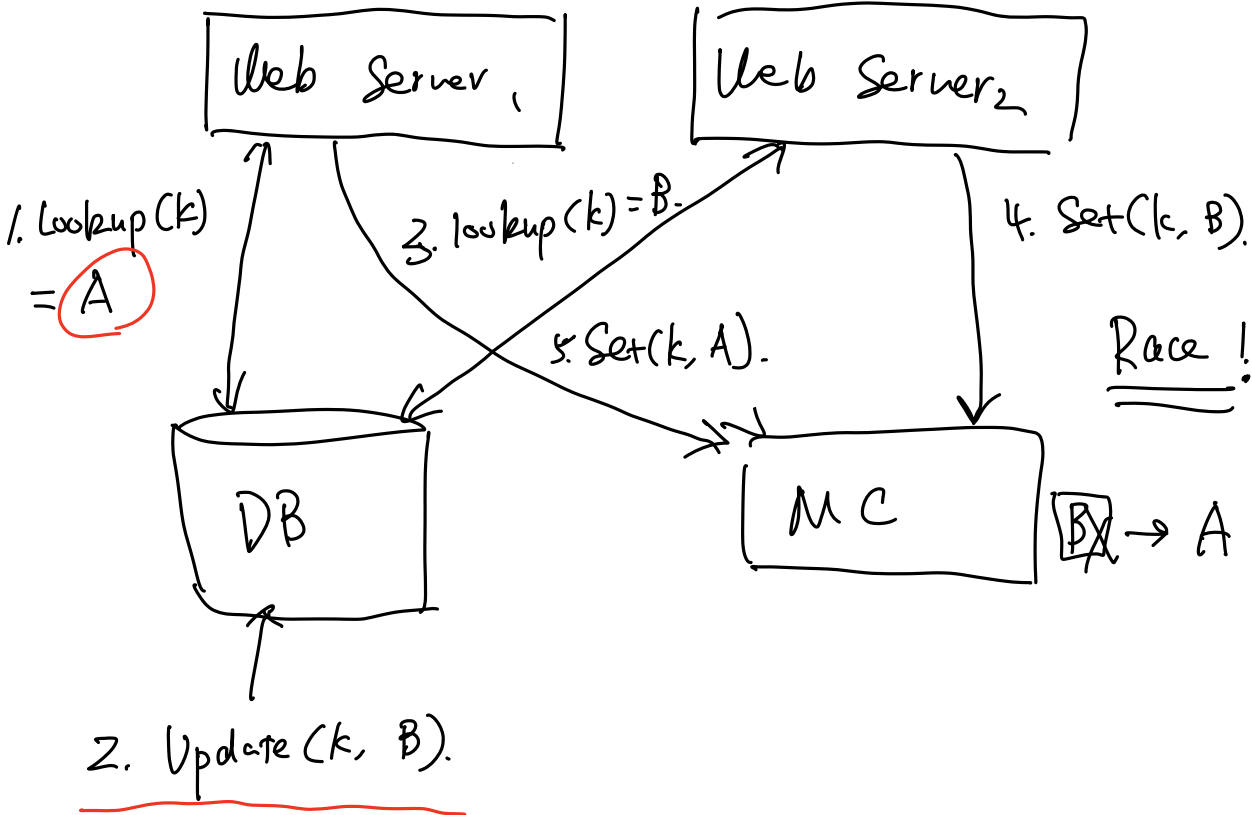
Problems with look-aside caching

Stale Sets

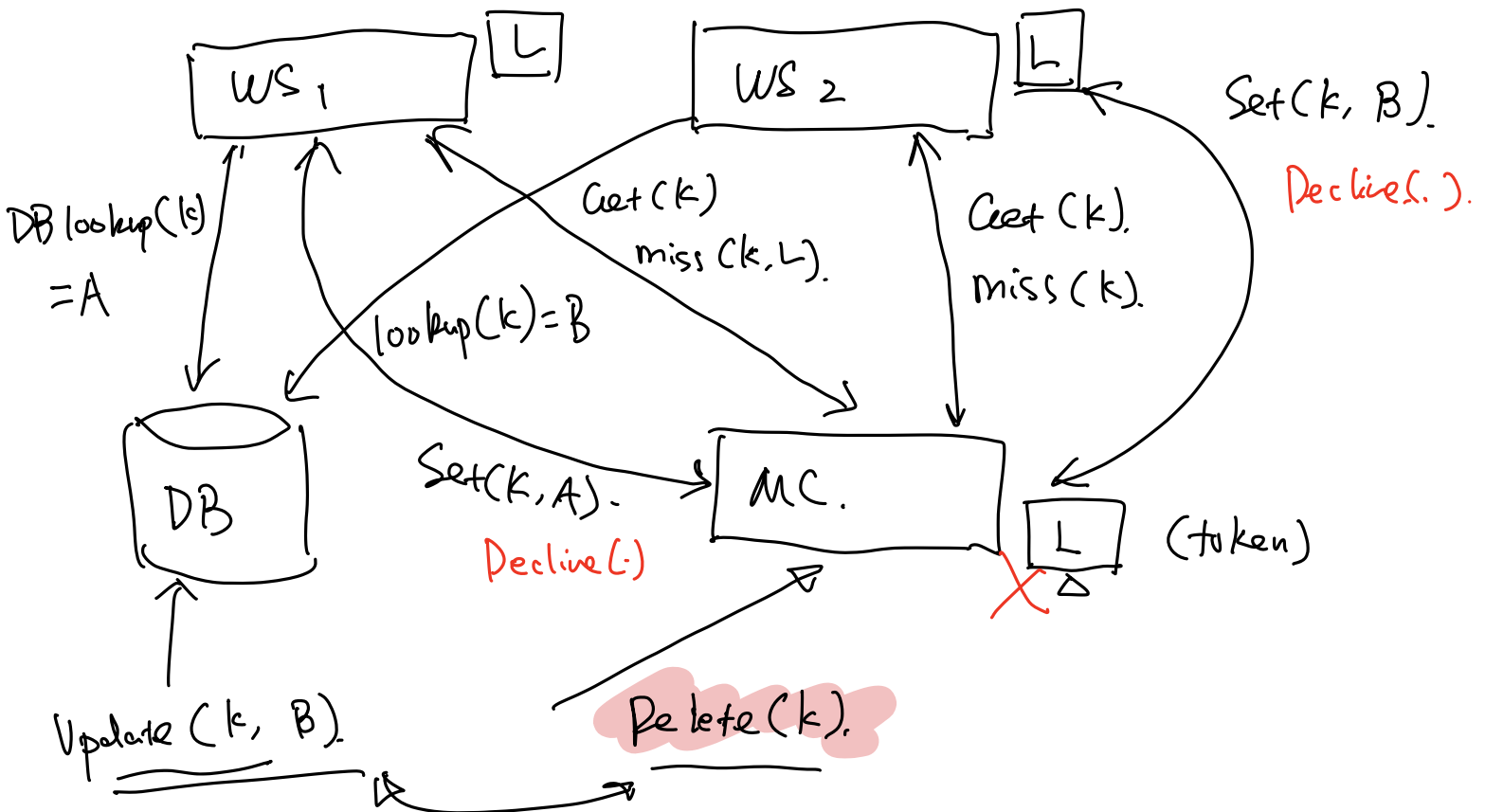


- Extend memcache protocol with “leases”
- Return and attach a lease-id with every miss
- Lease-id is invalidated inside server on a delete
- Disallow set if the lease-id is invalid at the server

Issue ① State Sets.

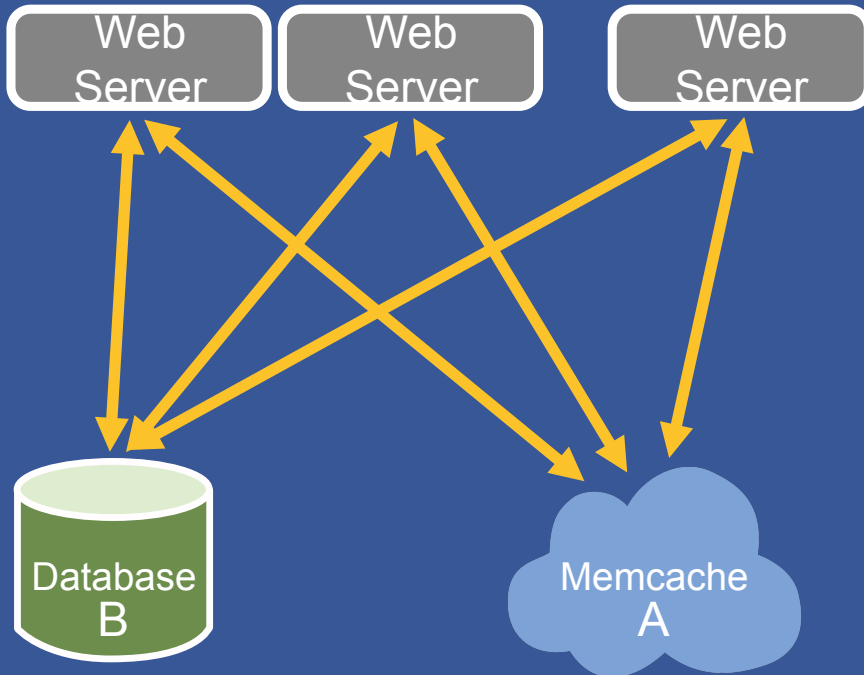


Fix: Leases. → best effort solution



Problems with look-aside caching

Thundering Herds



- Memcache server arbitrates access to database
 - Small extension to leases
- Clients given a choice of using a slightly stale value or waiting

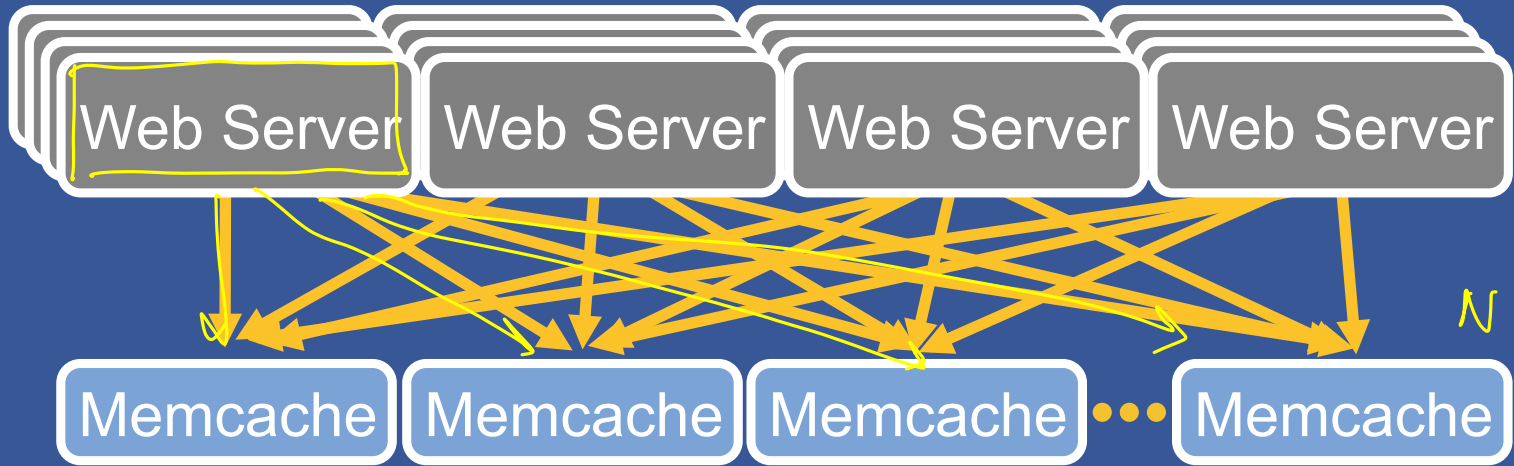
Scaling memcache in 4 “easy” steps

100s of servers & 10s of millions of operations per second

0	No memcache servers
1	A few memcache servers
2	Many memcache servers in one cluster
3	Many memcache servers in multiple clusters
4	Geographically distributed clusters

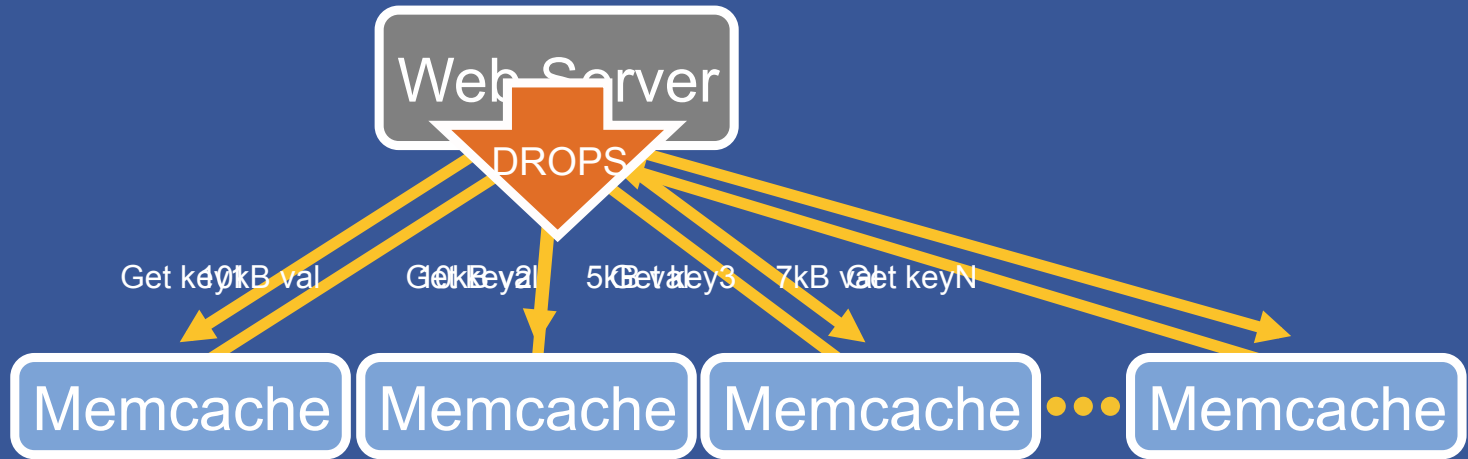
Need even more read capacity

$M \times N$
 M



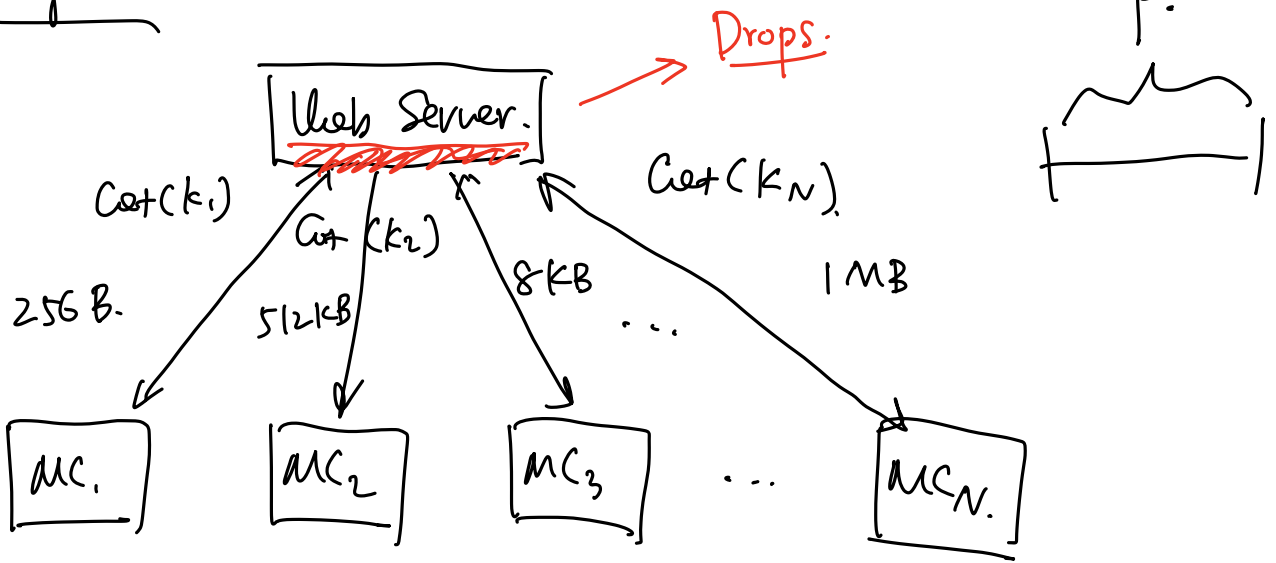
- Items are distributed across memcache servers by using consistent hashing on the key
 - Individual items are rarely accessed very frequently so over replication doesn't make sense
- All web servers talk to all memcache servers
 - Accessing 100s of memcache servers to process a user request is common

Incast congestion

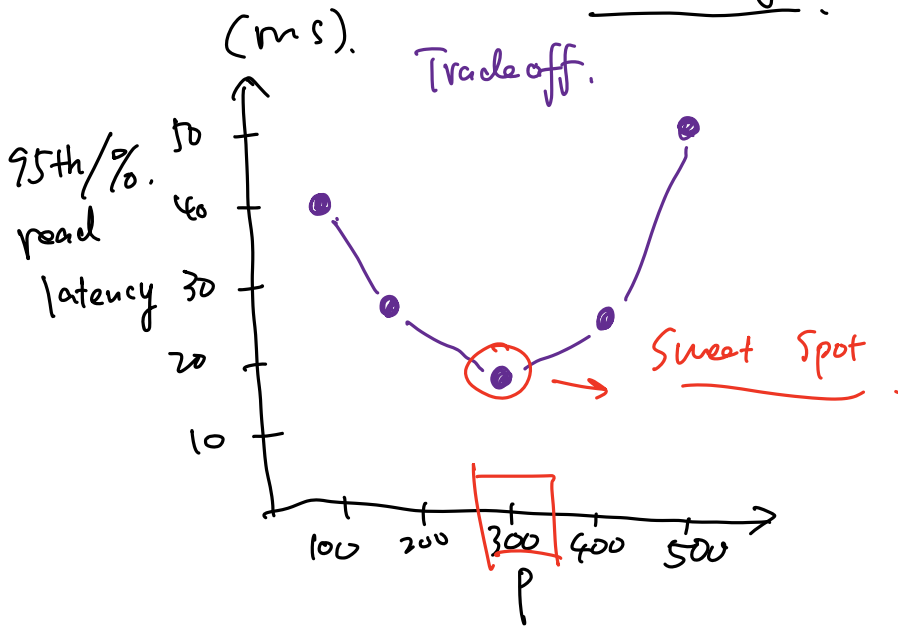


- Many simultaneous responses overwhelm shared networking resources
- Solution: Limit the number of outstanding requests with a sliding window
 - Larger windows cause result in more congestion
 - Smaller windows result in more round trips to the network

Incast Congestion:



throttling → sliding window.



p.

Performance Tuning.

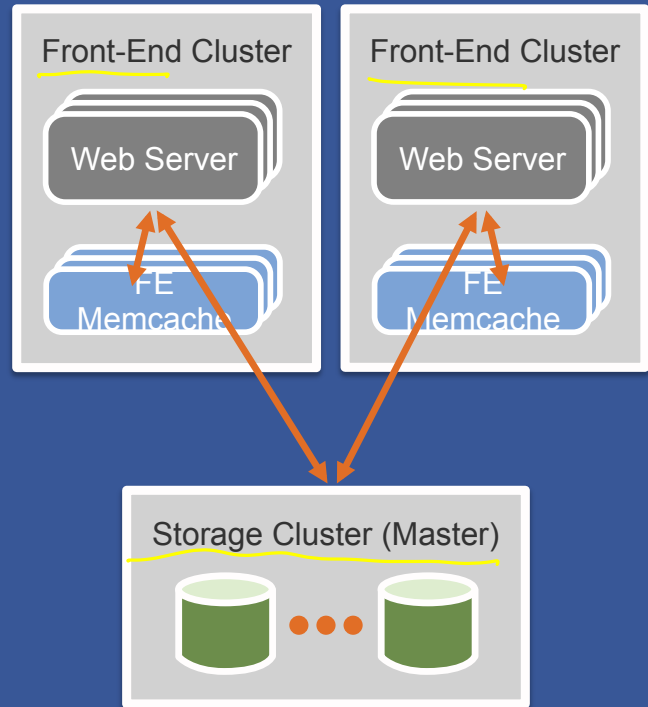
Scaling memcache in 4 “easy”

steps
1000s of servers & 100s of millions of operations per second

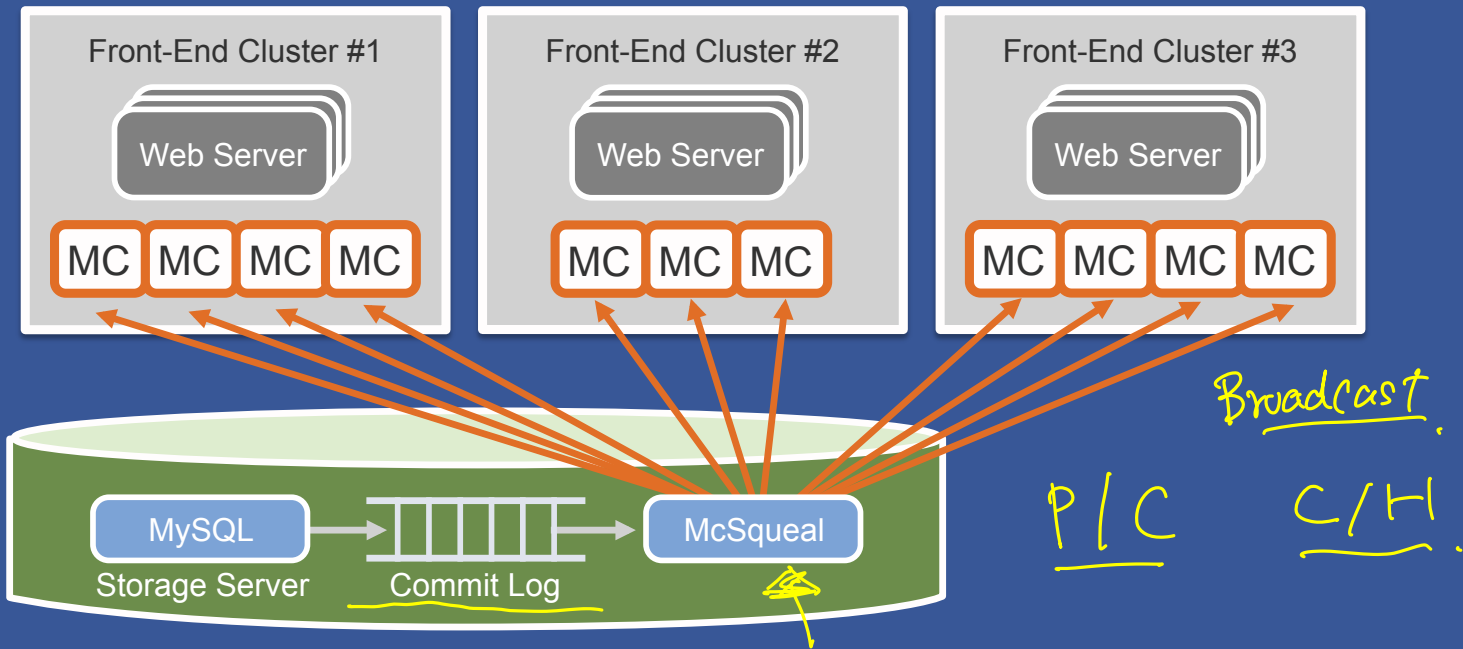
0	No memcache servers
1	A few memcache servers
2	Many memcache servers in one cluster
3	Many memcache servers in multiple clusters
4	Geographically distributed clusters

Multiple clusters

- All-to-all limits horizontal scaling
- Multiple memcache clusters front one DB installation
- Have to keep the caches consistent
- Have to manage over-replication of data



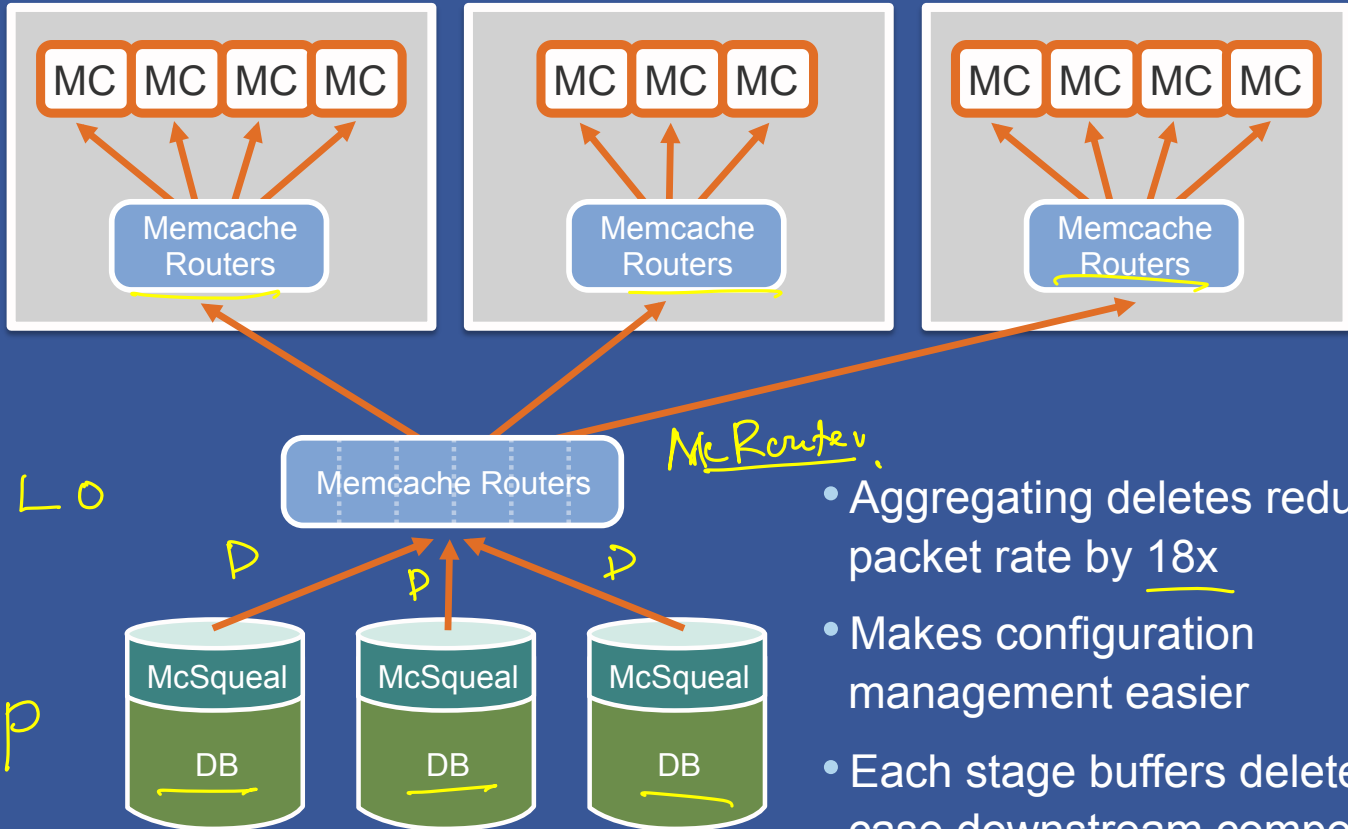
Databases invalidate caches



- Cached data must be invalidated after database updates
- Solution: Tail the mysql commit log and issue deletes based on transactions that have been committed
 - Allows caches to be resynchronized in the event of a problem

Invalidation pipeline

Too many packets



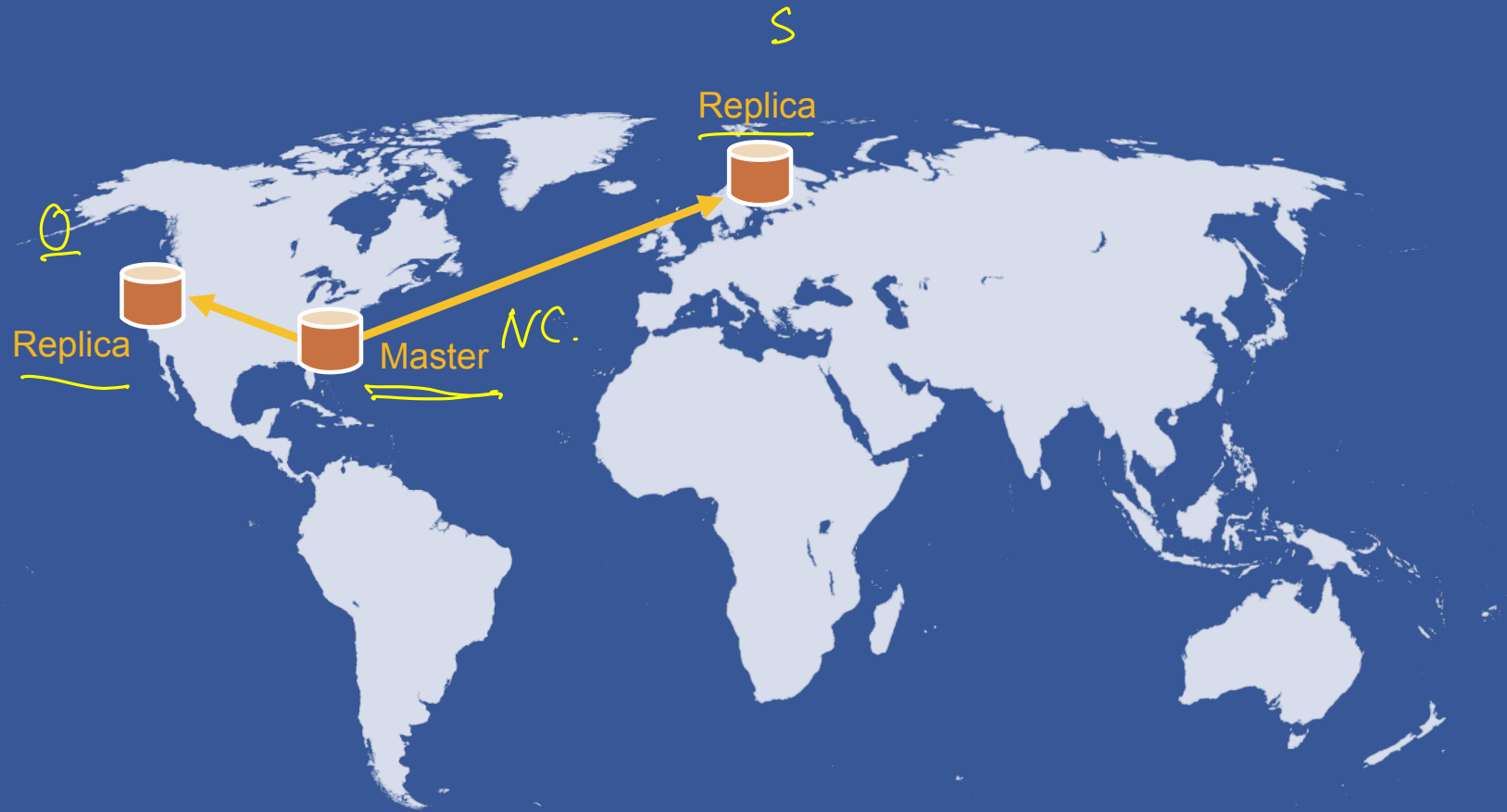
- Aggregating deletes reduces packet rate by 18x
- Makes configuration management easier
- Each stage buffers deletes in case downstream component is down

Scaling memcache in 4 “easy” steps

1000s of servers & > 1 billion operations per second

0	No memcache servers
1	A few memcache servers
2	Many memcache servers in one cluster
3	Many memcache servers in multiple clusters
4	Geographically distributed clusters

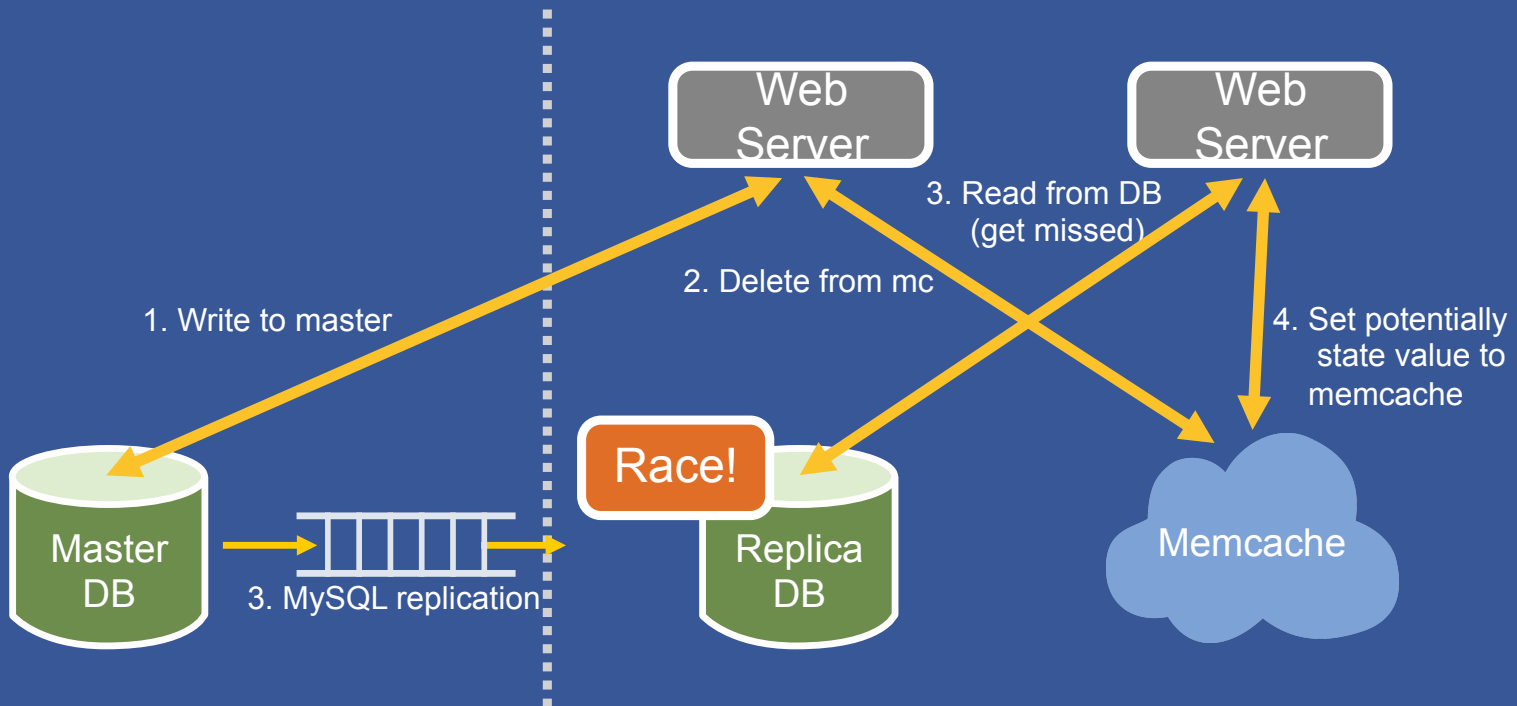
Geographically distributed clusters



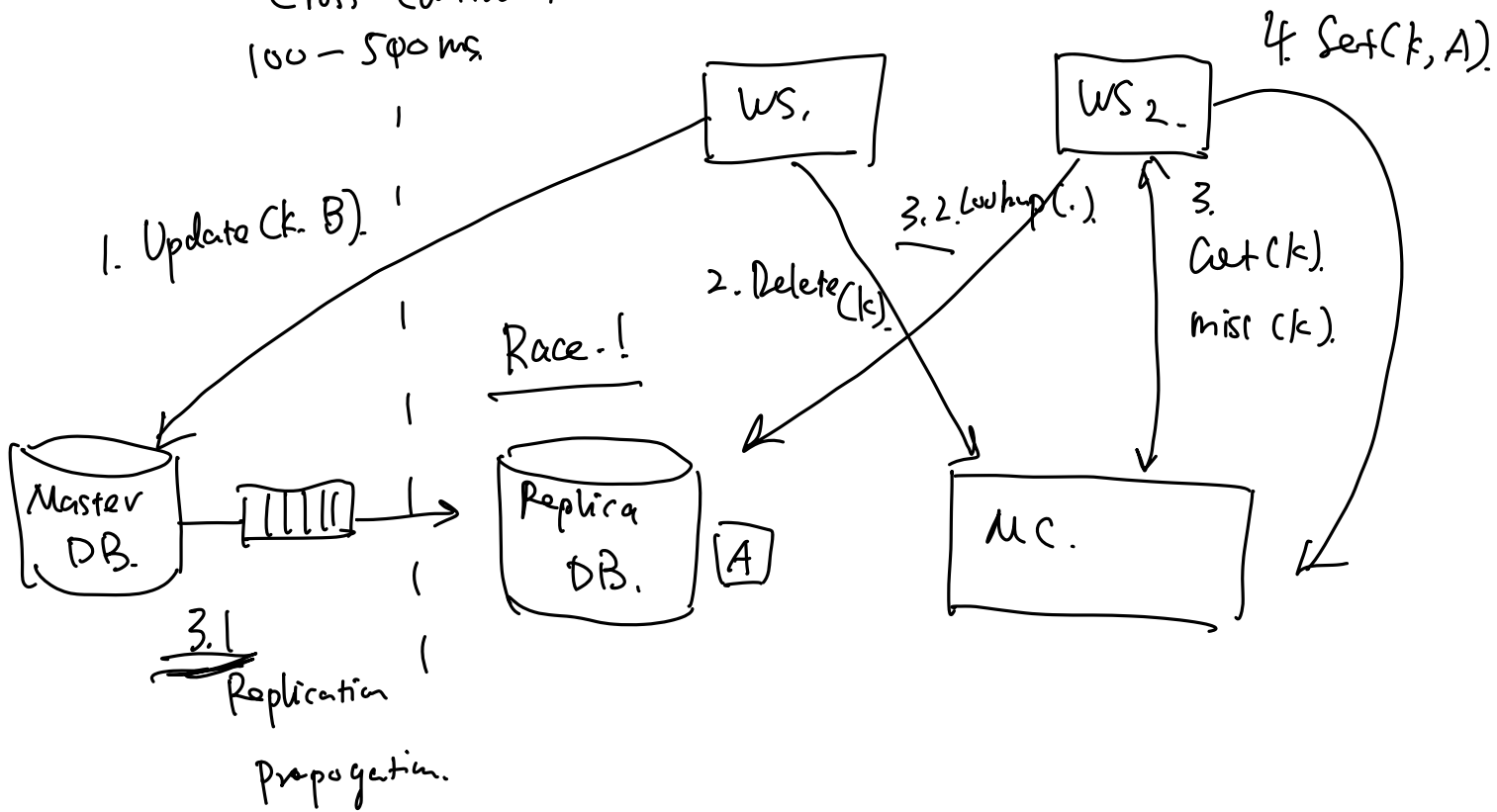
Writes in non-master

Database update directly in master

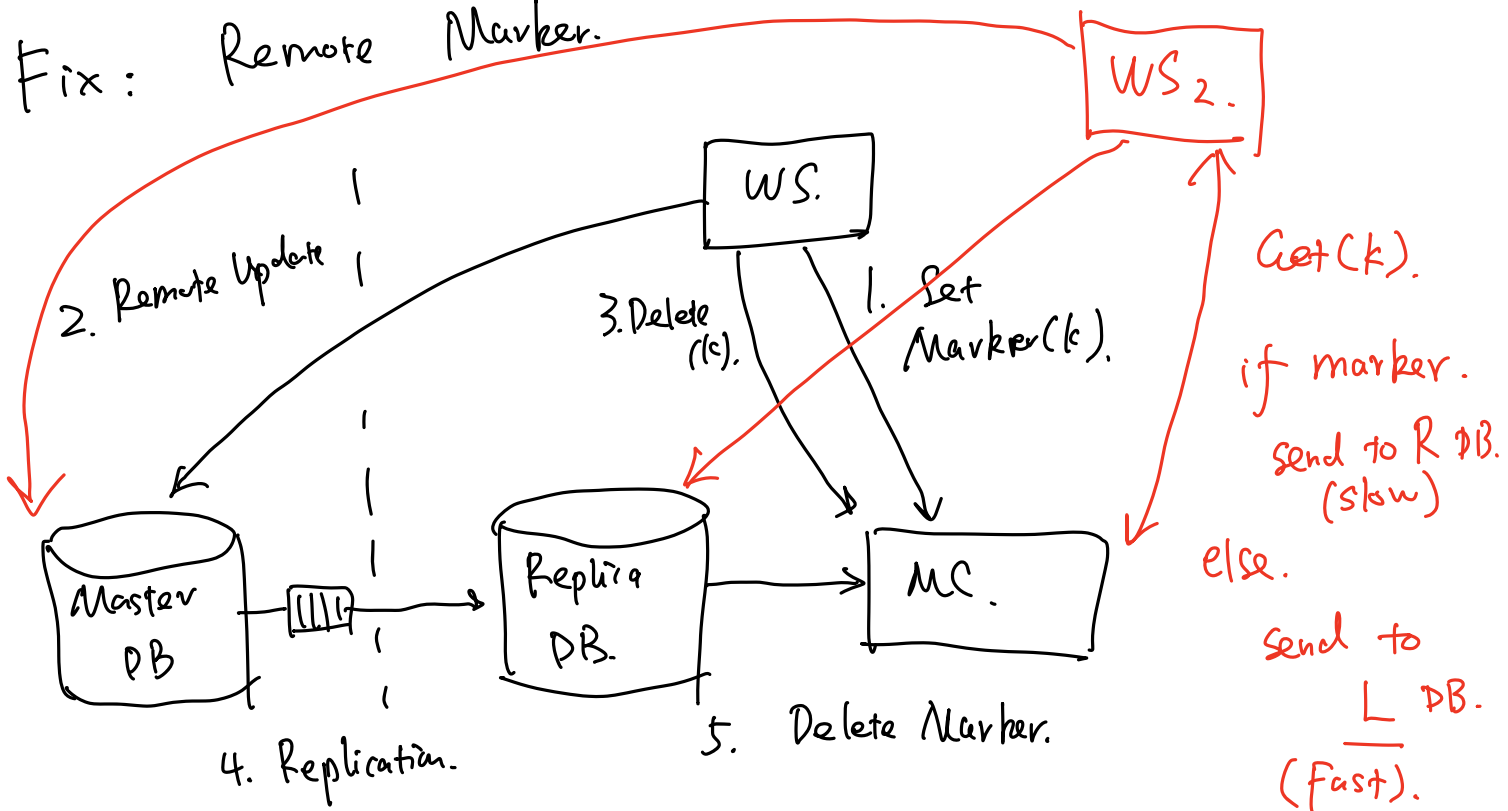
- Race between DB replication and subsequent DB read



Cross Continent
100-500ms



Fix: Remote Marker.



Remote markers

Set a special flag that indicates whether a race is likely

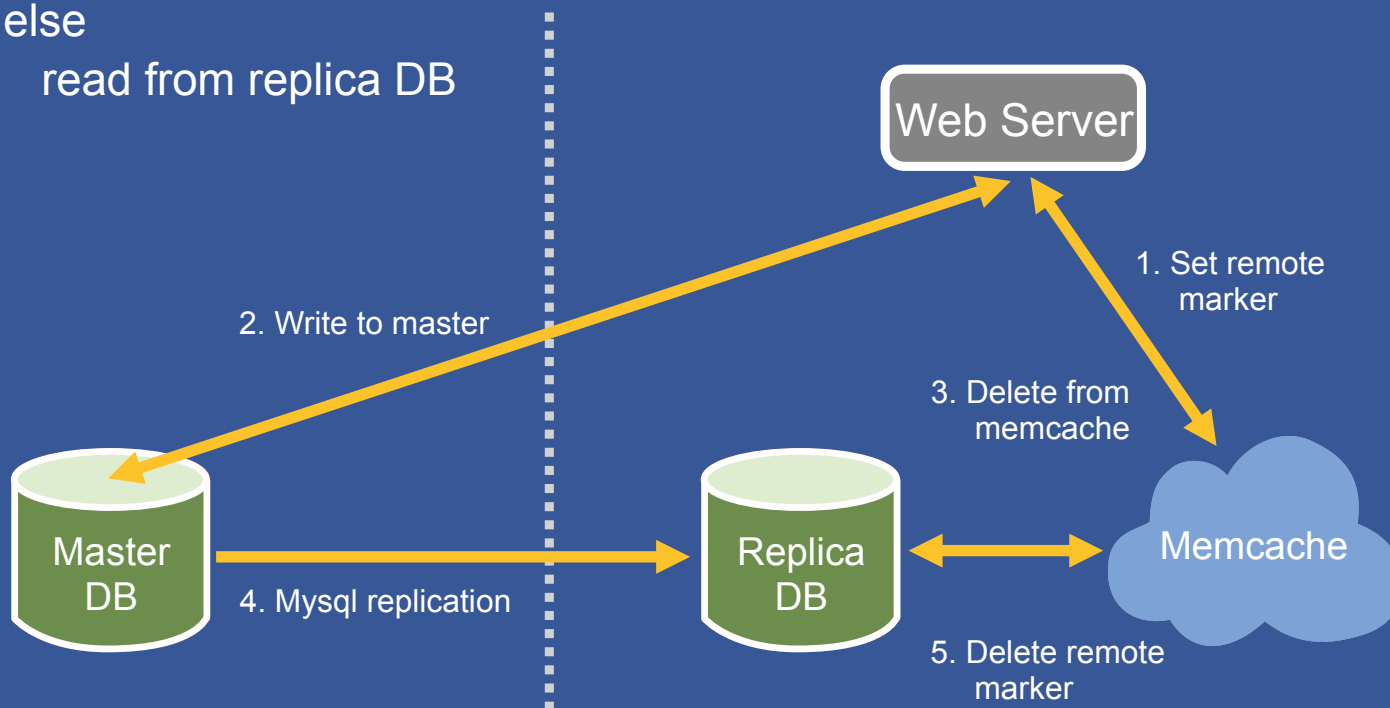
Read miss path:

If marker set

read from master DB

else

read from replica DB



Putting it all together

1. Single front-end cluster

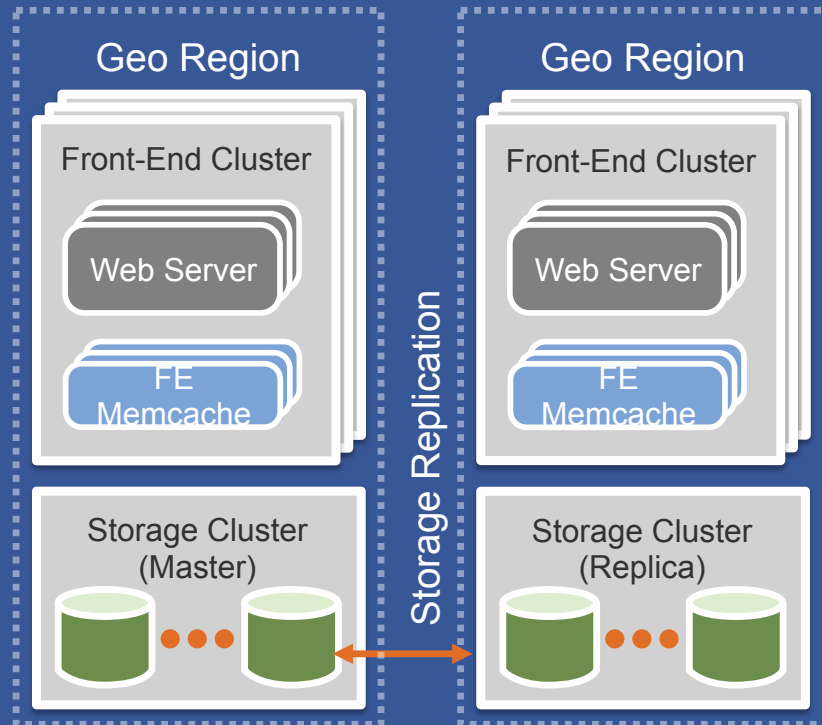
- Read heavy workload
- Wide fanout
- Handling failures

2. Multiple front-end clusters

- Controlling data replication
- Data consistency

3. Multiple Regions

- Data consistency



Lessons Learned

- Push complexity into the client whenever possible
- Operational efficiency is as important as performance
- Separating cache and persistent store allows them to be scaled independently

Thanks! Questions?

<http://www.facebook.com/careers>

facebook