

## FILESYSTEMS:- DISTRIBUTED FILE SYSTEMS.

1. Requirement: Durable Data storage + High Availability.

3 File Systems I've picked up out of many.

(1) GFS. (2) Colossus. (3) Facebook Tectonic.

\* most <sup>old</sup> powerful.  
\* seminal (original)

- + uses centralised Control plane (manager). Simple.
- Over the years, GFS Manager approached its scalability limits. However GFS cluster could survive.

So, they invented GFS-Colossus.

(2) GFS Colossus:- Evolution from GFS to meet

- ① scalability beyond Petabytes.
- ② enable more...

↓ example of evolution of system design

Principle: "Make the system only as complex as needed for present."

(3) FB's Hyperscale Tectonic System :-

+ Resembles Colossus.

+ Achieving high scalability w/ resourcefulness.

## GFS:- Google file System

Before GFS, there were single-node file system, network-attached storage, storage area networks.

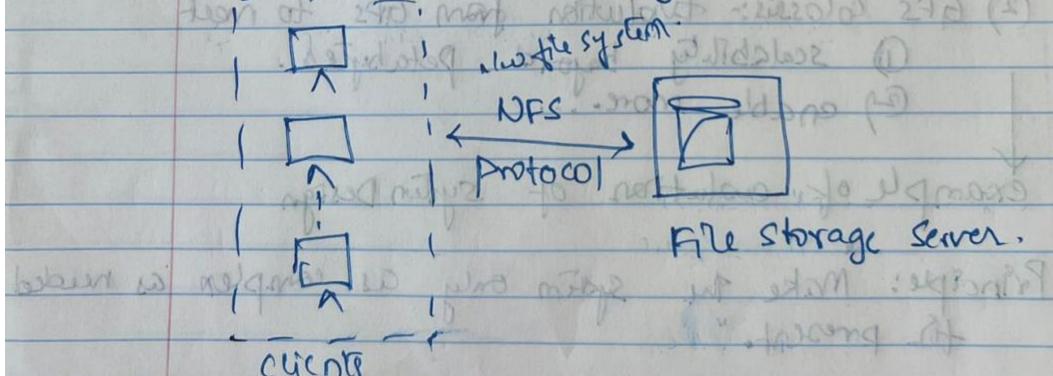
- 1. Single-node file system:- runs on a single computer and manages storage on it.

Limitations:  
1) storage capacity even if there is vertical scaling.  
2) I/O operations on disk per second.  
3) power of processor  
4) can become a single point of failure.

Only advantage: Low latency.

- > But high throughput is more valuable than low latency.

## 2. Network-attached storage (NAS)

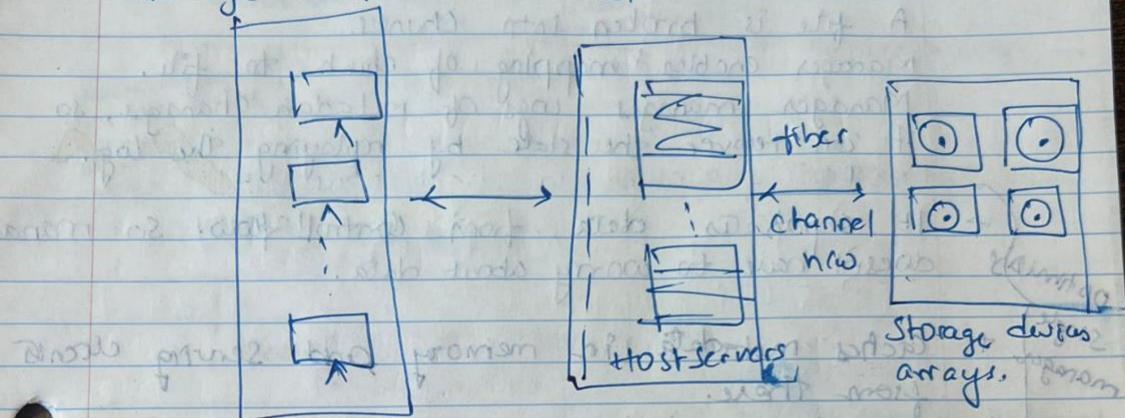


Protocols for NAS: CIFS, NFS, iSCSI

Data stored over network into a file server.

Same limitations as single-node file system.

### 3. Storage Area Network: SAN.



Adv: Easy to scale by adding more storage devices

Disadv: Difficult to manage the second network - i.e.

- Fiber channel.

- Difficult to troubleshoot any failure.
- Very costly.

### 4. GFS: + Scalability of units with cost.

+ Availability and consistency +

+ Fault Tolerance

+ Durability (backward status backed up)

+ Easy managing it

+ performance optimization

+ Relaxed consistency model

## GFS Architecture :-

2 major components :-

- ① manager node
- ② chunk servers.

A file is broken into chunks.

Manager enables mapping of chunk to file.

Manager manages logs of metadata changes, so it can recover the state by replaying the log.

- + It separates data from control flow. So manager doesn't have to worry about data.
- + Optimized single manager
  - + Caches metadata in memory and serving clients from there.
  - + Caching metadata on client side.

Q. How GFS handles single manager failure?

A: Can recover using Operation log.

↓  
this will be in persistent storage

Q. How is the time to recover low w/ Operation log?

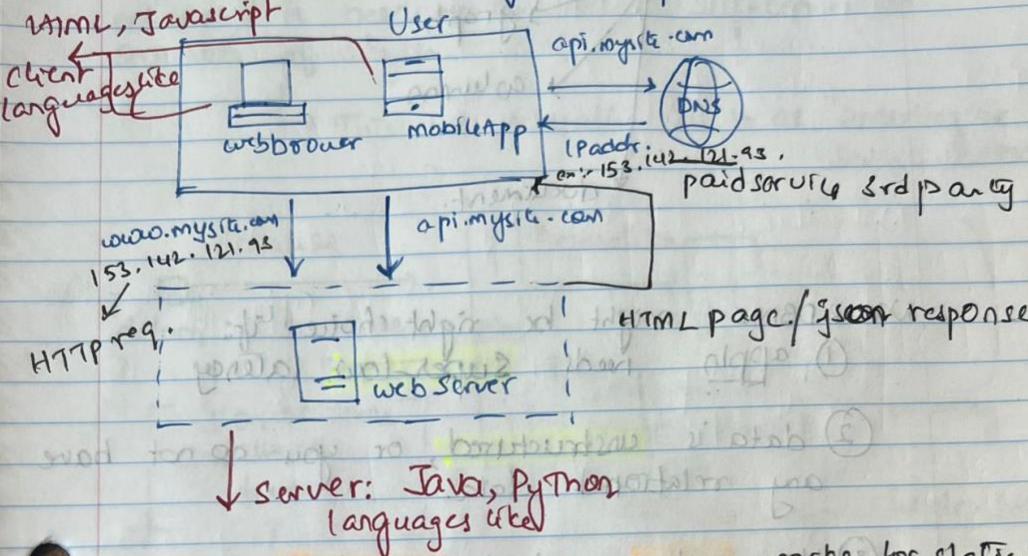
A: It checkpoints the manager state when it grows

(log)

beyond certain threshold size. Then during recovery, it loads the latest checkpoint & replay the logs after that checkpoint

We could achieve by a simple architecture, however,

HTML, Javascript



The above is the simplest approach. However, with the growth of user base, one server is not enough. So, we will separate DB from web server so that web/mobile traffic & DB can scale independently.

We have choices in which DB to use:-

- 1) Traditional RDBMS
- 2) Non-RDBMS.

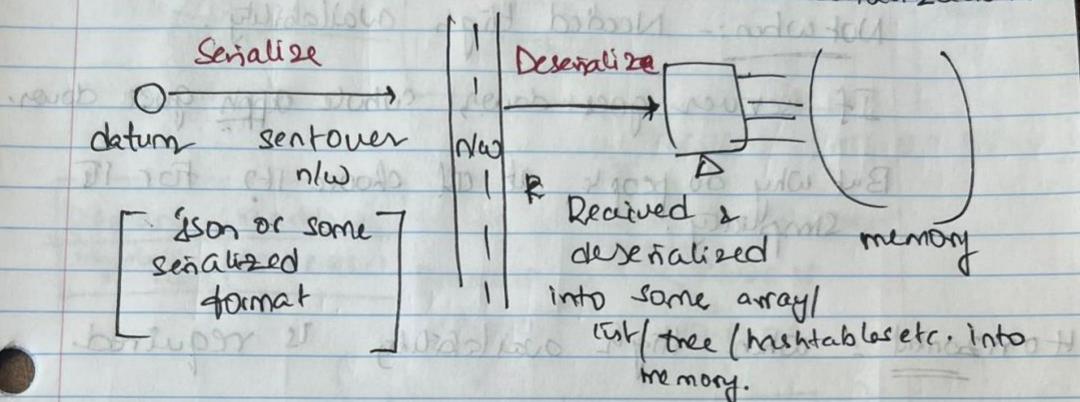
1) RDBMS: ex:- MySQL, Oracle DB, PostgreSQL  
→ rows & columns.

2) Non-RDBMS: CouchDB, Neo4j, Cassandra, HBase, Amazon Dynamo DB, etc.

(21)

But if we're sending data to a file or over a network, it needs to be in a self-contained sequence of bytes.

→ This we usually call it as "encoding or serialization".

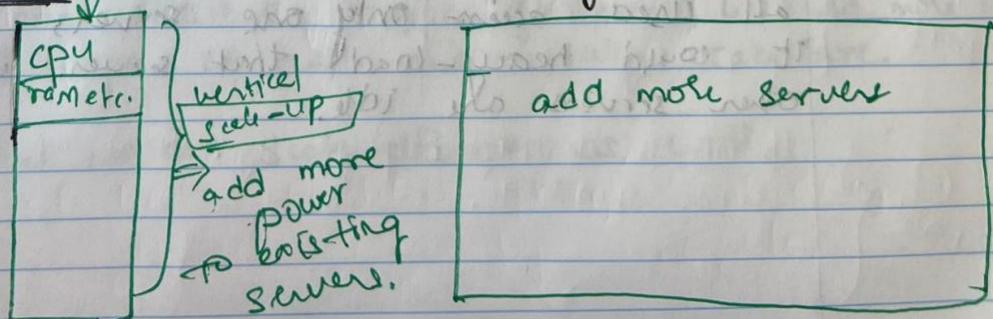


2 data Representations:  
① serialized / Deserialized ② in-memory form.

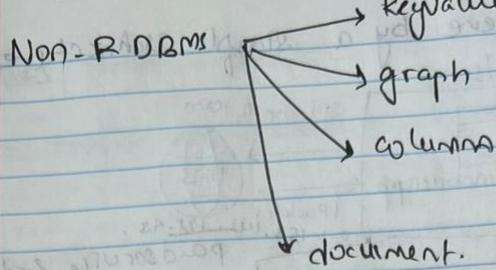
Encoding / Decoding.

SCALING      SCALING  
vertical = scale up.  
horizontal = scale out

Vertical vs. Horizontal scaling.



"Join" operations are generally not supported in non-RDBMS.



Non-RDBMS might be right choice if:

- ① appn needs super-low latency
- ② data is unstructured, or you do not have any relational data.
- ③ you only need to serialize and deserialize data (JSON, XML, YAML, etc.).
- ④ You need to store massive amount of data.

Above we mentioned "Serialization" & "Deserialization".

Serialization/

Deserialization: When we write data into memory, it is usually in some arrays, lists or hashTables etc.

In-memory, data is in arrays, lists or hashTables, etc. Why? = so that the data access is efficient & CPU can manipulate them.

(LS)

## When To Choose which scaling vertical vs. horizontal

Vertical scaling = when traffic is low

Not when:- Needed high availability.

If server goes down, whole app goes down.

But why do people at all choose it? For its simplicity!

Horizontal = when high availability is required.

Till now So, now designs have to a web servers  
we studied =

+ dB | scale it as needed

Now, when many users access web simultaneously  
access the web server, it reaches the  
web server's load limit, then it might give  
slower response.

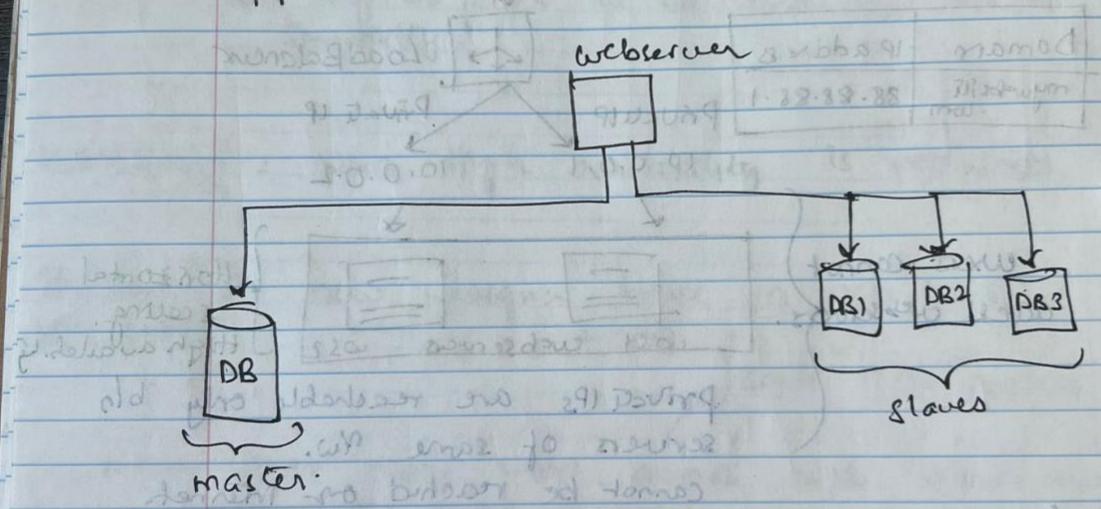
Even if file servers are there, what if  
all users access only one server, then  
it could heavy-load that server while  
other servers are idle.

## DATA BASE REPLICATION

1. usually master/slave relationship.

Master = only writes usually. = insert, update, delete.  
Slave = usually only read operations. ↳ select.

∴ Most applic. have high read: write ratios,  
people slaves are used.



Advantages: - all writes & updates happen in master node.

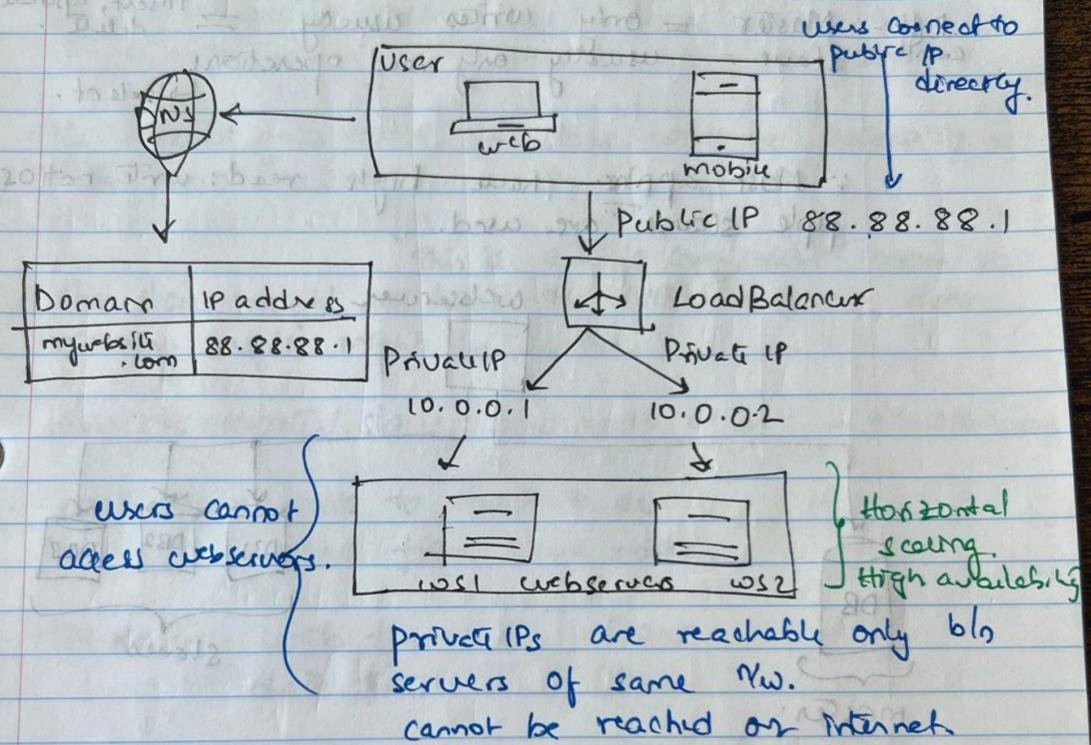
① Better performance { All reads in slave nodes.

Better performance { This improves performance, more queries in parallel.

② Reliability { If one DB is destroyed physically. There are other servers available. So no data loss.

To avoid this problem, we use a Load Balancer.

It evenly divides the traffic workload among web servers.



### Advantages:-

- 1) If one server goes down, we have ws2 which is still up.
- 2) LB can redistribute traffic. If needed, we add more servers. LB will automatically take care.

Now we can scale DB tier as well.

High availability

same as above.

performs

If all slaves go down, master takes care of the role temporarily till new slave(s) are added.

If master goes down, a slave will be promoted as the new master.

This is a little complicated process as the slave might not have the latest data. ∵ data-recovery scripts might be needed.

How do master/slave know each other? - Need to know.

> web servers write to master node.  
read from slave nodes.

TM Now, web tier  Done  
data tier  Done

Now, let us improve load / response time.  
For that we will add CACHE LAYER &  
shift static content (Js, Cs, Images, Video files) to  
CDN= Content Delivery Network.

## CACHE:

Temporary storage area.

Stores results of expensive responses or frequently accessed data.

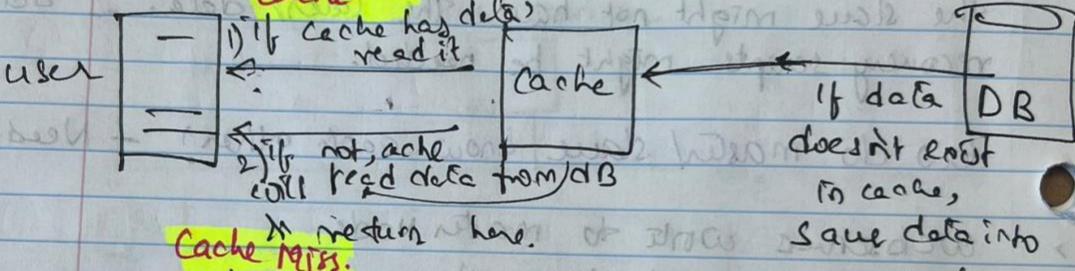
Why?

Better cache system performance,

ability to  $\downarrow$  dB workloads,

ability to scale cache tier independently.

**Cache HIT**



If data doesn't exist in cache,  
Save data into cache

**Cache Miss.**

Read-through cache strategy.



What other cache strategies exist?

Interacting with cache servers is simple. ::  
more cache servers provide APIs.

"CACHES are meant for PERFORMANCE."

Cache is a temporary data storage that can serve data faster by keeping data entries in memory.

Caches only store the most frequently accessed data.

→ non persistent data - keeps repeatedly read & written data

↓ To give lower latency.

DISTRIBUTED CACHE: multiple cache servers - so it's scalable & high degree of availability.

- CACHES:
- ① minimize user-perceived latency by already calculating & storing the results
  - ② store user session temporarily.
  - ③ serve temp data even if data store down for some time.
  - ④ Reduce network costs by serving data from local resources.

Cache is then further performed at different layers.

### Layers

Web      HTTP cache headers, web accelerators, k-v data store

①      CDNs etc.

② App      Local cache & (k-v) data store

③ DB      DB cache, buffers, (k-v) data store

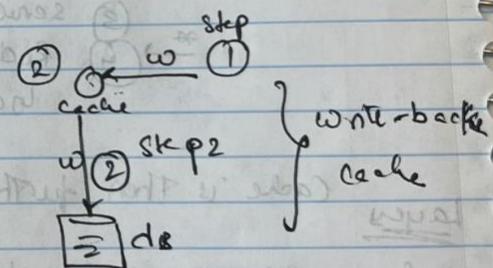
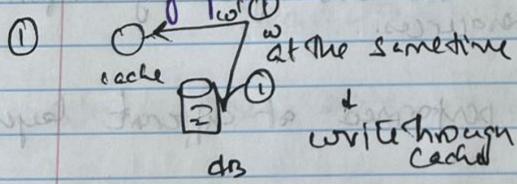
It can also be put in DNS & client side technologies like browsers or end-devices.

### Design considerations:-

to avoid

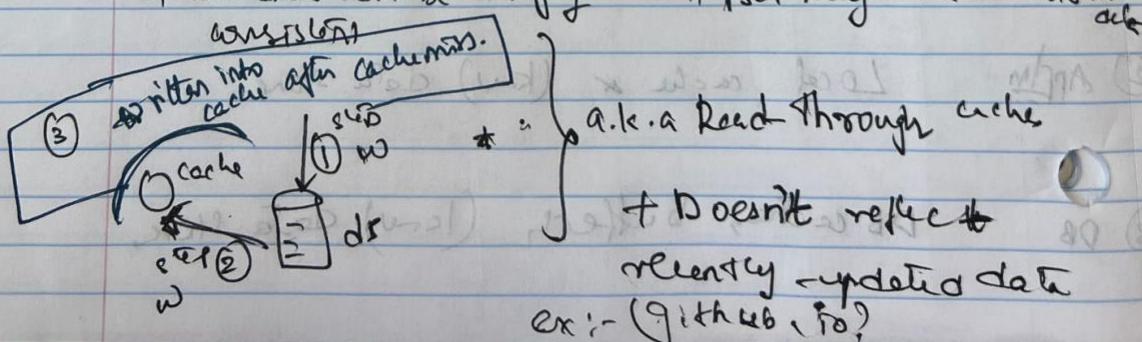
- (1) Scaling: spot (single point of failure).
- (2) If scaling, main challenge is consistency of data in DB
- (3) Eviction policy: remove unwanted data so the cache can hold the effective data
- (4) Expiration policy: Data expiration is not set, it will store permanently. So, exp date should be set, not too long or too short.
- (5) Writing policies: how to make cache have the data

### Writing policies:-



- + Writing time is more
- + But data will be strongly consistent

- + Writing time is less.
- + But may have inconsistent data



Eviction policies:- LRU, MRU, LFU, FIFO, MFU.

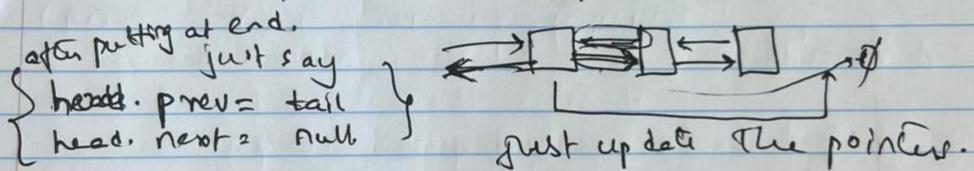
> Data structures while designing cache:-

Think what it can be we need  $O(1)$  <sup>read/write</sup> access +  
eviction policy  $\Rightarrow$  <sup>some order</sup> hash  
LRU/LFU Doubly Linked List

Why Doubly Linked List (DLL) & not SLL?

Ans In cache, you might want to move an element to the front (or end) to indicate its most recent use.

With a DLL, we can just update the pointers



Otherwise in SLL, we would have to traverse all the way in the list to find the previous node -  $O(n)$ .

Both SLL & DLL can be used.

→ has adv: for constant - removal & insertion at both ends.  $O(1)$ .

Overall: Cache Design Internals:-

- ① Hash Map
- ② DLL
- ③ Eviction Policy,

BOM are open-source.  
 BOM have client/server models.

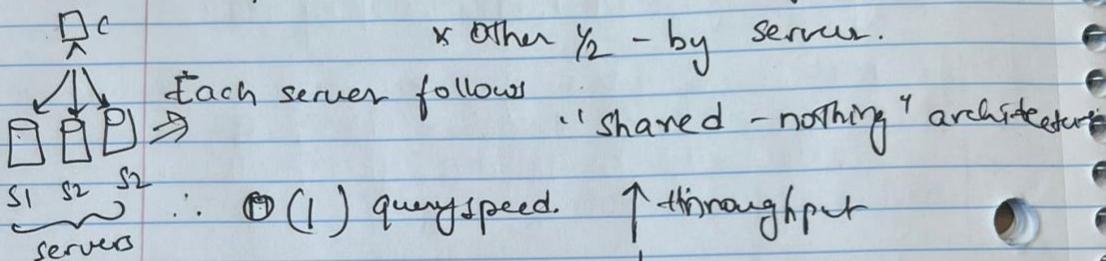
Memcache vs. Redis

① Memcache :- (k, v) store.

BOM "k" & "v" are strings. So, Memcache cannot manipulate it. This also means any data that will be stored has to be serialized.

client/server :-  $\frac{1}{2}$  work done by client

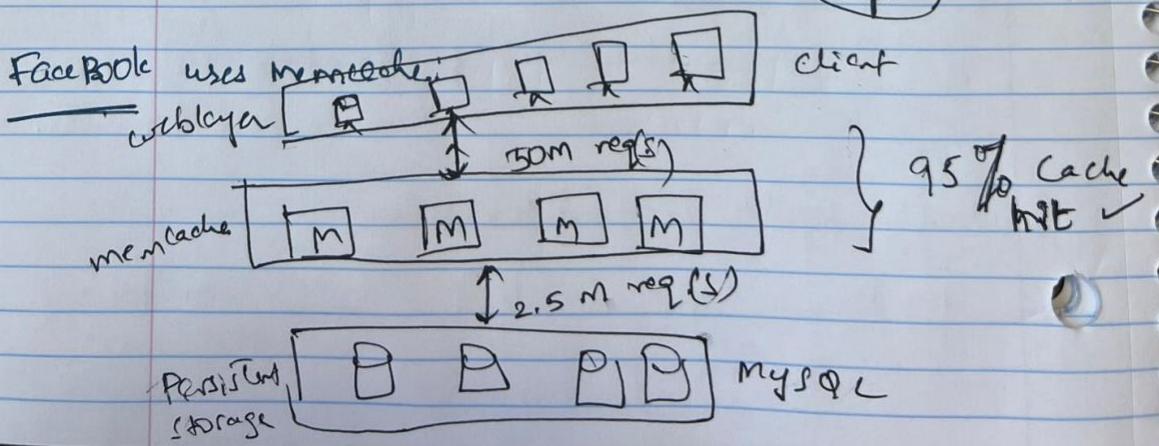
$\times$  other  $\frac{1}{2}$  - by server.



$\therefore$  ① (1) query speed.  $\uparrow$  throughput  
 $\downarrow$  low latency

✓ good for horizontal scaling

ex: code in memcache :- Cache.set("myKey", "Hello", 3600 \* seconds)  
 ↗ exp.



Redis: can modify data w/o waiting n/w bandwidth.  
stores data in (data structures) doesn't need serialization/deserializ  
both Redis & Memcache belong to NoSQL  
family. ↳ Cache / database(k,v) / message broker

in-memory storage.

Memcache is preferred for smaller, simple read  
read-heavy systems. whereas Redis is useful  
for both read & write-heavy systems.  
good for streaming videos as well.

Memcache vs. Redis  
data stored (strings only) (data stored in data structures)

- |   |  |
|---|--|
| + can't modify data<br>∴ n/w bandwidth wasted in uploading & downloading. | + can modify data<br>∴ n/w bw not wasted                         |
| + wastes time & effort in serialise / deserialise.                        | + saves time & effort by avoiding & serializing / deserializing. |
| + servers do not know each other.   |  |

## CDN: CONTENT DELIVERY NETWORK

CDN: is a n/w of geographically dispersed servers to deliver static content. like images, videos, css, js files etc

cache of pages

→ static (images, videos, css, js)  
→ dynamic (based on req. p aM)  
we don't cover querying things etc.

CDNs are caches of static content of web pages.

- "Globally dispersed servers."
- (+) The closer the CDN server is to the user, the faster the response is.
- (+) The further away, slower response.

Amazon provides CDN. Akamai also provides CDN.

Adv + static assets (JS, CSS, images etc) are fetched from CDN  
Hence DB is lighter.

STATEFUL vs. STATELESS - stateless user session

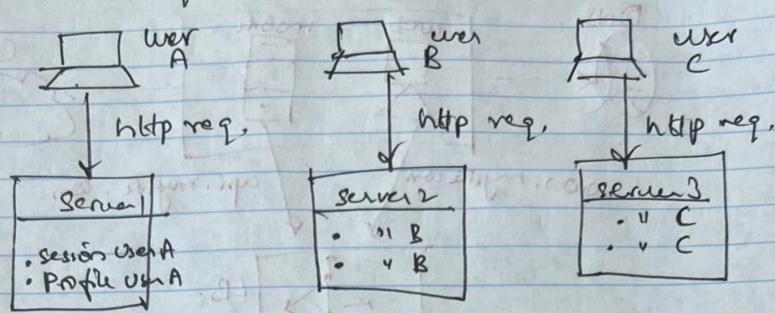
- (1) Stateful server: remembers client data (state) from one req. to other.

Stateless: keeps no. info. It is usually stored in persistent storage  $\Rightarrow$  RDBMS or NoSQL

- (Q) Who does it remember/not remember?  
Ans: Server.

## non DB based on DA

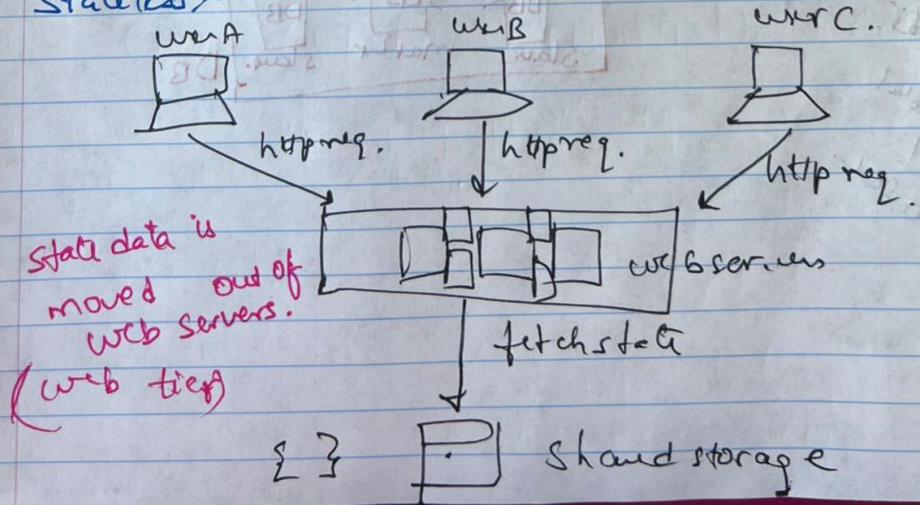
ex:- Stateful :-



- ① Problem is :- If User B by chance (Via LB) sends req. to Server 1, it doesn't contain the data. Then either it fails or needs to be rerouted to Server 2.  
Overhead.
- ② What if all requests come from User B? Then only Server 2 is kept working, while rest are idle.

∴ Stateless.

Stateless



All we learnt till now.

