

# System Design

↳ Gaurav Sen Playlist (Quick Intro,

Overview (up  
Topics)

\*

Horizontal Scaling:  
(Pros) { Scales well  
Resilient

Vertical Scaling:  
(Pros) { Fast inter process  
communication  
Data consistent.

We make hybrid systems taking big boxes  
(natively scaled machines)

which scale well horizontally.

\* Load Balancing

used in distributed cache.

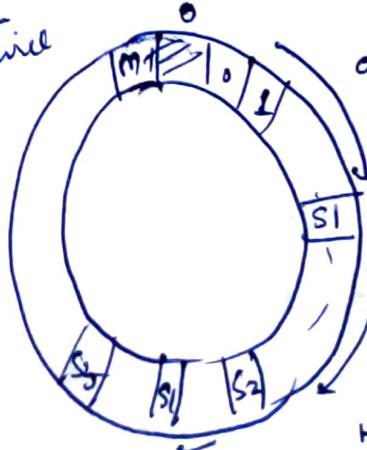
Consistent Hashing

assuming a  
load balancer,

uses a hash function value  
to map to a server.

while adding a new server,  
we need to make balance remains  
proper, so that old request ids,  
which were mapped to previous  
servers, keep getting served from  
old ones so that cache  
remains in a good state  
for the servers. Hence  
existing consistent  
hashing.

\* each  
server will  
have  
K hash  
functions.  
Hence  
they occur  
on ring  
K times.  
(each)



all rep.  
here  
mapped  
to S1

$0 \rightarrow M-1$

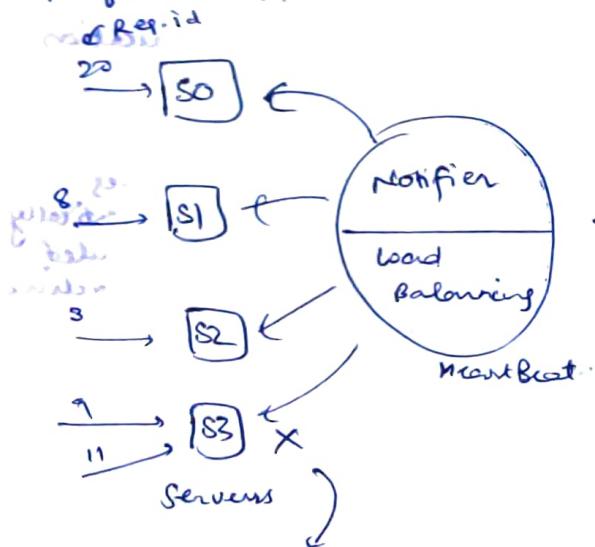
request Ids hash  
 $h(r_i) \% m$  value

\*  $S_1 + S_2$  all servers mapped by K hash functions.  $S_1 \rightarrow$  server value hash wrapped

$S_2 \rightarrow h_2 \% m$   
Hence they occur on ring K times.  
We now have more consistent setup.

## Message Queue

\* Example dominos pizza → for all requests received, add them in a list. Thus client + pizza shop can work asynchronously and have faster performance.



S3 crashes. Heart beat check will confirm its dead. Route them again. Here 3rd is also incomplete, but load balancing will ensure 3rd does not go again as it is being done by S2.

Ties combination of Notifier, load Balancing + HeartBeat

= Message Queue

e.g. Rabbit MQ, zero MQ, JMS

ID	Contents	Done
1	Pepperoni	Y
2	—	N
3	—	N



DB

## Monolith vs Microservices

Separation of responsibilities → Microservices

### Microservice

→ Good for smaller teams

→ Complicated

→ One service goes down, others still up → hence high availability.

→ Parallel development can be done.

### Monolith

→ More context needed

→ Complicated deployment

→ All calls are

Remote procedure  
hence faster

If S1 only calls S2

S1 → S2

its better to include S2 inside S1 and make monolith

↳ Eg of monolith  
↓  
stack overflow.

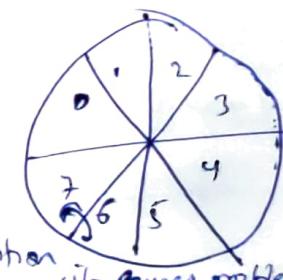
## Database Sharding

↳ Partitioning data  
\* You should (partition) data on a attribute say id/location, etc.

### Problem

↳ JOINS

When you need to run a query flat w/o sharding data from other partition



↳ Causes problem (expensive)

\* To have dynamic shards (flexible), measured.

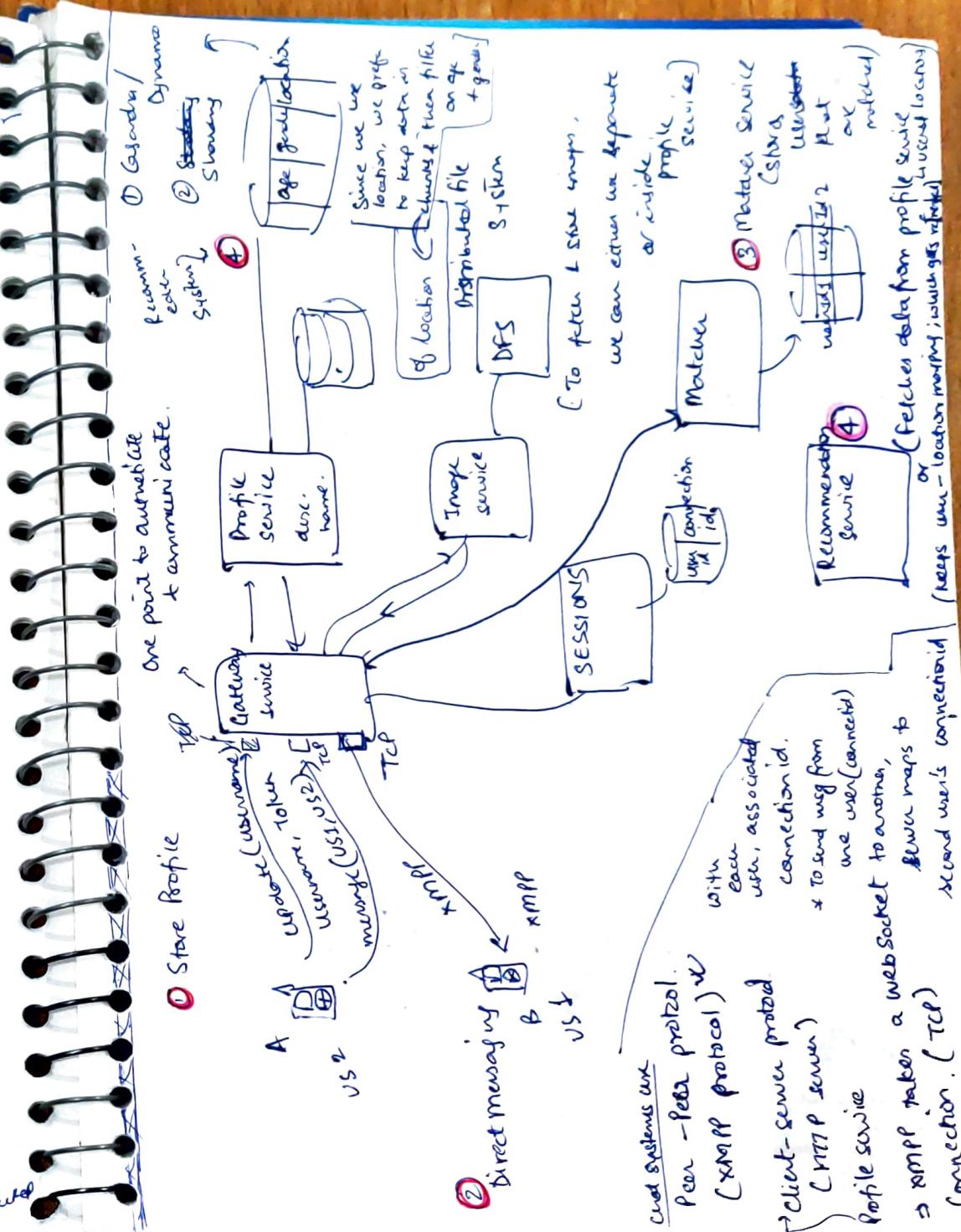
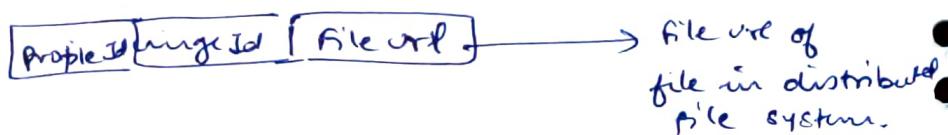
Advantages that a DB gives you over file system -

- Mutability [Ability to change] but why they are useful for images] Not needed
- Transaction guarantee → will be risky since if select \* performed huge chunk will be moved.
- Indexes (Search)
- Access control

↳ one gets same using file system.

You will  
not be searching on the  
content of the file.

- Advantages of file system I,
    - cheaper
    - faster ( ~~cost~~ since separate place of large files)  
can be done with vertical partitioning using DB but again slow & nightmare !!
    - Content delivery network  
(~~cost~~ since these are static files, they can be hosted on a CDN)



# WhatsApp System Design

will be designing these features -

- ① Group messaging
- ② Sent + delivered + read receipts

- ③ Online / last seen

- ④ Image sharing

→ similar to  
Tinder design

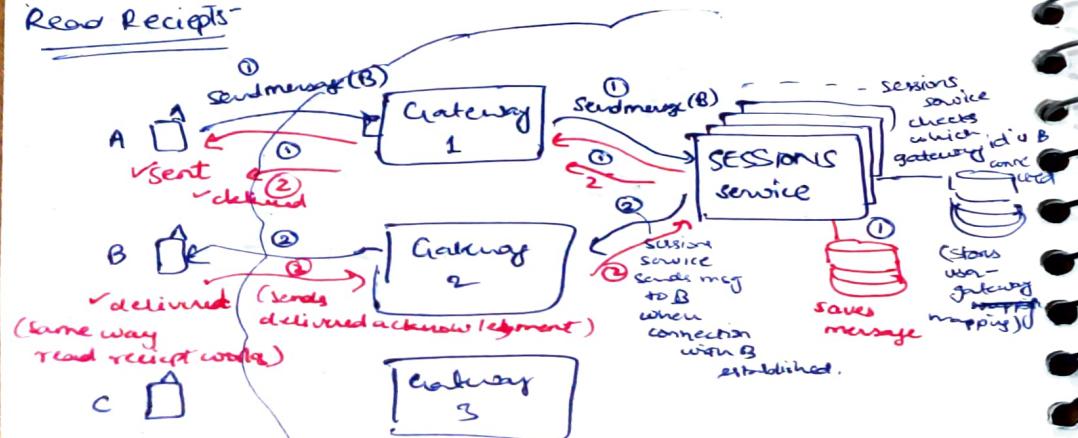
- ⑤ Chats are temporary / permanent

- ⑥ One to one chat

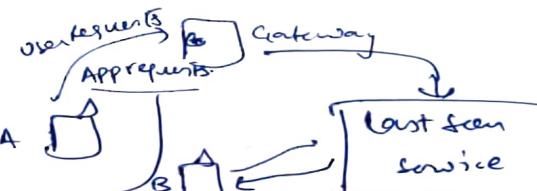
\* TCP connection (BROKEN over  
web Sockets)

used as server needs to  
communicate with client.

One to one chat  
+  
Read Receipts



## ③ Online / last seen.



A

B

like log policy etc.

make sure you  
don't update

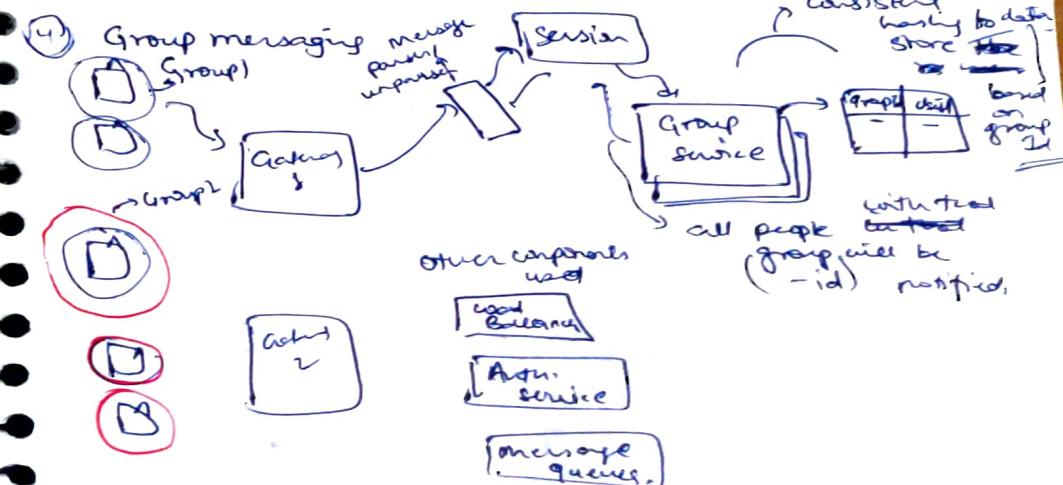
last seen Table  
for app requests

User	lastActivity

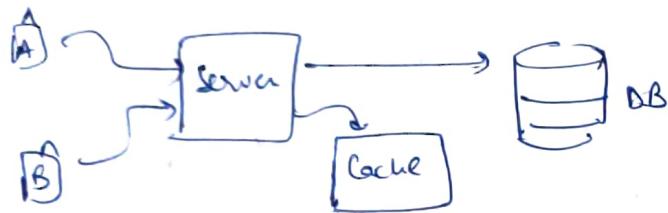
\* when user exists  
in main app,  
a report is sent  
to update  
last seen  
table.

\* B is checking last seen of A

\* Tell user is online, the user-reports  
one periodically sent which keep the  
last seen table entries updated (to current  
timestamp).



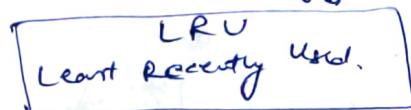
## Distributed Caching



### Benefits -

- Save Recomputation
- Reduce load on network (less calls)
- Reduces load on DB
- You can't store everything on cache as cache hardware is expensive → SSDs compared to commodity hardware of DB.
- Large data on cache will ~~increase~~ <sup>increase</sup> search times at cache.
- Hence only relevant data (that seems reusable) should be kept on cache.

Hence to clear cache ; [Cache eviction]  
we have Cache policy



↳ oldest entry  
Keeps getting pushed below, until it is evicted from the cache.

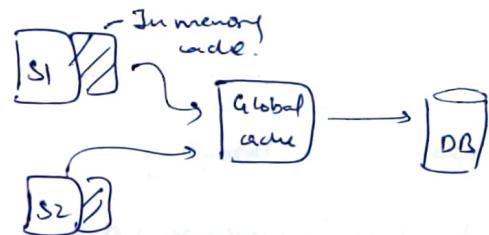
↳ eg-  
Celeb posts. Kept at top of cache. Everyone views it.  
Slowly keeps getting pushed below. Finally removed.

\* If cache not used correctly it is harmful

everytime before going to DB,  
check in cache

data keeps getting added /

removed from cache without ever getting used



Cache can be placed in-memory at server level  
or can be a global cache.

In - memory

→ Faster.

→ Not resilient (server crashes, cache gone)

→ Not accurate/  
Cache consistency an issue.

(Cache at all services need to be at sync)

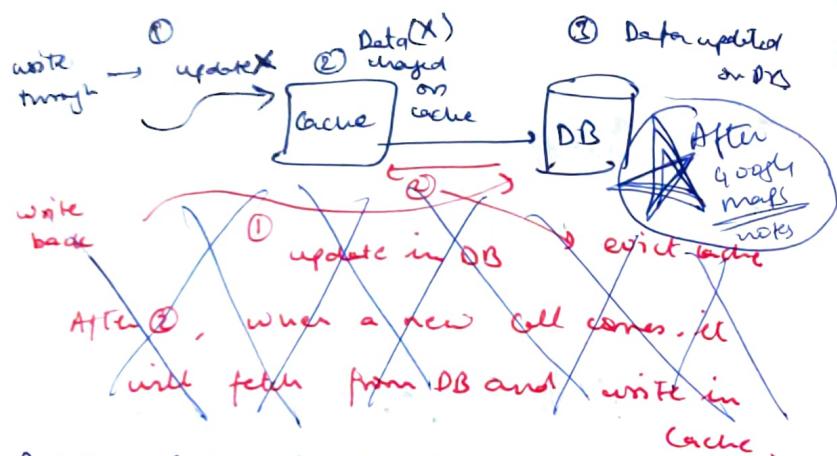
Global cache → Slightly less fast

→ Accurate

→ Can be scaled independently of servers.

Characteristics of cache - determines accuracy & cache data

- ① write through
- ② write back



Problem with write through ]

- \* If it is in-memory cache, write through isn't bad.

Write Back ]

- \* With write back, for each update we invalidate cache which is expensive.

What can be done ]

Hybrid!!

If it is some non-critical data, we can have data updated on cache (write through), then after some defined time, bulk change can be done at DB.  
For financial data, we only have to use writeback.

## API design

Imp things →

① API name

② Parameter defining

③ Response object defining

Send only what is required

④ Errors

Throw relevant errors

client asks for first 10, then next 10, ...  
so on

→ If you have large response, good way is to break it by either

→ Pagination

client says return first 10, then next 10, so on

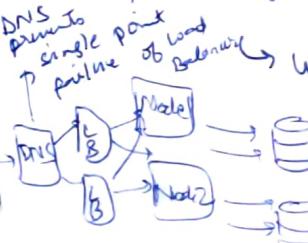
→ Fragmentation

Server gives out broken responses

one-by-one with an avoided piece say TCP no. 10, 11, ... last one -1.

## Single Point of Failure

- \* If your system has a single point of failure it means system is not resilient.



We can have more nodes (servers) master slave (DB) system in multi-region.

→ YouTube Capacity Estimation

## → NOSQL Database

SQL → Structured Query Lang.

No - SQL

represented as

ID  
name (ID)  
address  
Age: 1  
Gender: 2  
Name

ID2  
name  
address  
Age: 2  
Gender: 2

By default supports horizontal partitioning (easy sharding)

Data stored as →

ID	Name	Address	Age	Gender
123	John Doe	123 Main St	23	M

ID	Country	City
1	India	Delhi

Adv. of NoSQL over SQL

① Insertions & retrievals easy  
(requires JOINS in SQL)

② Schema is easily changeable.  
(In SQL, would need to add data for all. Hence costly querying as DB locked for that interval.)

③ Built for scale

④ Built for metrics/analysis/aggregates

Not necessarily same

JSON

Disadvantages →

① Not good for updates (deletes/inserts)  
(consistency across nodes is an issue)

② Read / Fetch of a value is slower.  
(goes to a block, finds key (say age) then returns)

③ Data doesn't have relational property.

④ Joins are hard (to join on a property, scan through each block then merge ----) in SQL, easy as it is built for joins.

Cassandra (Ex of NOSQL DB) →

Key space (0 → 100)

hash (repId)

requestId

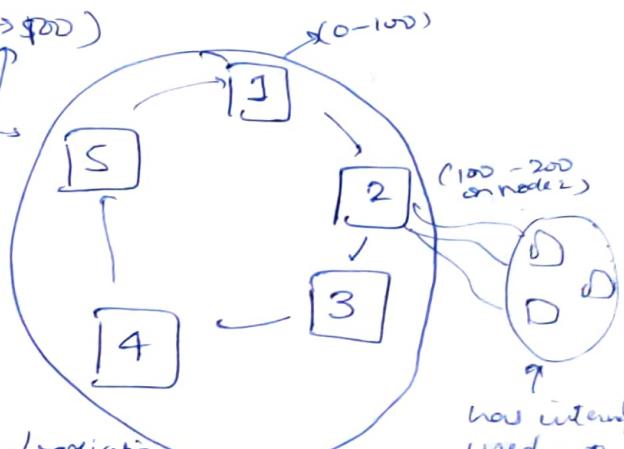
→ Data is replicated with RF of 2.  
⇒ Replication factor

∴ when data added

• Hence load balancing good.

• RF provides redundancy/replication.

→ eg data at 2 replicated 3rd cluster ↑  
★ Cassandra has concept of consensus. → achieved by quorum,



how it can be used another cluster

If quorum = 2, with RF = 3,

→ Two nodes must agree for a value  
out of 3,  
before sending it to user.

else → DB error -

→ Another feature of Cassandra which stands out is way of writing to DB.

so, it maintains a log file

key	value
key1	-
key2	-
key3	-

so, write just happens where the pointer is currently at.

→ NO scanning of last index ~~for insertion~~ <sup>to</sup>

Secondly, after a period of time  
this log file is moved to a SSTable.  
(Sorted String Table)

Keys are sorted.

key	value
key1	initial value

immutable

Say initially

a write for key 123 → name = John Doe

again later,

key 123 → name = John Doe - NC

So, multiple keys in DB.

No problem for consistency, we can use latest timestamp.

But storage ↑.

So, Cassandra does compaction.  
(+ Elastic Scan)

Two SSTs are merged.

Merge sat → 2 sorted arrays merged.

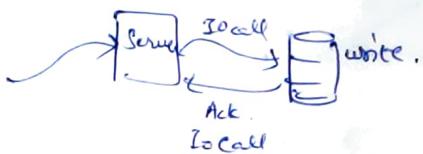
Space complexity  $\mathcal{O}(\min(m, n))$

Time  $\mathcal{O}(n)$ .

For delete, Cassandra has concept  
of Tombstone. So basically when  
key 123 deleted, tombstone (flag) <sup>kind of</sup>

is set. Before read, you verify  
tombstone to be unset.

Scaling write operations.



- Lesser IO calls  
Condensed write queries.

- Write using linked-list.  $\mathcal{O}(1)$  time.  
↳ Log files use linked list DS to store things.

Hence Adv →

- ① Lesser IO
- ② fast writes

disAdv →

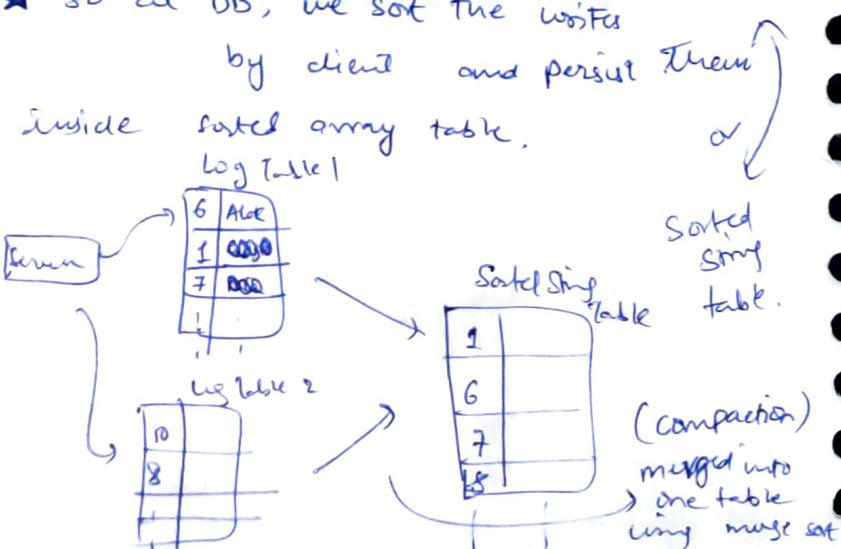
- ① slow reads
- ② Additional memory consumed at server (for ~~compacting~~ condensing queues)

kind of ok (Not big problem)  
problem!!

To improve reads, we can use  
Sorted Arrays.

Hence for write →  $O(1) \rightarrow$  linked list  
read →  $O(\log n) \rightarrow$  sorted array table.

\* So at DB, we sort the writes by client and persist them inside sorted array table.



\* Deciding when to merge,  
↓  
merging every small file is a heavy operation. To insert 6 size table, in  $10^6$  records, required  $O(n \log n)$  time!

→ Hence we merge when compaction time is reasonable.

Also merge intelligently,



another 6  
another 6 } merge, then merge to get 24.

6 → 6

12 → 12

To speed up ~~read~~ reads, another step can be done is to keep BLOOM FILTER.

e.g. →

14	12
----	----

Do a search → Helps reduce storage.  
comes. Do not set → Not present

When a word like CAT, CORONA present, flag for this set.

so for each sorted string table, we have a bloom filter.

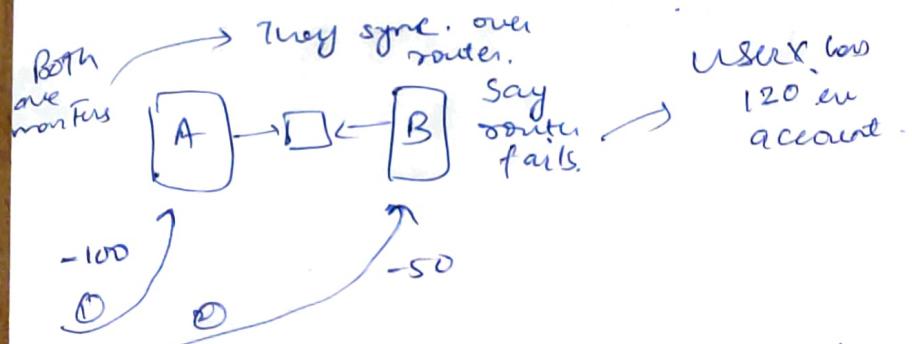
when merged → a new bloom filter created.

\* before searching in actual SST, if not present in bloom filter, it will not be present in SST.

### - Distributed Consensus & Data Replication Strategies -

\* For Databases, we use Master-Slave arch.

why can't we use Master-Master arch?

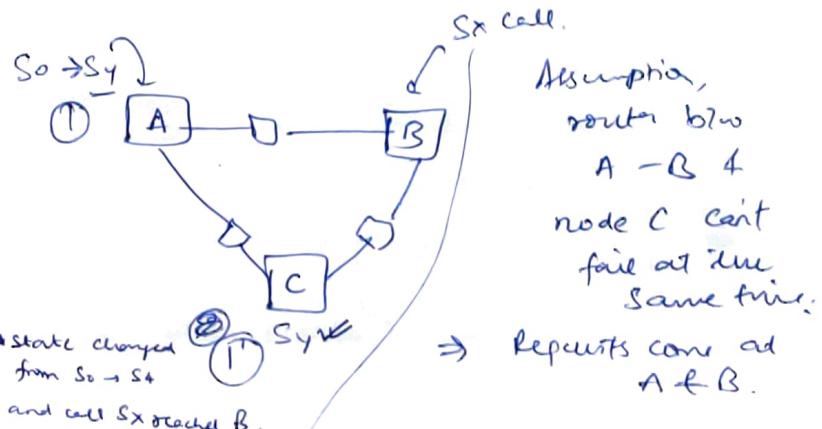


call ① to -100 comes, A assumes himself to be master (as connection with B) lost.

call ② to -50 goes to B, (B has entry of 120 still) argues B master, deducts 50. So in balance 120, 150 was deducted

\* This problem is called split-brain problem.

To solve this we can have another node (C),

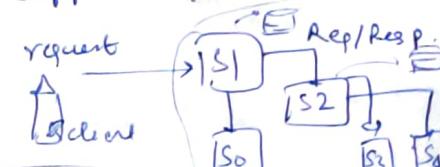


when cell at B for Sx, B calls C to update state no longer Sx, it's Sy. So transaction fails for user. B then pulls state & C. (i.e. Sy)

### 27/06 Publisher Subscriber Model

Need →

Suppose system like



\* Here problem

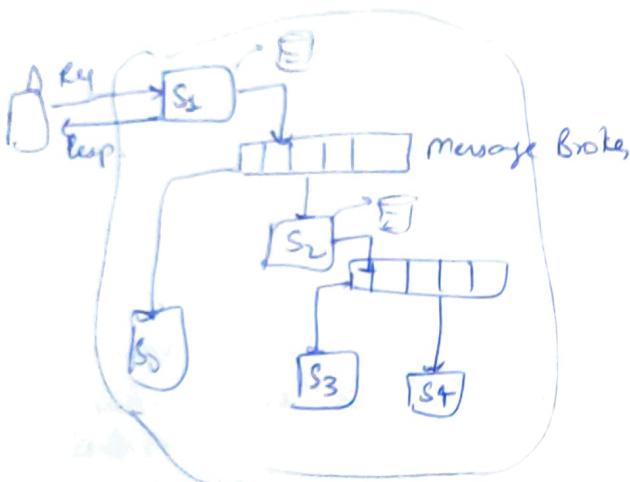
→ long wait times at S2. ⇒ If S2 fails, then returns to S1, data is inconsistent at S1 & S2.

S1, S2 → series

S1 → S2 → Async calls

S0, S2 waits for T resp. from S3 & S4 then returns to S1.

So, better way to use message brokers  
 ↗  
 Kafka, RabbitMQ



S1 publisher, S2 & S5 are subscribers to message broker.

S1 works to message broker/message queue,  
 (publishes)  
 and returns rep to client without  
 any wait.  
 Resp of message broker that  
 request will  
 be sent to S2 & S5.

Advantages ↗

- Generic message/event can be sent by message services to message broker.
- Transaction guarantee ↗  
even if S2 down, message broker ensures that at some point of time (when it is live) these messages (in queue) will be replayed on S2.
- More Scalable.

### Disadvantages ↗

① Poor consistency → Say financial system

So → invoice service

(deducts bank commission per transaction)

→ S2 down.

→ for each rep from client

S1 will be processed.

S2 rep. in message queue. Hence inconsistency in amount.

② Not idempotent ↗

Say when S2 writes on its message broker,

Example ↗

Gaining services  
 (analytics)

→ TWITTER uses it ↗

Precise use case.

Tweets get published, Subscribed

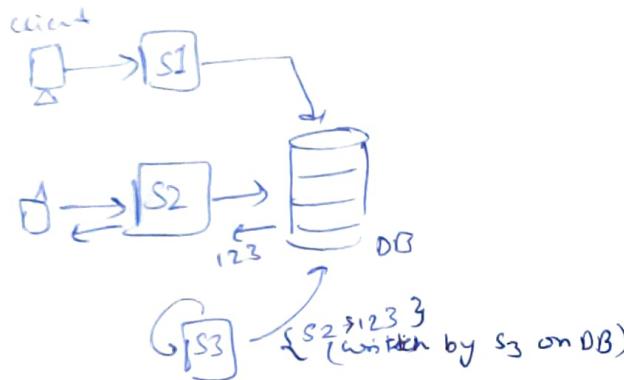
Accounts get notified ↗

→ 4/7/20

— ANTI-PATTERNS —

↳ Using DB as message queues is considered (Should be avoided) an Anti pattern

## Disadvantages of using DB as message queues



- ① Polling intervals ~~needed~~ to be correctly.  
Short intervals → high load on DB  
Long intervals → ~~inefficient~~ / bad user experience.
- ② Read / write messages from DB will have problems like locking etc.  
Also DB's are optimized for either read or write.
- ③ Delete (of read messages) will be an overhead.
- ④ Scaling difficult.

So generally, if system is a small service → DB fine for message queues.  
If scalable → message queues.

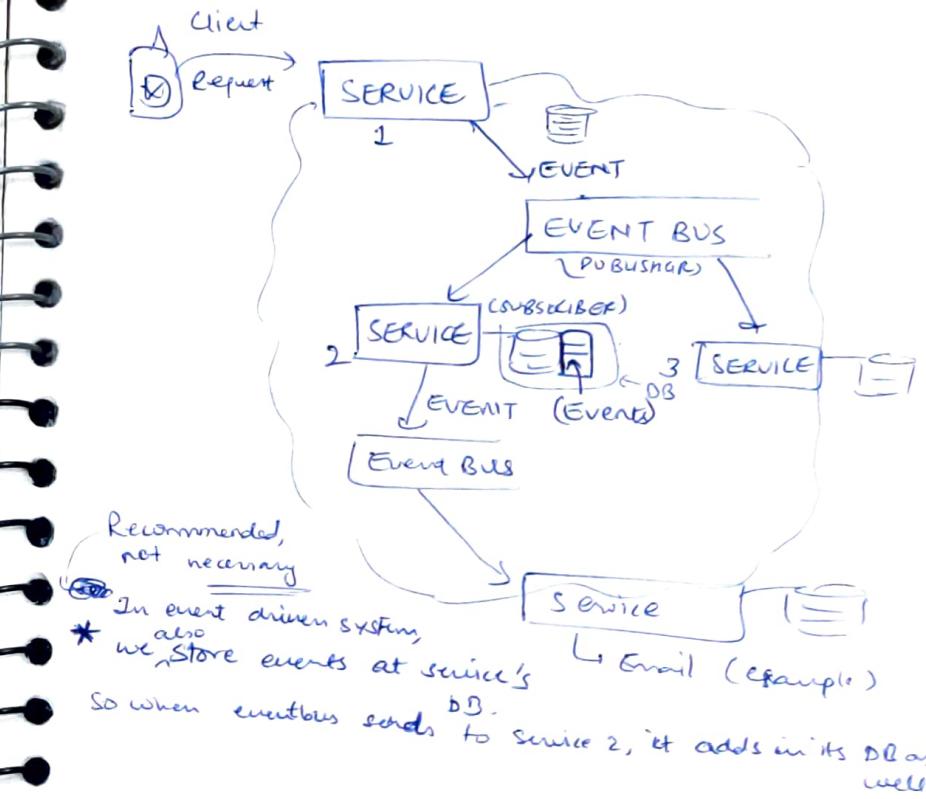
## Adv of message queues →

- ① No polling - when new message for a service, the message queue pushes it to the service.
- ② Read and write / delete of msgs all responsibility of message queue.  
+ Scale for w., it is a black box.

## Disadv of message queues :

- ① For not a large server, DB already handles. So unnecessary complexity added.
- ② Infra setup costly; requires trainings.

## ⇒ EVENT DRIVEN SYSTEM



Examples of Event Bus Systems → GIT, Read, Game services (Pubg, CS)

Advantages of storing events in DB

Overall advantages of E/S

- ① Availability: for some reason, message queue / parent service unavailable, if event stored in DB (not yet consumed), service can consume/work on it. Hence high availability, but we have poor consistency. (Message queue)

## ② Easy Roll back

↳ If some bug happens in service after execution of an event, if event stored in service, we can rollback to previous state.

## ③ Replacement → (Say Service 1 wants to replace Service 2) -

So you can replay all events done by Service 2 and then subscribe to corresponding event bus.

## ④ Transactional guarantee -

If an imp. task like sending invoice mail done by service, if we store event at service, we ensure that it will be executed

## ⑤ Stores intent

↓  
With intent of each event stored, you can select / remove set of events to create desired data state.

## Disadvantages →

- ① Consistency (as in adv first pt.)  
② ~~Version~~ Replacement not possible for Gateways. (Creating new service from a Gateway → service difficult)  
③ less control (as request-response, you can control timeout, etc. here message bus as a component reduces control of service)

→ Since there can be many events coming, (re-applying events, rolling back can be difficult)

Hence it is good to squash events in DB.

## Other Practical disadvantages (practically faced)

- ① Since S1 ~~resp.~~ is just Event Bus, it's difficult to determine complete flow of input → email.  
② Migration difficult (to move from a fluent & listening/sending service, to a rep-resp. service is costly migration)

NUTSHELL In Event based systems, services send events if they think needed by other service.

In reg/req, a service gets request, fetches data and completes execution. All about data. Around this point,

# INSTAGRAM SYSTEM DESIGN

- ① Store / Get images → Same as Tinder.  
(Filesystem, CDN, etc.)
- ② like / comment  
element
- ③ Follow someone.
- ④ Publish a news feed.

## ② (like / comment)

### Entity Relation (ER Diagram) -

Like

Like	Activity ID	User ID	Type	Time/Fair	State
1	POSTID		comment		0/1
2	comment	30	POST		

Comment					
ID	Activity ID	TGT	User ID	Time/Fair	State
1			10		

LIKES (AGGREGATED) ← for get all likes  
(we don't want to do)  
Select \* on like table

Imagine done  
on news feed  
for every post

POST (2 way)				
ID	Author	User ID	Time/Fair	Activity ID
1	user1	user2	Time	Activity ID

ACTIVITY TABLE	
Relationship	Type
comment	POST

③ Follow → (my followers, who am I following)

User	Followers
Alok	X
Dave	Y
Mike	2
X	Alok
X	Y
Y	Alok

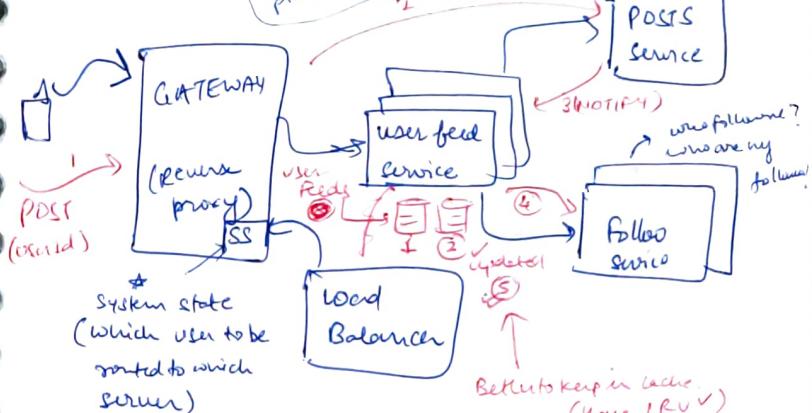
Let's say a celebrity makes a post, now we are updating cache of each user's list of posts. This is put on top of each user's cache.

But if million followers, it means many changes!

We can say avoid user can update it by given user.

## ④ News feed I

### Design



To get news feed I

get UsersFollowed by (User Id)

← returning set (User Id)

get Posts by userId (set(Post))

← returning set (Posts)

\* But for every request for news feed, if we query POST service, very heavy.

So better to PRE-COMPUTE. whenever we get POST (userId) request,

TO update the news feeder

of all followers of user id.  
(uses follow service internally)

② Notifies user feed service

## GOOGLE MAPS Algorithm:

Step Designing a location based database

- We need → measurable distance  
 (Input param) ① b/w 2 points
- ② Proximity  
 (people close to you)
- uniform assignment of locations  
 ↓  
 Scalable granularity of locations

## Cache:- [Grokking]

### Write through:

- Data is written in cache and thereby in DB. Single operation.
- Adv:- consistency maintained.
- Disadv:- Since for a single report, 2 writes happen ↴ Cache ↴ DB.  
 High latency.

### Write Back:

- Data written in ~~book~~ at cache alone, then periodically or after some event, all writes moved to database.
- Adv; very fast operation. low latency.
- Disadv: If cache fails, data lost.

## CDN:- Content Delivery Network

↑  
 eg of Cache.

→ Helps in serving static data.

## Indexing in NoSQL

→ Indexing in NoSQL is based on a partition key.

↑  
 Key (or hash result) on basis of which data is partitioned.

- \* You can create custom index-key to reduce search times. e.g. if on data we want to have all rows where country = India + gender = female ; Better to have index (Sort-key) of country + gender + City.

Now, Select \* from users where sort\_id like "INDIA + FEMALE %"

↑  
 Regex

- \* Thus we can create custom keys if required by user. NoSQL can't handle data normalized at all places.