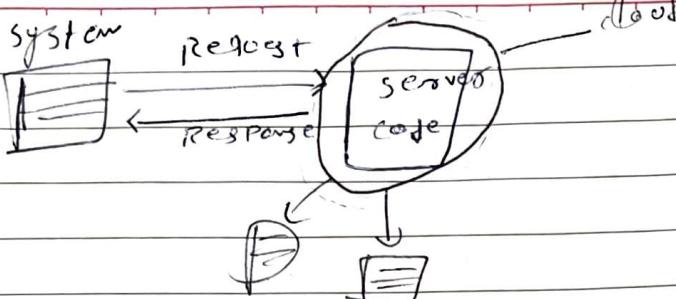


System Design

Day 1



- ① Buy Bigger Machine → Horizontal scaling
- ② Buy More Machines → vertical scaling
- (Request distributed)

Horizontal

- ① Load balancing require
② Resilient.

- ③ Network call
can require to call
one system to other
④ Data consistency issue
⑤ Scalable well
ex. users increase

vertical

- ① N/A
② single point failure.

- ③ Inner process communication

- ④ Consistent
⑤ Hardware limit

we can't
stop every
machine
when fail
call fail
50's
now
or
system
running
same
machines.

Day 2

Distributed System (High level Design)

- Recruitment
- separation of concerns
- Fault tolerance

vertical
scaling

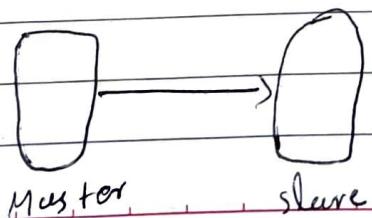
- ① Optimise process & increase throughput using same resources.

resilience
(Recovering
from failure
quickly)

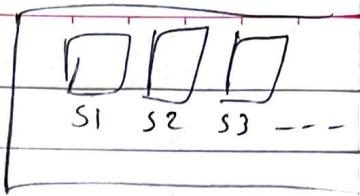
- ② Preparing beforehand at non-peak hours.

backup

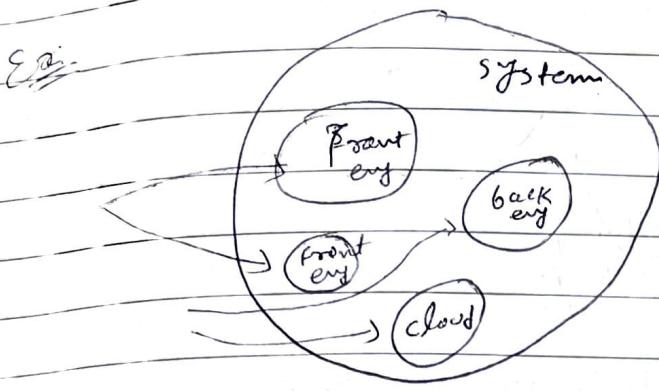
- ③ Keep backups & avoid single point of failure



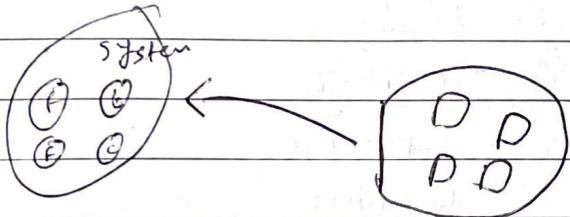
(4) hire more resources } horizontal scaling.



(5) Microservices

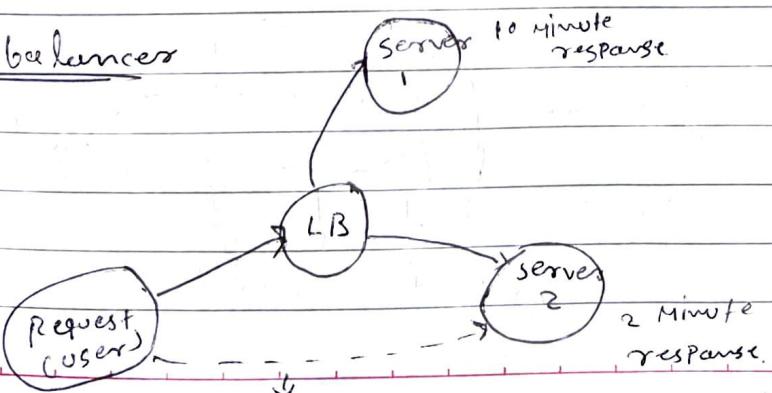


(6) Distributed system :- when there is an issue in main system, the other part of the system are ready to serve.
- It's a fault-tolerance.



For ex. Facebook has the local server everywhere, so people from India, get the connect with the Indian server, any response quickly, instead of taking time to response in connect with USA server.

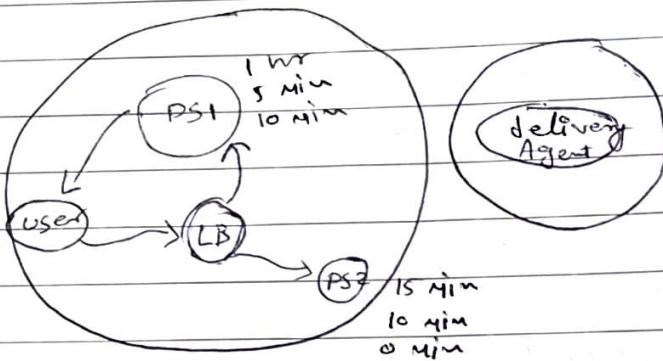
(7) Load balancers



This request goes to server 2 due to less response time.

(8) Decoupling

~~Ex~~ Pizza Shop.



→ we separate delivery & order, so we can
scrutinize, Analyze, Audit, report, etc. -

(9) Logging & Metrics

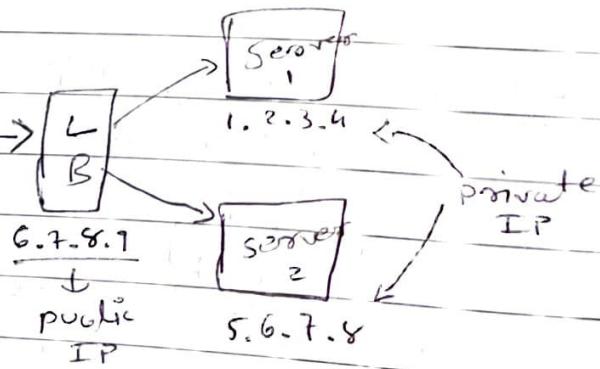
⇒ to check every system we decoupled.

- (1) Analyze
- (2) Auditing
- (3) reporting
- (4) Machine Learning

(10) Extensibility :-

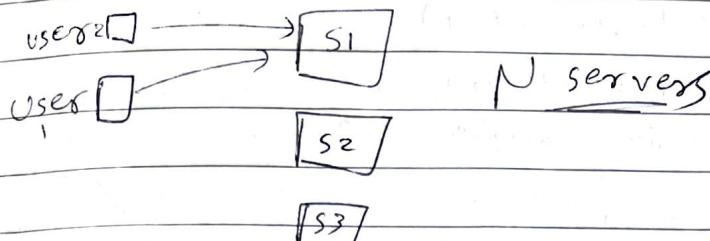
& Load balancer :-

An
- Scalability
- Availability
- Flexibility



~~Day 3~~ Load Balancing :-

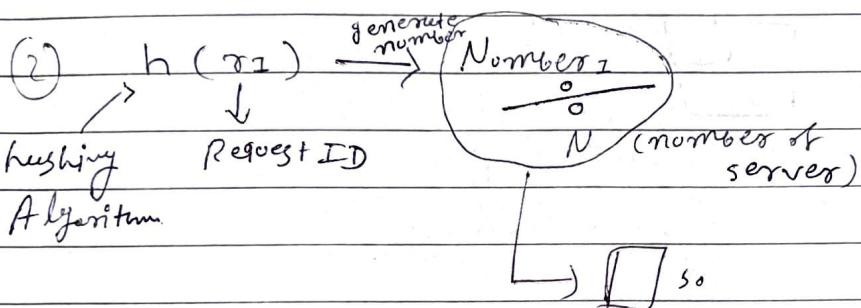
* consistent Hashing



Load balancing.

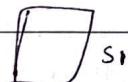
(1) Request ID \rightarrow 0 to $M-1$

↳ when user send a request it uniformly generates number from 0 to $M-1$

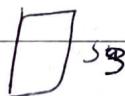


e.g. $h(10) \rightarrow 3 \text{ } \% \text{ } 4$

= 3



Request \rightarrow Server S2



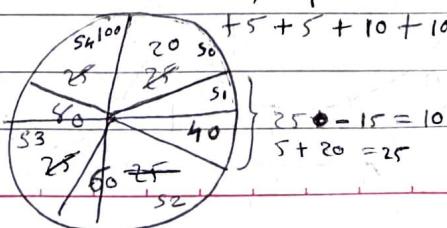
$h(20) \rightarrow 15 \text{ } \% \text{ } 4$

= 3

Now, if we want to expand the system, but it is more biggest issue why?

4 - servers \Rightarrow 25% each

5 - servers \Rightarrow 20% each



loss out

1
1

odd lost

1
1

$$+ 5 + 5 + 10 + 10 + 15 + 15 + 20 + 20 = 100$$

$$25 - 15 = 10$$

$$5 + 20 = 25$$

change
entire
workspace

It's a problem because
for each

we have $h(\text{userid})$ by user ID.

Some user request some $h(\text{userid})$
 server again & again, so server store his data from database to local cache, so time is less.

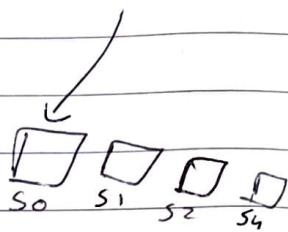
$h(\text{darsheel 2784})$

$h(\text{faisal 2784})$

$h(\text{faisal 2784})$

so on
11
11

$$20 \times 4 = 0$$



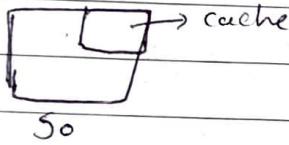
So, the same userid always points to the same server, and we set the cache in the server.

but, we add new servers in the server.

Cache and server changes.

It's an issue.

(useful cache is jump due to system change)



so

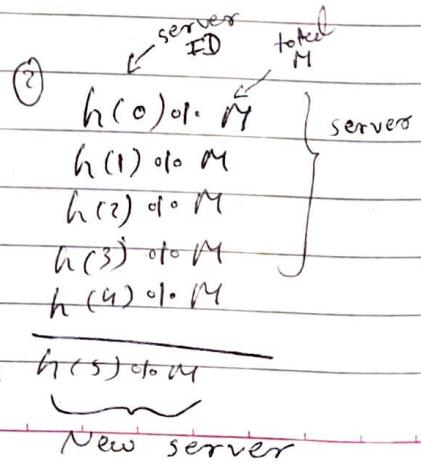
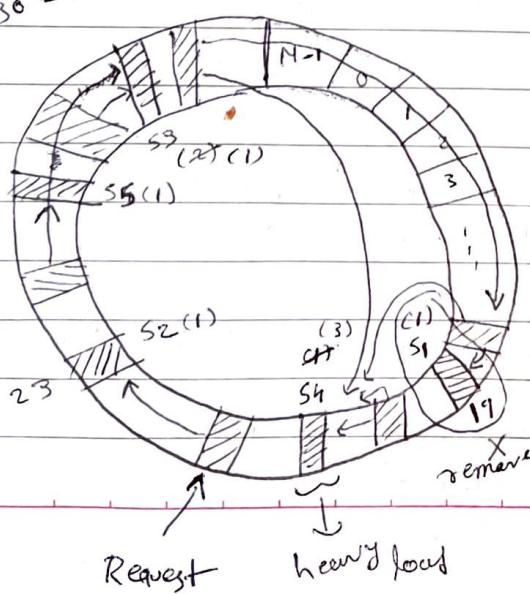
* Solution

① Transitional Methods:

① Request ID $\rightarrow h(\infty)$

$$h(0) \oplus 30 = 19$$

$$= 19 \oplus 30 = 11$$



We can do much better way.

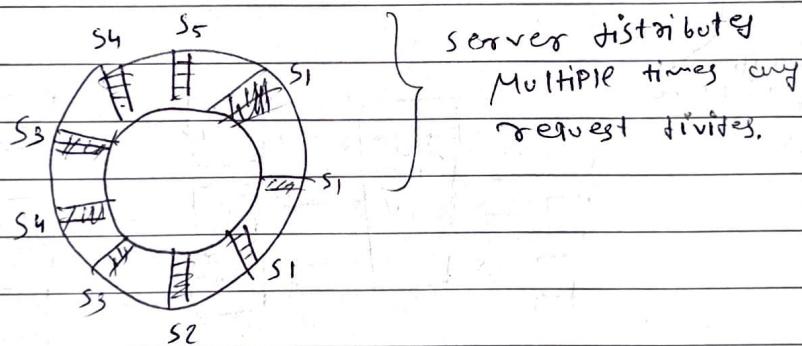
Instead of hashing only 4 servers, we hash each servers no. of server times.

$h_1(0) \cdot 1 \cdot M$	$h_2(0) \cdot 1 \cdot M$
$h_1(1) \cdot 1 \cdot M$	$h_2(1) \cdot 1 \cdot M$
$h_1(2) \cdot 1 \cdot M$	$h_2(2) \cdot 1 \cdot M$
$h_1(3) \cdot 1 \cdot M$	$h_2(3) \cdot 1 \cdot M$
$h_1(4) \cdot 1 \cdot M$	$h_2(4) \cdot 1 \cdot M$

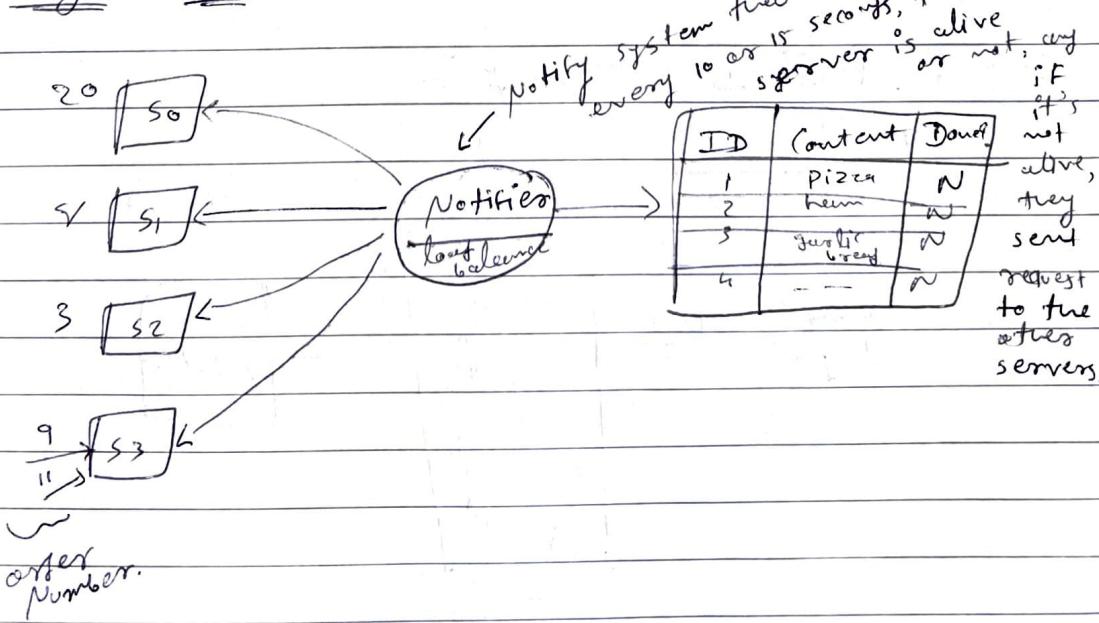
S_1

S_2

So, our



* Message Queue



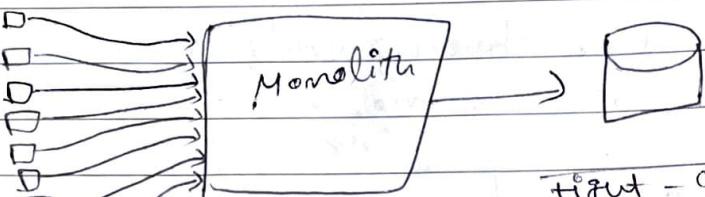
~~Ex~~ of Monoliths is we deploy our system in AWS database and live our system.

DOMS	Page No.
Date	/ /

~~Day "n"~~ Monolith vs Microservices

Adv.

- less team
- less complex
- less duplicate
- faster
- ~~this~~
- More context required
- complicated deployment
- single point of failure.

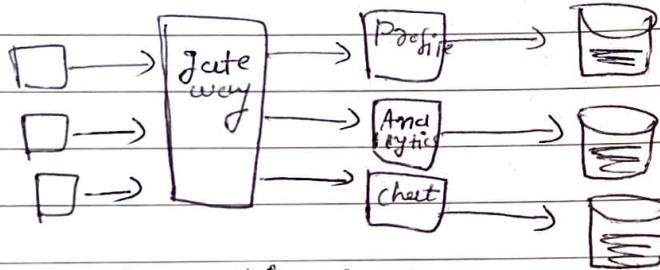
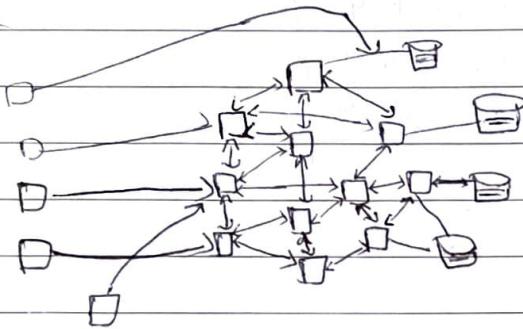


tight coupling in development.

Ex stack overflow

Adv.

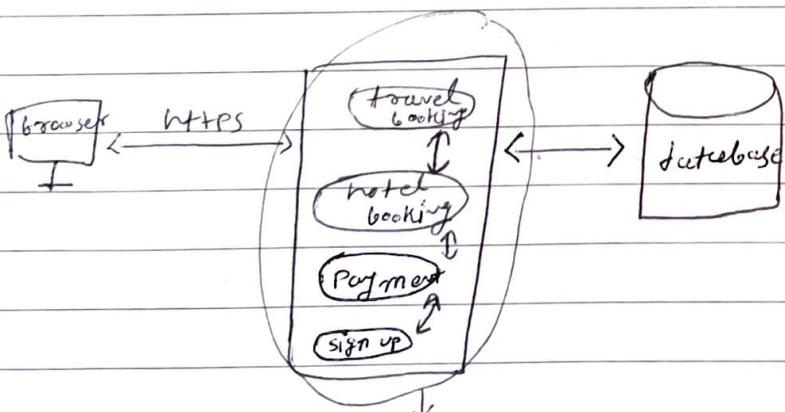
- easy to scale
- easy for new members
- working in parallel
- easy to reason about
- ~~this~~
- Need skill Architecture



Ex. Google, Microsoft, Amazon, --- etc.

~~Day "n"~~ Database sharding.

* Monolith :-



All this features is tightly coupled.

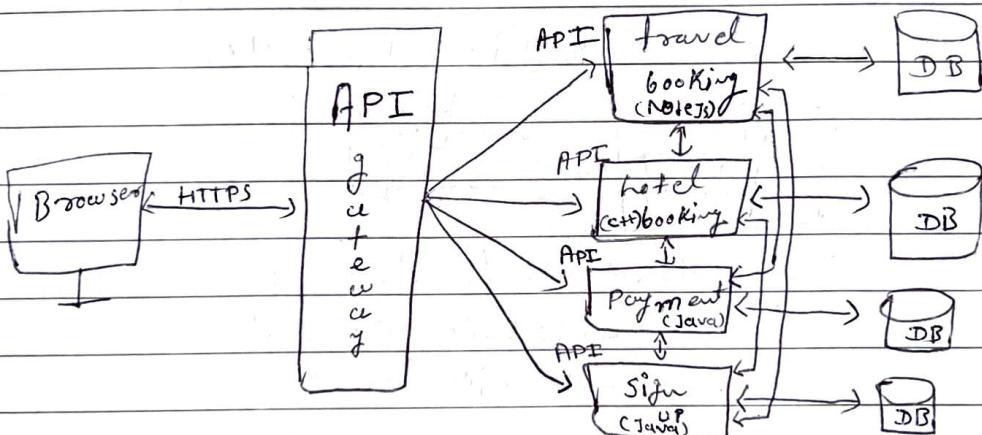
\Rightarrow The problem is that for ex we want to change the features of hotel booking any change the database, then we can't do it, we need to change everything, every features. So, it's scalability issue.

Adv

- low complexity
- easy to deploy

Dis

- Rigid
- hard to scale
- slow performance
- single point of failure
- slow continuous deployment

* MicroservicesAdv

- loosely coupled
- Agile & flexible
- Independent Development
- Independent Deployment
- Fault Isolation
- Mixed technology stack
- Granular scaling

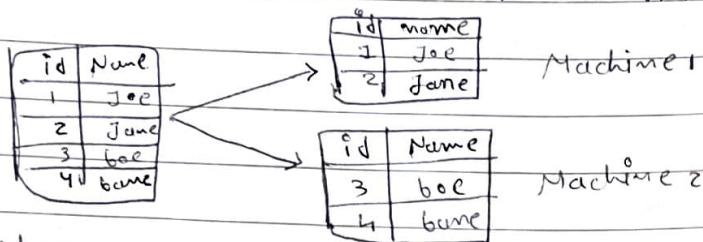
Dis

- high complexity
- consistency
- Automation
- Debugging

Day = 5

Database Sharding :-

- Splitting data into smaller chunks or shards, where each shard/chunk can be on different machine

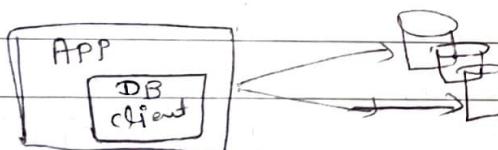


(complicated solution.
so, do if it's needed.)

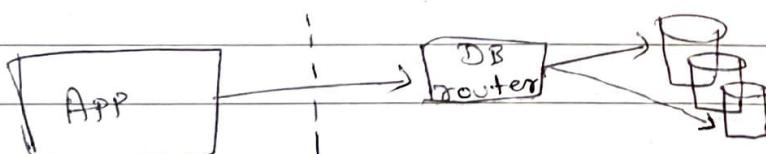
- why
 - data may be too large to fit in single machine.
 - Provides horizontal scaling of database.
- Distributed processing to speed up computation.

How is it achieved?(1) Application level sharding

- Database client picks relevant shards for reading/writing data.

(2) Database level sharding

- Router/config present in database and not with client

Sharding strategies - Hashing

- computes a hash of the sharded key field's value. Each chunk is assigned a range based on the hashed sharded key values.

- while range of sharding key may be "close" their hashed values are unlikely to be on the same chunk.

sharding strategies - Range

- divide data into ranges based on the sharded key values. Each chunk is assigned a range based on the sharded key values.
- A range of sharded keys whose values are "close" are more likely to reside on the same chunk.

Ex Order Management system

- order data of users in e-commerce site

```
Order {
    long orderId;
    long userId;
    Date date;
    Status status;
    OrderDetails details;
}
```

- sharded on basis of orderId or userId.

-

=> Database with sharding support

- Cassandra
- HBase
- MongoDB, etc...

* Issues

(1) Hotspotting → uneven distribution of data

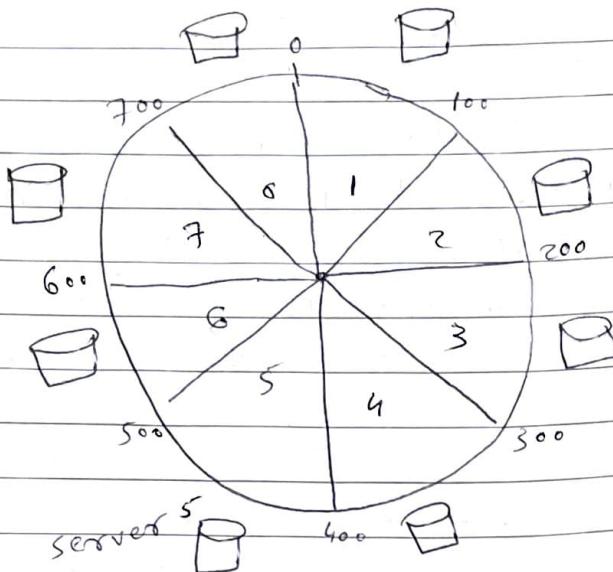
- one key have large amount of data

(2) Redistribution of data

- New machine in cluster

(3) Network Partition

- network partition between DB machines when using DB based sharding

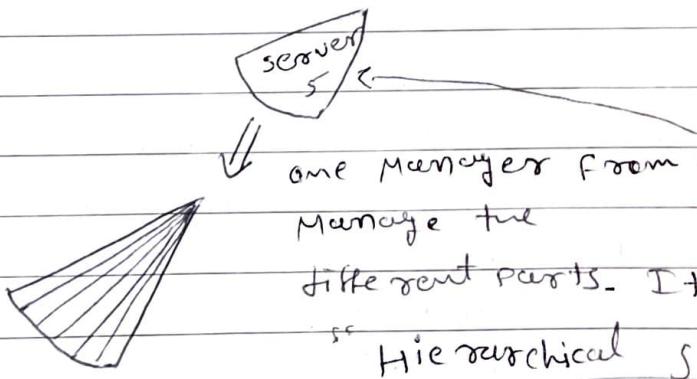


- (1) Horizontal Partitioning with one key like user ID
- (2) Vertical Partitioning with column.

Problem:-

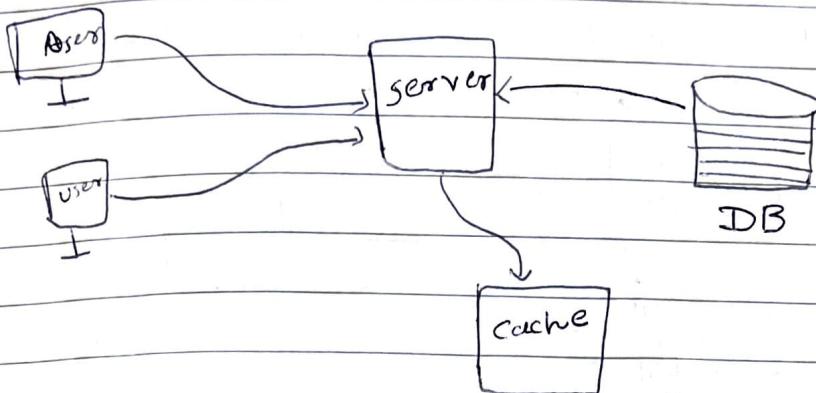
- ⇒ if we want a data from two different shards/server, then we need to pull out data and do Joining.
- ⇒ Memcached
- ⇒ Fixed number of shards.

II Solution.



- ⇒ If the any one server fails, we use Master-Slave, so slave pulling data from master, so we have slave to get the data.

Cheer is Distributed caching?



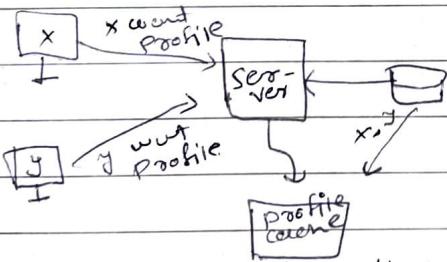
- (1) Reduce N/w calls
- (2) Avoid recomputation
- (3) Reduce DB load

* Cache Policy (Maintain cache Properly, can remove data, can update)

- (1) LRU :- remove old entry from the cache.
- (2) LFU :-
- (3) sliding window policy

+ Disadvantage.

- (1) Extra call
- (2) cache thrashing (when we have less cache)



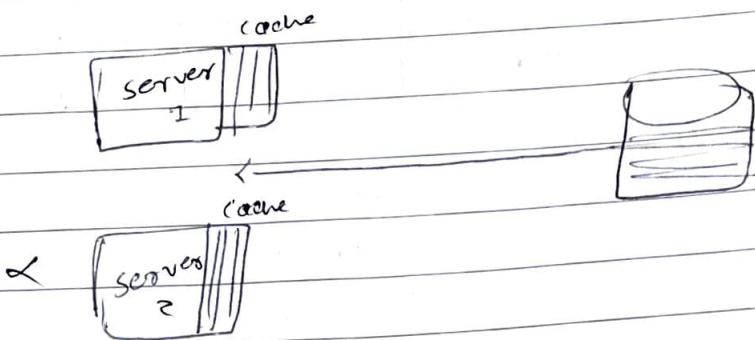
X, Y, Now again when X get the data, it's finding in cache, it's thrashing.

- (3) Data consistency.

Ex) when user update profile in database, and when they tried again to get updated profile, why it's not changing in ~~cache~~, then it gives old profile data to users.

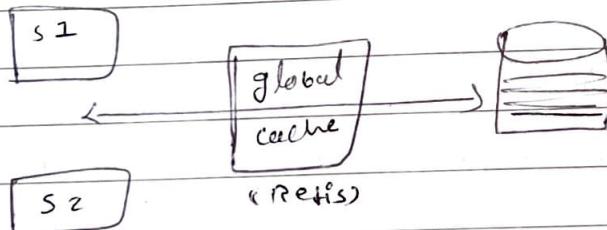
* Where we can place cache?

- (1) it's close to the database
- (2) it's close to the server



⇒ if we put the cache in memory of the server, then when server fails, cache also fails, and the cache of server is also not sync with server 2, so we can't get the all information that user wants.

↓ solution
(global cache)



* global cache

- ⇒ faster
- ⇒ Accuracy

* How consistent data in cache?

- (1) write-through
- (2) write-back

① write through



first we

change user X's
password in cache
before updating to
database.

② write back

⇒ we change update in database, any time after some time we change update in cache

problem

if we keep updating cache any database every time, it's a performance & timing issues

solution

we keep changes in cache for certain entries like 50 to 100, any then after some time in one network call we update the database.



~~Design~~

API Design

① Naming

② Parameters defining

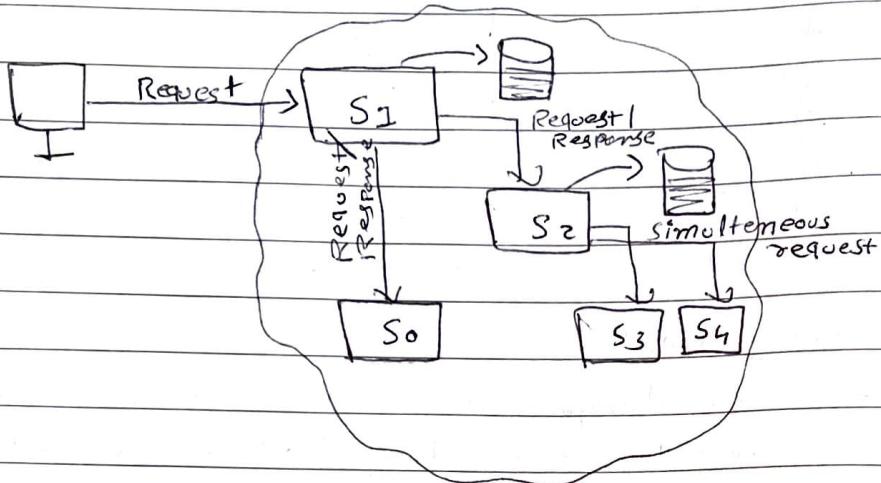
③ Response object defining

~~Ex:~~

www.xyz.com/check-Mess/getAdmin/v1

* Publisher - subscriber Model

⇒ what's the issue in microservices?



⇒ S₀ to S₄ is a services

⇒ client sends request to S₁, after S₁ has done it's changes, S₁ needs to send two messages to S₂ & S₀. any after S₂ process, it's also needs to send messages to S₃, S₄. so, it's a "strong coupling."

⇒ here S₁ needs to wait before sending response, S₁ waits for S₂, S₂ waits for S₄ response. So, it takes long time.

⇒ Now, for ex., if S₄ fails, it sends update to S₂, S₂ sends to S₁, why S₁ say fail response to client. Now, client again sends to S₁, S₁ previously made database changes, but now it's also do for 2nd time. so, it's "inconsistency of data."

* Solution:-

Publisher / subscriber Model

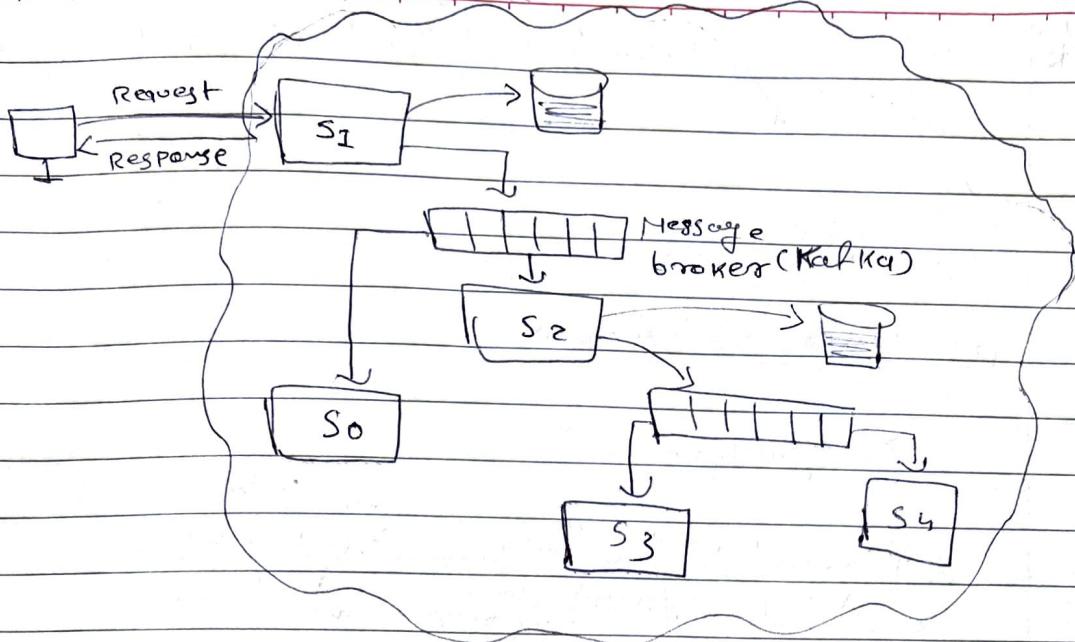
event driven services

DOMS

Page No.

Date

/ /



$\Rightarrow S_1 \Rightarrow$ Message broker (it keeps request)

MB $\Rightarrow S_2$ (if S_2 fails ~~can't ready to send response~~, then ~~MB seg is so~~, MB keep request, carry when S_2 is ready, then send them a request.)

$S_1 \Rightarrow$ client (success).

Adv :-

- Simplify interaction
- Single point of failure
- easily scalable (we can add other services easily and attached to MB.)

Dis :-

- Poor consistency (we can't use for the secure services like payment and all)
- Idempotency still required.

~~Ex~~ Twitter usage public-subscribe model.

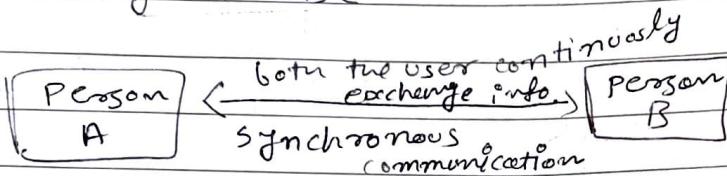
Day 9

* why do database fail? Anti patterns to avoid

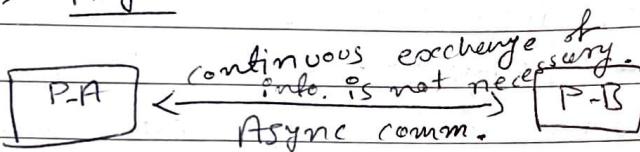
Message Queue

① Sync vs Asyne

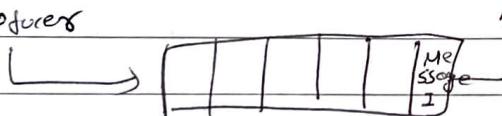
→ Synchronise



→ Asyne



Producers

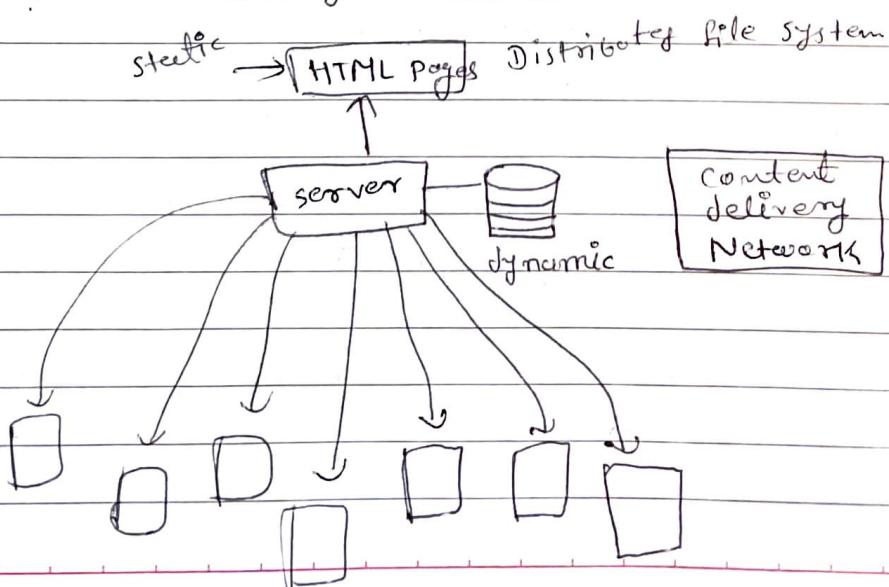


Consumers

- C₁
- C₂
- C₃
- C₄

Day 10

Content Delivery Networks



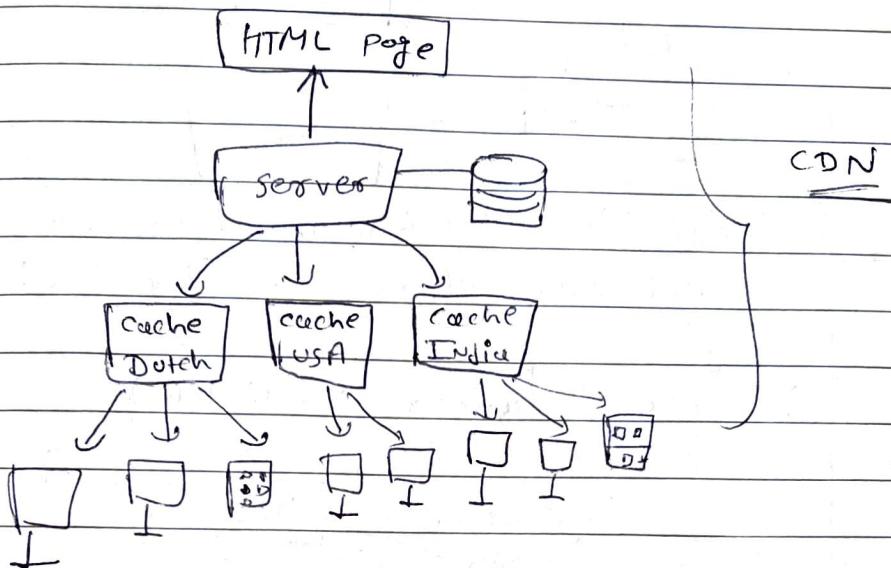
~~we need those things in our website.~~

- ① caching
- ② customized data (device, location)
- ③ fast



Solution

* solution 1 but it has some drawbacks.



difficulties

- Available in different countries
- follows regulations
- serves the latest content.

* Content delivery Net.

- hosting boxes close to users
- follow regulations
- Allow posting content in the boxes via UI

* CDN Services => S3

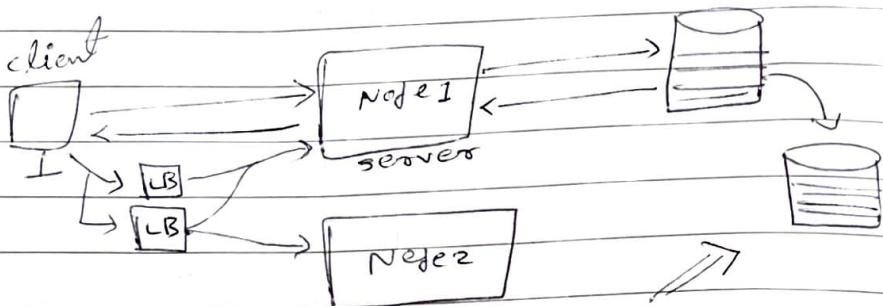
Adv of CDN

- speed
- low
- uptime
- security

Day - 11

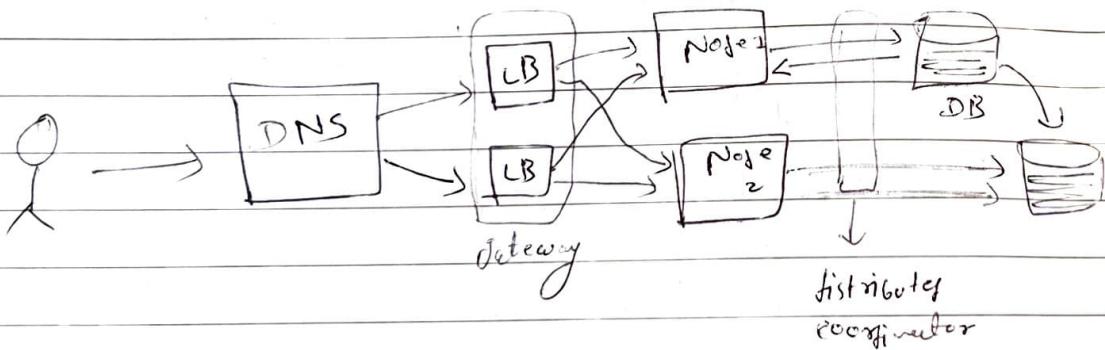
How to avoid a single point of failure in distributed system?

- Add more nodes if server fails



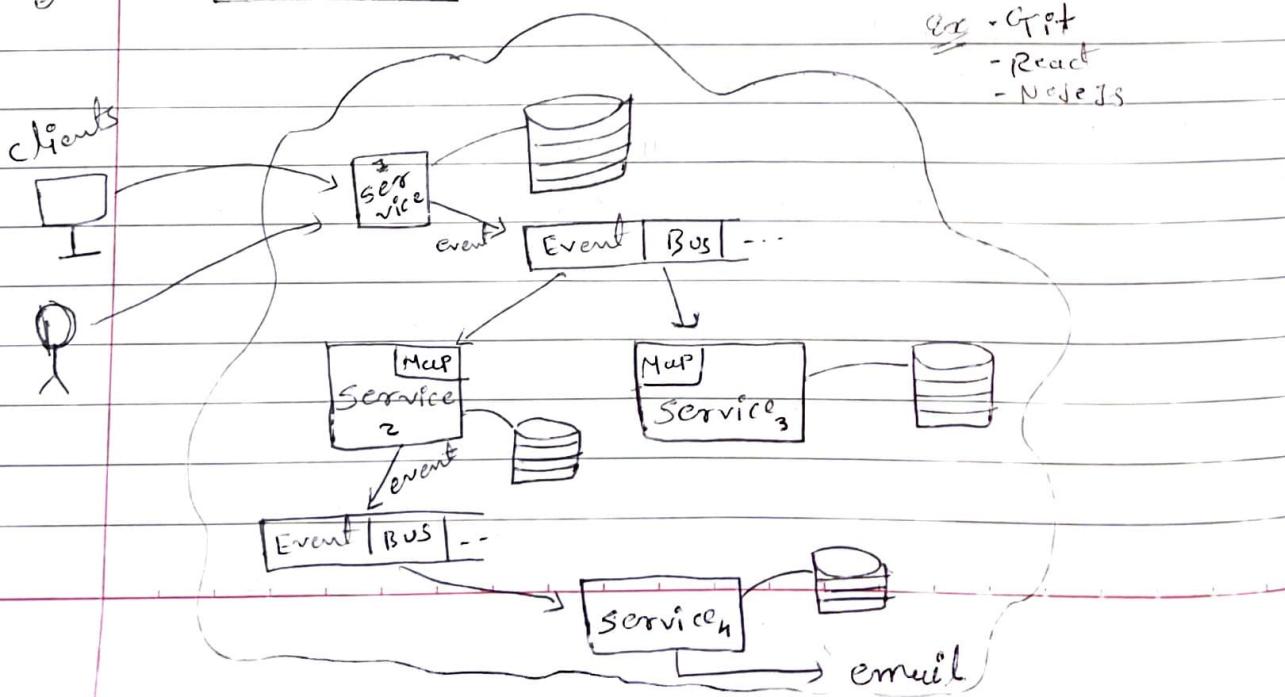
- Add other database if Main database fails.
(Master-slave Architecture)

- Load balancer (Multi LB)

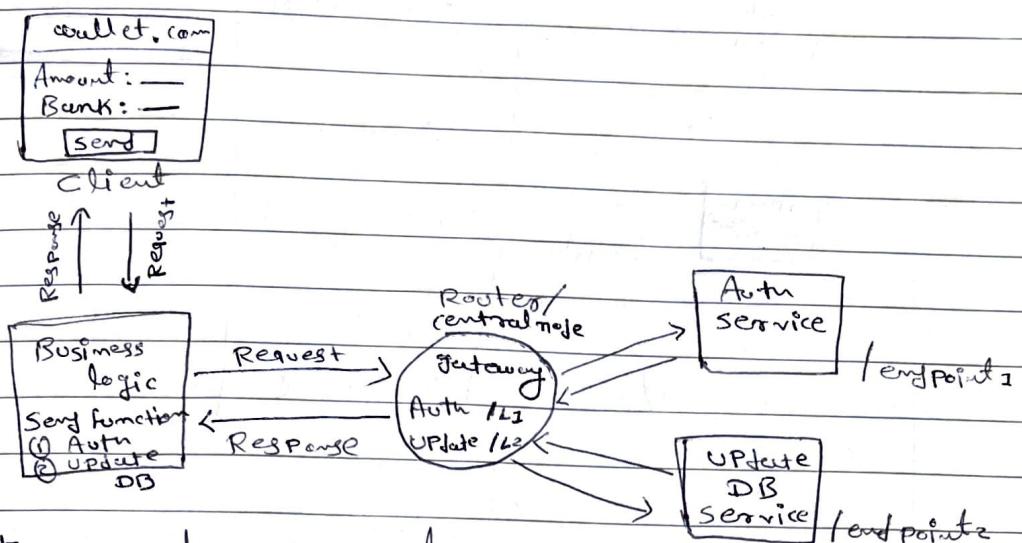


Day - 12

Event-driven Architecture



⇒ Let's understand Request driven model (Microservices)



⇒ The problem is that this is synchronous, so after completing one service, it takes & executes other service.

⇒ The issue is when we need to save the token of the previous service.

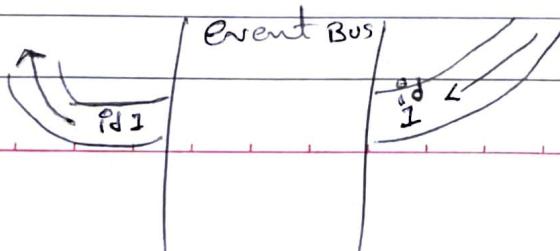
For ex in a wallet, if I send money to other user, it Authorise me, update in database, and the next service fails, so, here it update the database, but money not sent due to the other service fails, and it says the error to user it's fails. So, what about updates database?

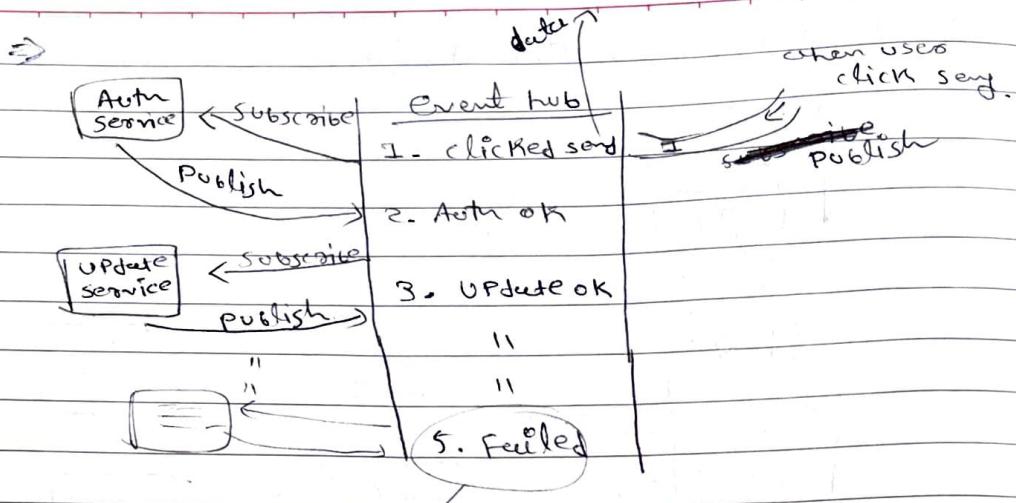
- Recovery is difficult in Microservices.

Now, let's implement EDA in microservices.

- ① Let's understand public-subscribers model.

Subscribers Publishers





Adv.

- ⇒ EDA has automated testing & recovery.
- ⇒ Availability
- ⇒ Easy roll-back
- ⇒ Replacements
- ⇒ Transactions
- ⇒ Stores Intent

dis.

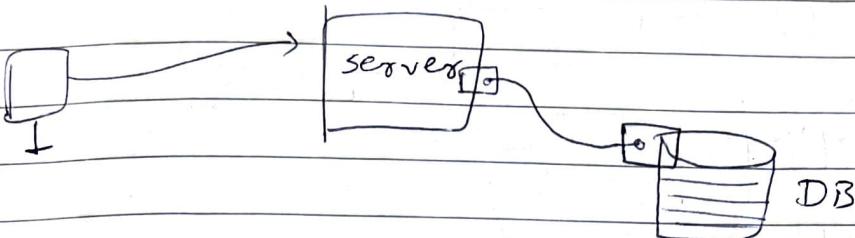
- ⇒ Consistency
- ⇒ N/A to gateways
- ⇒ lesser control

$$\begin{aligned}
 \log_2(10) &= \log_2(2 \times 5) = 10 \\
 &= \log_2(2) + \log_2(5) \\
 &= 1 + 2 \cdot 3.219 \\
 &= 3.219
 \end{aligned}$$

DOMS	Page No.
Date	/ /

Day = 10

How database scale writes:



- ① Confuse data queries into single queries.
- ② Linked list $O(1)$ write (log)
- ③ Sorted array $O(n)$ search time
- ④ we use linked list + sorted array in database to improve read fast.

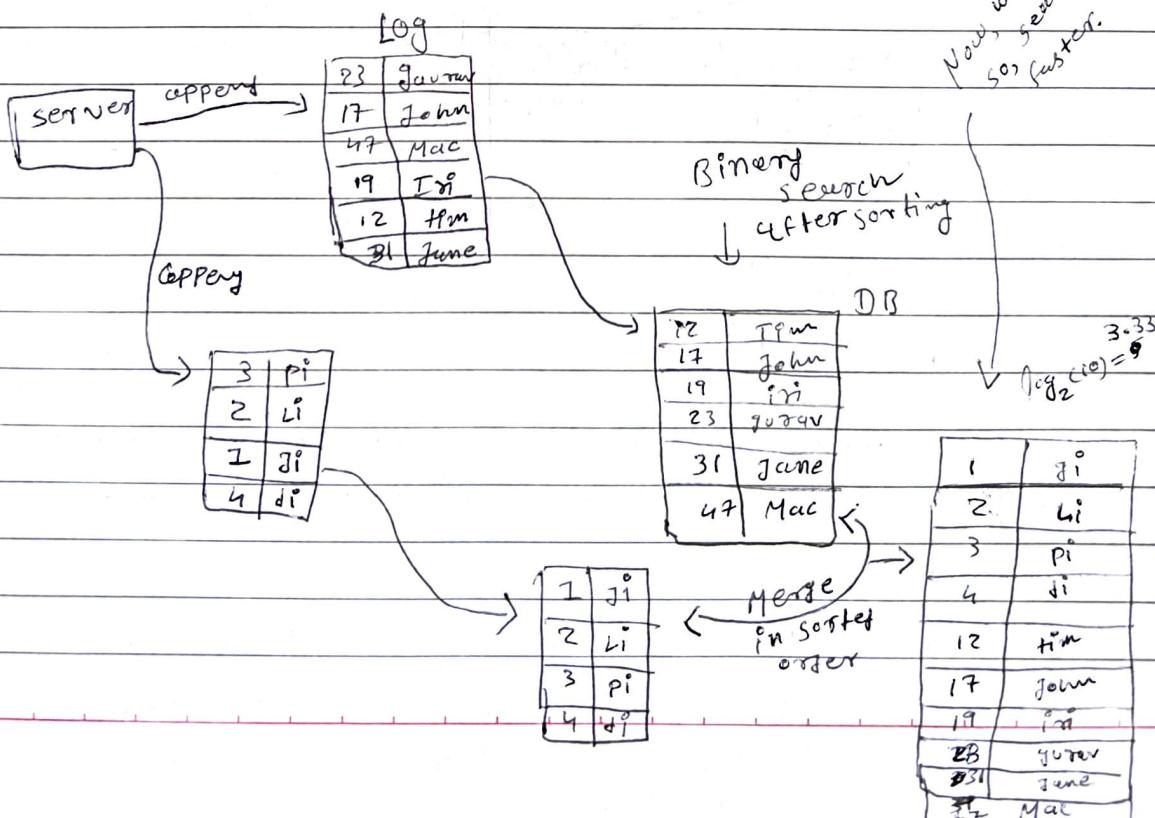
Avg.

- ① Lesser I/O operation
- ② fast writes

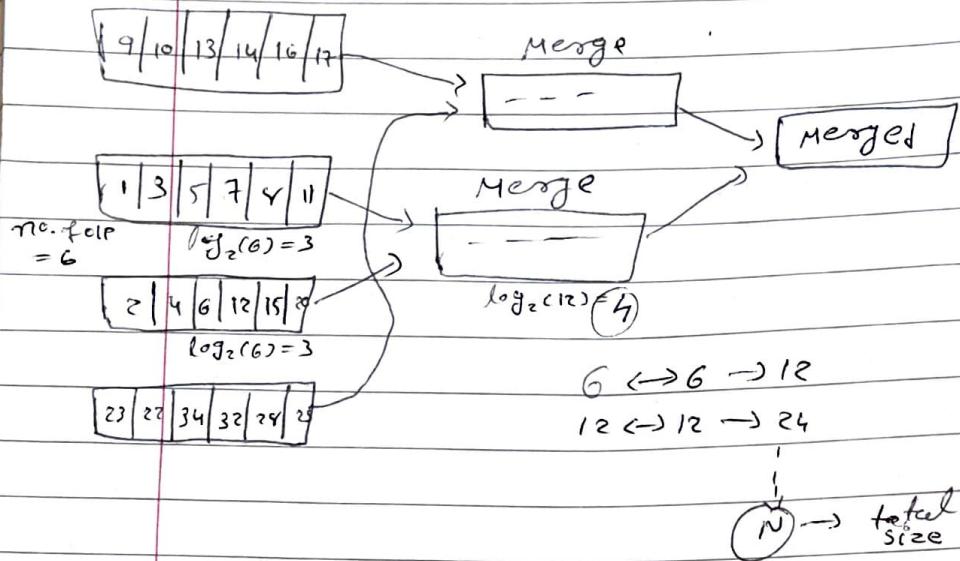
This

- ① Additional memory
- ② slow reads.

⇒ So, first we sorted the data, and then persist.



So, in the background.

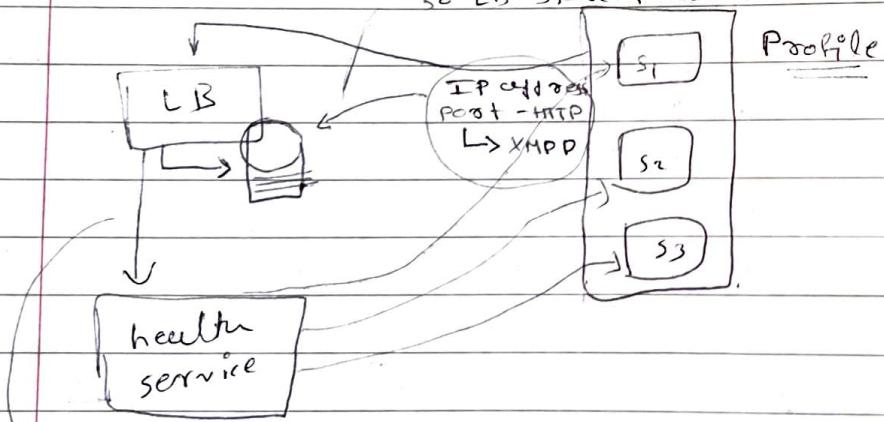


Day = 14

Service discovery & heartbeats in MicroServices

⇒ When the data pipelines are slow & we need to process the data, then we use service discovery.

This snapshot send to LB every 5-10 sec, so LB store the snapshot in its memory.



→ Based on the snapshot health service communicate with particular service & checks that is alive or not.

* Health check is important because if you have a fleet of 10 servers instance behind a LB, and one of them becomes unhealthy but still receives traffic, your service availability will drop to, at best, 90%.

+ Methods for single point failure in database

- Sharding

- master-slave (distributed consensus)

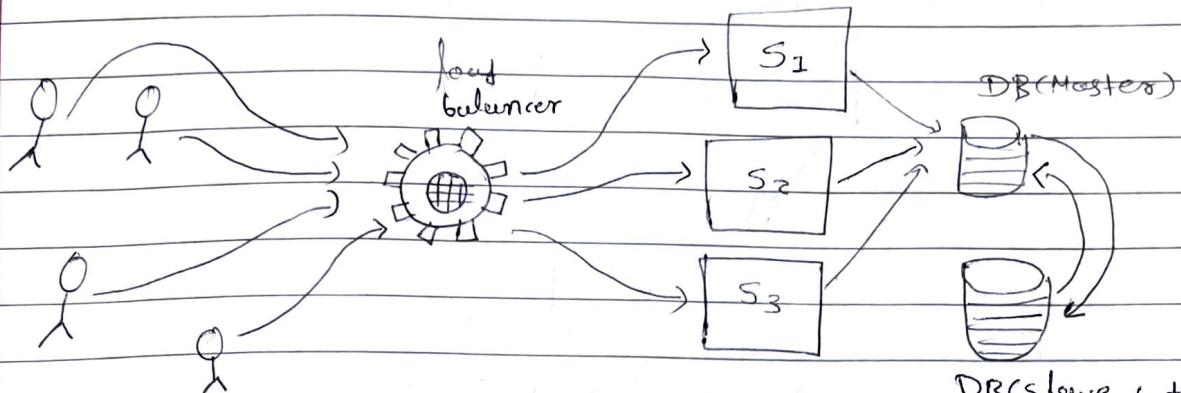
DOMS

Page No.

Date

/ /

* Distributed Consensus & Data Replication strategies on the server.



DB(Master)

DB(Slave, but work as Master)

⇒ Slave can copy data from master in two ways:-

- (1) synchronously
- (2) asynchronously.

(ii) disc

If disc

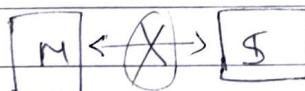
- taking time cuz load on master database

- because it's async so, master's transaction may not reflects instantly if it's fails.

⇒ Now, when master DB fails, we have slave DB to get the request from the server, & perform operations.

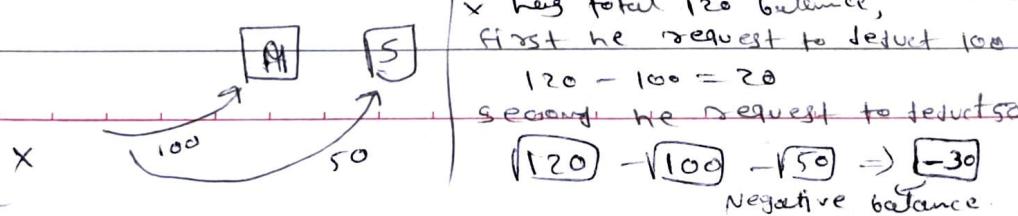
⇒ As slave work as a master, so, it also takes write operations, and updated to the original master.

⇒ This is going well, but the problem is that what about when the link or the connection between the master & slave fails.



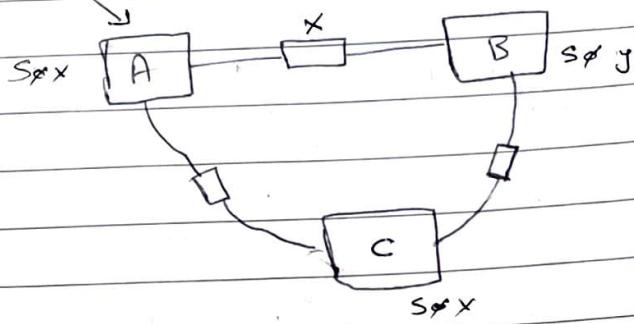
⇒ They both get assume, and take control that I'm a master, and work as a master.

This is called
split brain problem



+ split - Brain

A gets write request, we assume initially state is So, when A get write request it's state is Sx



⇒ when A gets request, it propagate to C, so both's status is Sx.

⇒ when B gets write request, its status is Sy, and B tried to propagate request to the C, but it's not happen because C ask to B about it's (B's) previous state, and it was So, so, C tells your status is not like me, so, you (B) need to update your state, then you can do the transaction.

⇒ So, now B rollback it's transaction from Sy to So, and sync up with Sx from C. So, new transaction of B is not happened.

⇒ Now, B can get new request, and perform the operation to changing it's state.

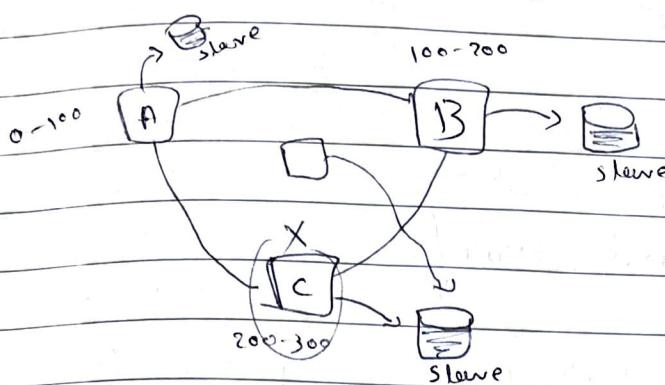
This is called "Distributed Consensus".

Multiple node agree on a particular value.

Cats

→ 2PC, 3PC, MVCC, Saga
PostgreSQL

⇒ we can also use sharding with Master-slave:



⇒ central co-ordinator redirect request to c's slave when c sharding is full.

Day = 16 How to avoid cascading failures in a distributed system

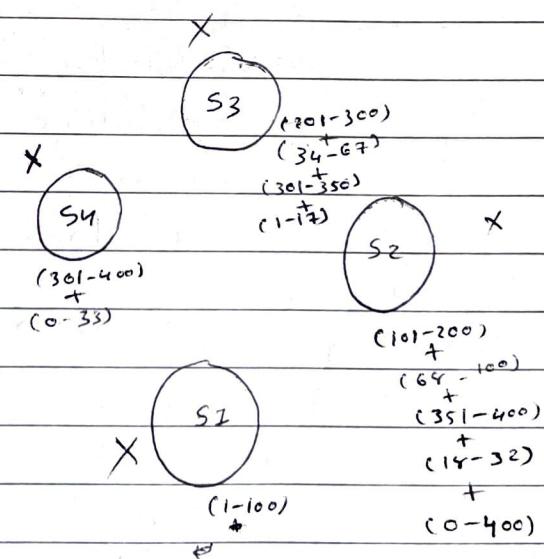
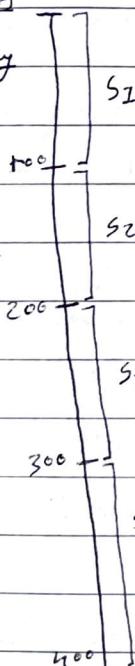
① cascading failure ↴

② viral / black friday

③ Job scheduling
(email ex.)

④ popular post

⑤ predictable load
increase



⇒ when S1 crash, it's requests divide into S4, S3, S2

⇒ when S4 crash, it's request divide into S3, S2

⇒ if S3 "", all the requests handle by S2, and it's a high chance that S2 will also fail/crash.

This is called "cascading failure problem"

* Solution

- (1) Re-scale
- (2) Auto-scaling
- (3) rate limiting
- (4) batch processing
- (5) Approximate statistics
- (6) caching (key-value pair)
- (7) Graduate deployment
- (8) coupling

~~Ex 1~~ if you want send mail to 13 users, then you sent by batch.

$$[1000] - 1 \text{ min}$$

$$[10,000] - 2 \text{ min}$$

"

"

"

$$[10,000] = 10000 \text{ min.}$$

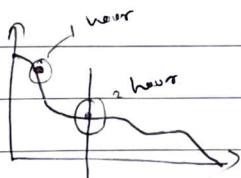
~~Ex 2~~ youtube use fitter for shows the views count on youtube videos

1 hour \rightarrow 1000 likes
views comments

2 hour \rightarrow 1000×1.5

$$1500 \rightarrow 1700$$

Approx. Reality.



* Capacity Planning (rules & tips)

→ capacity estimation → no. of transactions
 Amount of data to be stored.
 Availability

→ Metric system → Millions, billions, trillions

$$1 \text{ M} = 10^6$$

$$1 \text{ B} = 10^9$$

$$1 \text{ T} = 10^{12}$$

→ Storage capacity

$$1 \text{ B} = 8 \text{ bits}$$

$$1 \text{ KB} = 1024 \text{ B}$$

$$1 \text{ MB} = 1024 \text{ KB}$$

$$1 \text{ GB} = 1024 \text{ MB}$$

⇒ Memorizing Exercise

- 1 Million / day = 12 / sec, 720 / min, 4200 / hour.
- 10^{12} / day = 120 / sec

⇒ Rounding off



~~5440~~

1000

~~960~~

1000

256

256

256

256

64

64

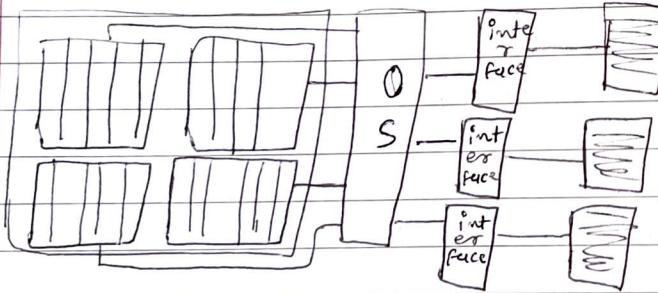
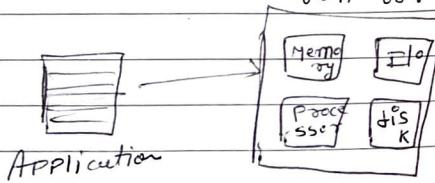
⇒ Latency numbers every programmer should know

⇒ capacity estimation of cluster?

- no. of reads & writes
- no. of transactions
- huge data transfer
- New bandwidth
- no. of requests a server has to serve
- no. of servers needed
- how a LB can handle these many requests

* Container & virtualization in cloud computing

we need these resources to run our application.



⇒ we have multiple programs which are going to be taking up slices of the resources that OS can provide, and the remaining is going to be unused.

⇒ The problem is that when the memory runs out with one program, then every one else has the issue to perform their operation.

⇓ Solution

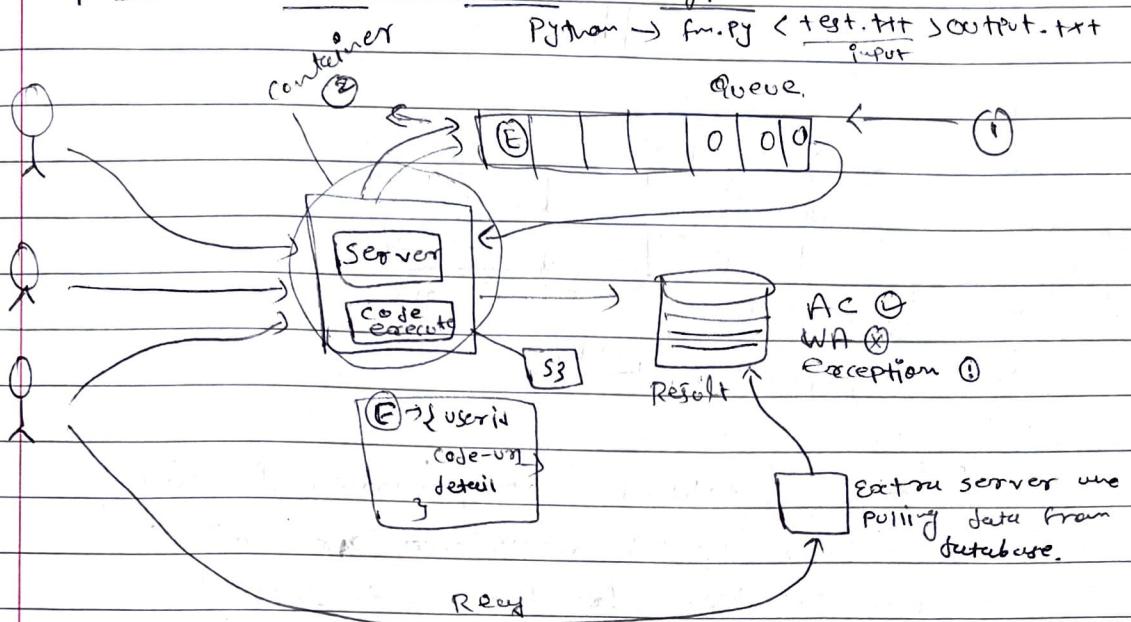
- ① VM
- ② container

- database \Rightarrow NoSQL, SQL (dynamoDB, cassandra, etc.)
- heartbeats \Rightarrow zookeeper
- LB \Rightarrow ELB (elastic load balancer)
- Message Queue \Rightarrow Rabbit MQ, Kafka
- CDN \Rightarrow Akamai, open connects

DOMS	Page No.
Date	/ /

- * to write the data in real fast like tiktok videos, we need to use Amazon S3.

* Remote code execution engine



we need only two things

- ① Job scheduler
- ② container

* Data consistency & tradeoffs in distributed system

① what is consistency?

if we have two servers then,

Any user's data should

be same across both servers.



User DB
(server 1)

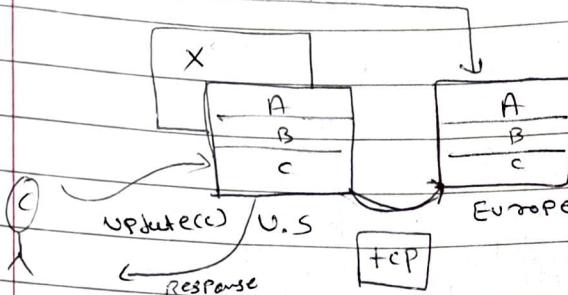


User DB
(server 2)

* Issues of single server!

- single point of failure
- cost of vertical scaling is high (impossible)
- high latency

↓ Solution



⇒ Copy of data in every server of every region. So, when US server fails, user can get the data from the Europe's server.

⇒ if user c want to update, there is some mechanism between all the servers to update c in their server.

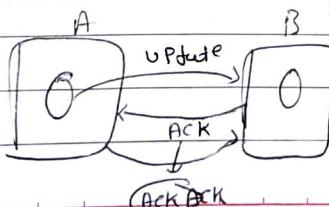
⇒ Data available everywhere, so latency will low.

⇒ we use TCP protocol for updating data from 1 server to others.

$C \rightarrow C++ \rightarrow US \rightarrow Europe$

* What happens if the message transmission fails?

- Two general problem



"Eventual consistency is better than absolute consistency"

DOMS

Page No.

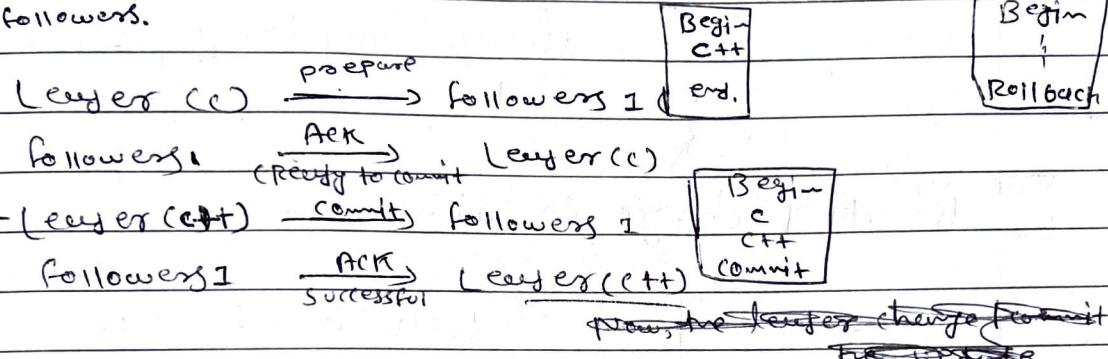
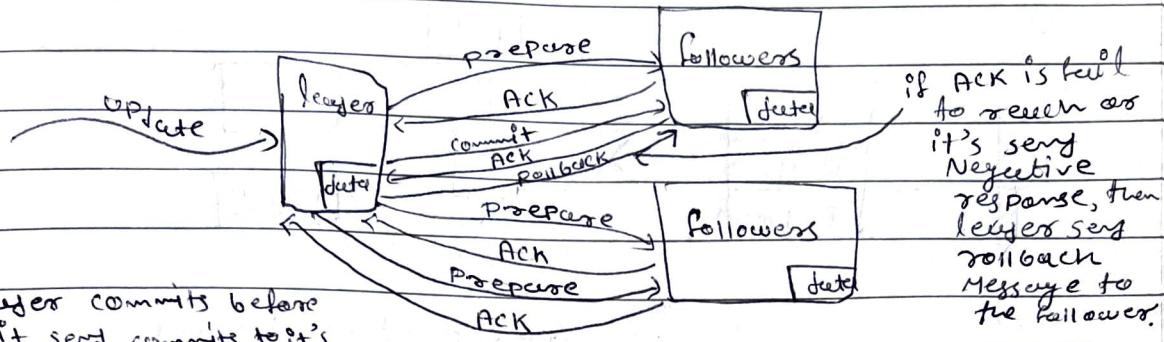
Date

/ /

⇒ Solution to the consistent state in distributed system?

* 2 Phase commit. (Absolute consistency)

- 1 Phase (Prepare)
- Retry (commits)



⇒ if the rollback happens

in Followers 1

**Begin
C
C++
Rollback**

then it should ACK

Rollback statement to Leader & Leader also Rollback. So, both Leader & Followers 1 get C.

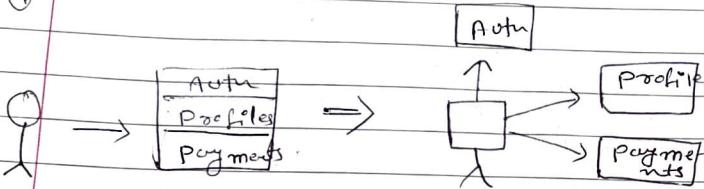
Only problem is that,

when L send commit, but it doesn't go through to F1, then we block that row in L, F1, F2 (all server's data), why it's not availability, so then user tries to access C, it will not available due to block.

Monoliths to Microservices

- ⇒ small team ⇒ Monolith
- ⇒ large team ⇒ Microservices

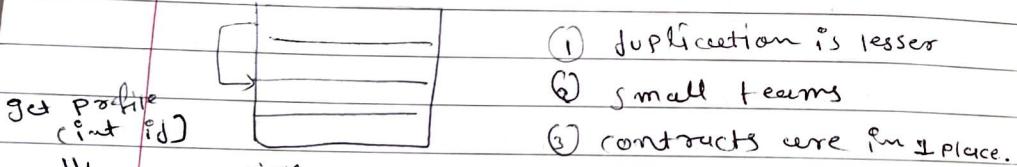
(1)



- (1) separation of codes
- (2) Engineering coding is easier
- (3) deployment is easy

(2)

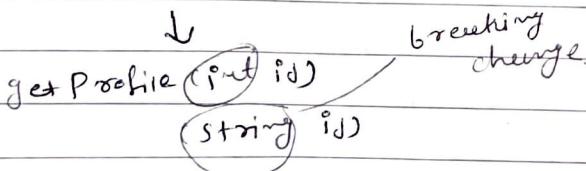
Breaking change in microservices



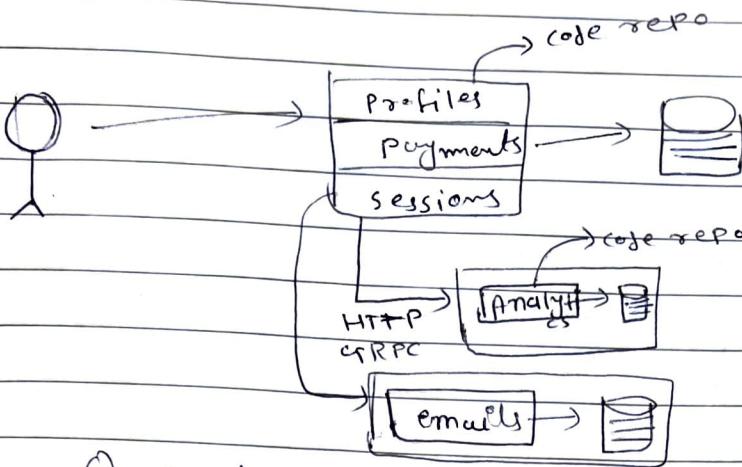
!!
Only one function
easy to handle
changes from
module A to B

Microservice

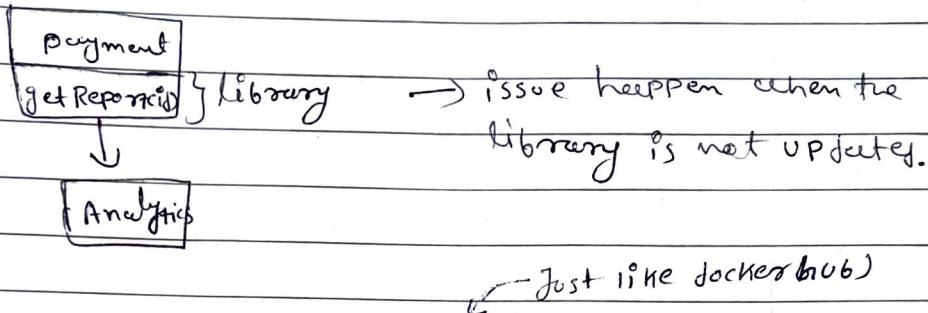
`getProfile (int id)`



+ steps for Migration



① Contract



② Router (LB or service Registry)

③ Simplify Deployment (Autodeployments)

- | | |
|--------------------------|-------------------------|
| (1) SCP repo - dev cloud | (1) CI CD pipelines |
| (2) ssh box | (2) containers |
| (3) run repo | (3) Service Deployments |

~~ex~~ Jenkins

④ Communication between services can be varied.

⑤ Logging