

= System Design

(Beginner's Version).



Design Requirement

Availability =

Uptime

Uptime + downtime

- * Measured in term of 9's, like
 $2(9's) \Rightarrow 99\% \text{ availability}$
 $3(9's) \Rightarrow 99.9\% \text{ ...}$

99.999% is considered really good, because after it, improvement is of about .5 mins.

- * generally defines Service Level Objectives (SLO)

SLA / SLO
(Expected Availability) / (Actual)

Reliability :- * $P(\text{system won't fail})$

- * Can't be calculated
- * like single server failing is more probable than double servers.

* Horizontal scaling is more beneficial when compared to vertical scaling

fault tolerance :-

If one server stops to work,
will our system continues to serve requests or not?

if a part of the system goes down, will our system continues to function or not?
or

Redundancy :-

Something that is not entirely necessary.

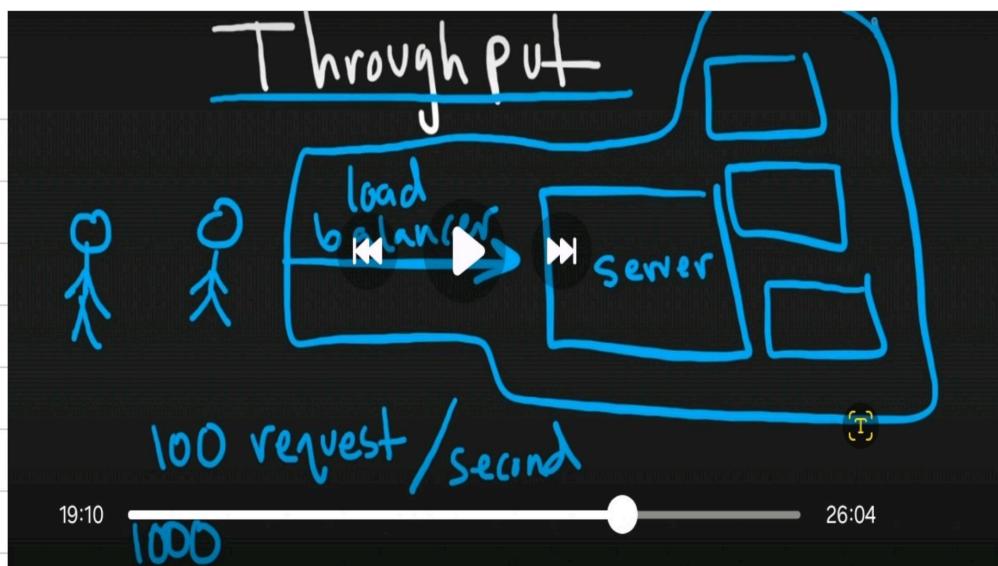
Suppose :- one server is enough to run our code but we have added another server so in that case the second server is redundant.

But it is generally good to have, because in case of a fault, we will have multiple copies.

Through put :-

Number of concurrent requests our system is able to handle per second.

Generally, we do horizontal scaling to increase the throughput, but the downside is, we require to balance it using some load balancer which turns to be complicated



But horizontal scaling becomes even more challenging in case of a database because then some data will be at one place, other at another.

QPS (Queries/second) :- Number of queries our database can handle.

Bytes/second :- Used mainly in case of a data pipeline

Latency

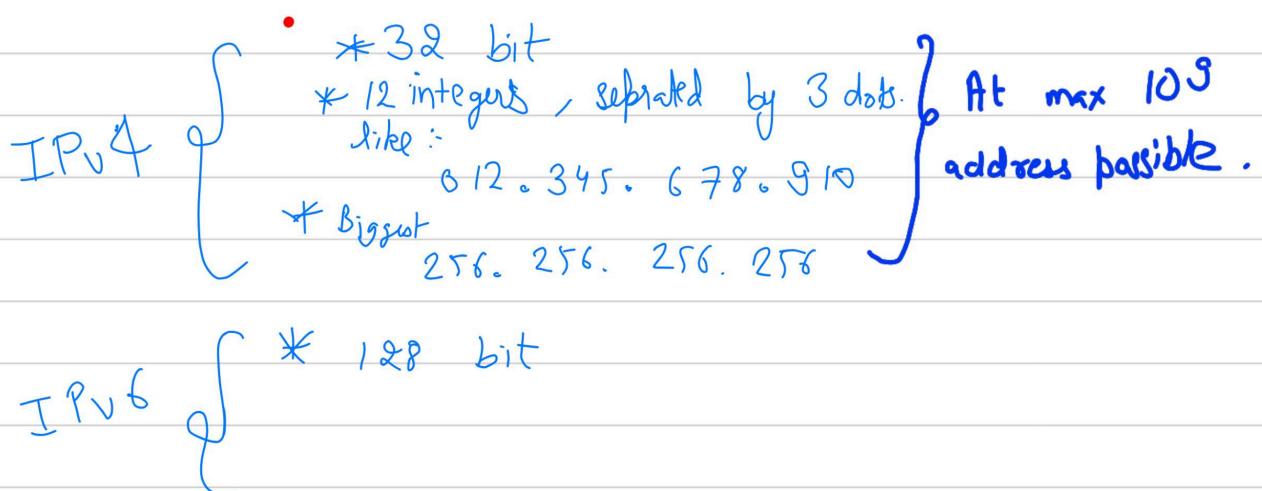
:- End to end :- Total time taken to complete from user to server and then may be back to the user

~~much better~~

To reduce Latency, generally the servers are set at the different parts of the world or Using some CDN.

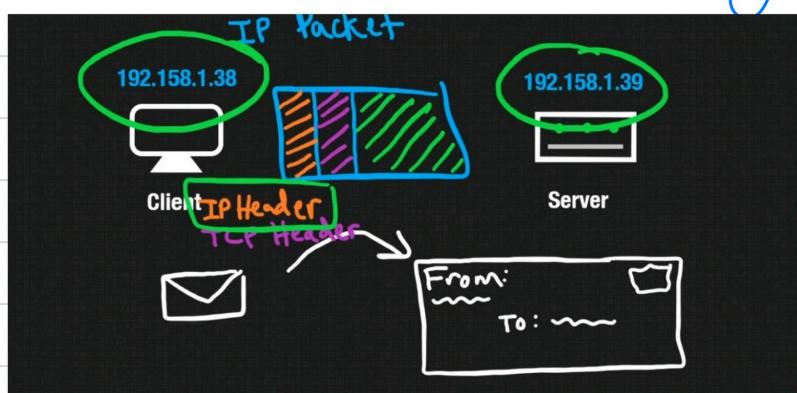
Network Basics

IP - Address :- Address to uniquely identify machines.



data is sent in form of packets where:-
it contains :-

- (i) data
 - (ii) sender's IP
 - (iii) Receiver's IP
- } Meta-data / Header



Public IP address :- for a server to be publicly accessible, it should have a public IP.

Private IP

Port \Rightarrow single application on each port.

TCP & UDP

TCP

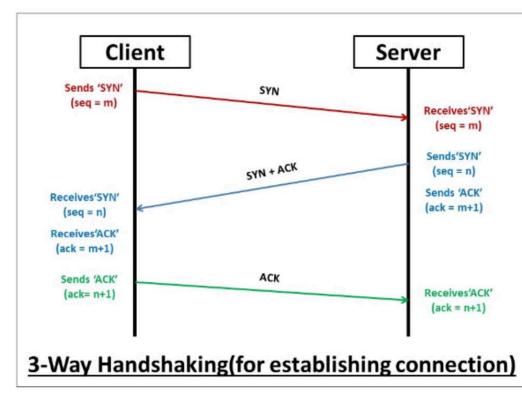
* Order is maintained

* Reliable

Retransmission of lost package. If package not arrived, then re-send it

* 2 way
3 way handshake

* Slower



UDP

* Not maintained

* Non-reliable.

* Nothing is sent

* One way

* faster

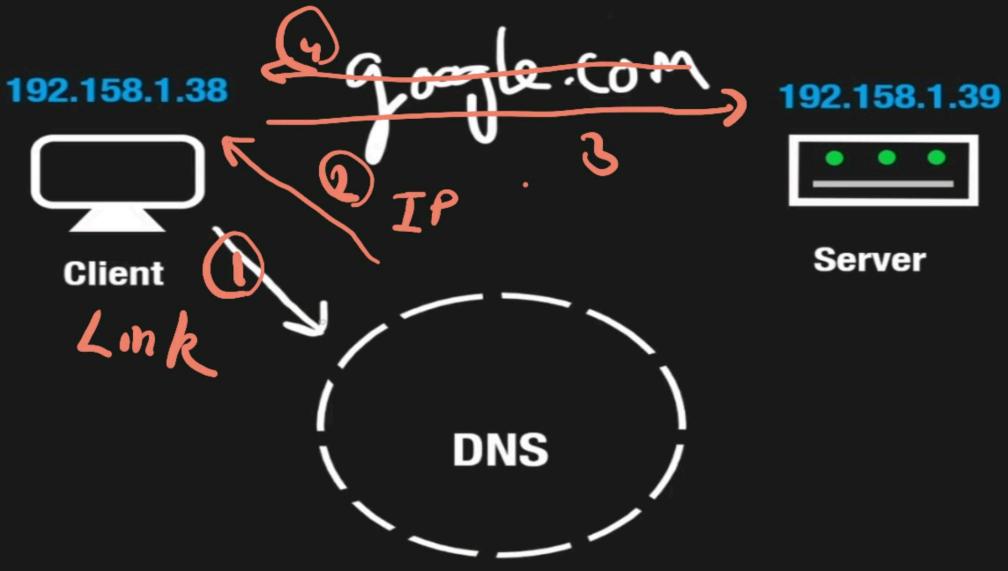
* Live streaming, we might loose one or two video frames.

* Also, used in DNS.

DNS

DNS

Domain Name System



- 1 => Link enter
- 2 => Ask DNS for the IP
- 3 => Go to IP
- 4 => Get the response

I CANN => Internet Corp. for assigned names and numbers.

* Owner of all domains in the world.

DNS records store the info. about the domain names.

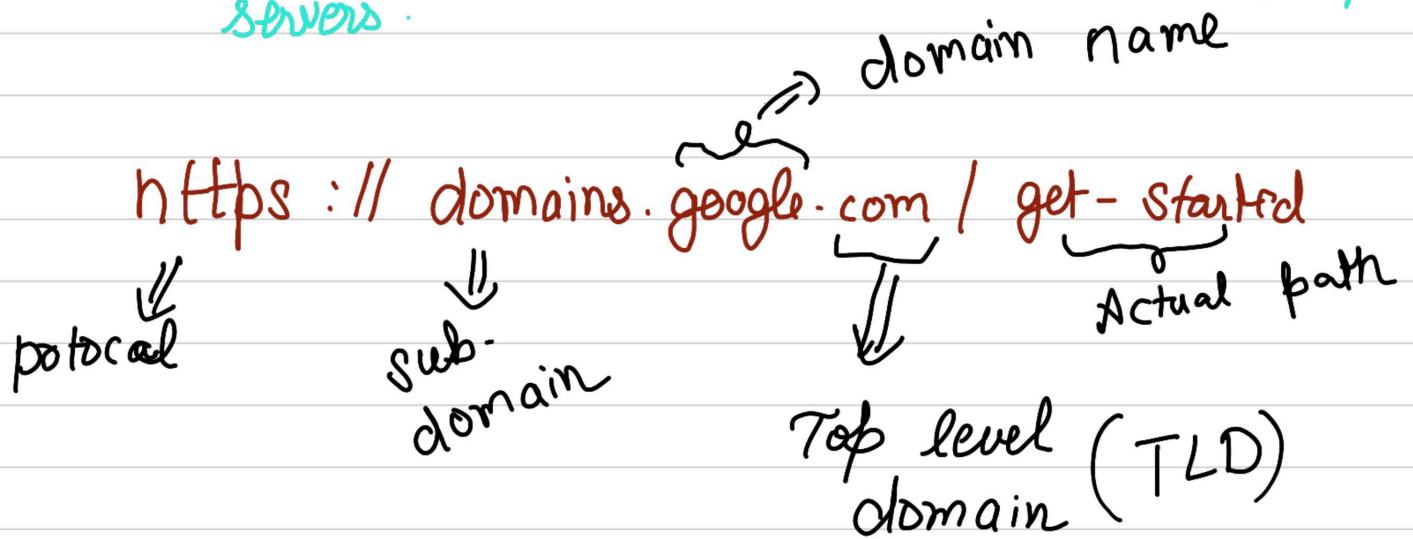
The IP address can be both static or dynamic.

But in general, IP addresses do not change very often, therefore they are cached in the user's computer.

Facts

firewall on the public server allows the req. to reach and ask.

but do not allow in case of private servers.



www is also a sub-domain

Application

Protocols

Client may be a server.

RPC

- :- Remote procedure call
- * When we req. to move data b/w machines which data we are req. to move. So, this is done via an RPC.
- * Like a fn call.
- * Think, the code is on the server, we are req. to execute it on the server. So, we call it through an RPC.
- * Since, loosely defⁿ, most of the network calls are considered as RPC.

HTTP

- * Request response protocol
- * Built on the TCP.

TCP
UDP

Y discussed Earlier.

SMTP, FTP, Web-socket, SSH.

Various Methods :-

Get

- ⇒ Return some data
- ⇒ End point
- ⇒ Retrieve Info.
- ⇒ Do not have body. So include info in the Header. or the url.
- ⇒ More cache prone.

Post

- ⇒ Do something like creating a user account.
- ⇒ End point
- ⇒ Creating an info.

Put

⇒ update

Delete

⇒ delete something
⇒ Idempotent.

Status Codes:

Informational responses (100 – 199)

Successful responses (200 – 299)

Redirection messages (300 – 399)

Client error responses (400 – 499)

Server error responses (500 – 599)

SSL | TLS

- ⇒ part of HTTPS
- ⇒ encrypted
- ↓
 - Transport Layer Security
- ↓
 - * Secure Socket

Layer
* Old * New

HTTP is not really secure in case of
an MITM (Man in the middle) attack.

Limitations of HTTP:-

* Can't be used for live chat. Only way to make it possible through HTTP is polling where we continuously send a req after a fixed interval.

*

Web-socket

* for chatting.

* Steps :-

(i) Set up HTTP

(ii) Upgrade to web socket using

the 101 Ack. from server.

- * Both didn't
- * Made obsolete after the HTTP v.2 which uses Streaming method. but are still used :)

API :- Application Programming Interface

Rest
Graph QL
gRPC

API Paradigms

Most Commonly Used.

ReST

- * Representational State Transfer
- * Loose restriction which are applied

TO MIRROR.

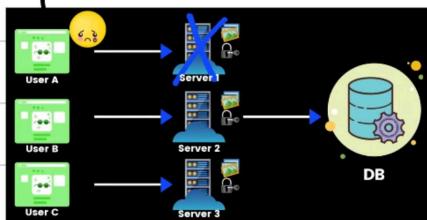
- * Scalable
- * JSON (data format)
- * Commonly used

* Stateless

- * Closely aligns to HTTP

Stateful (data)

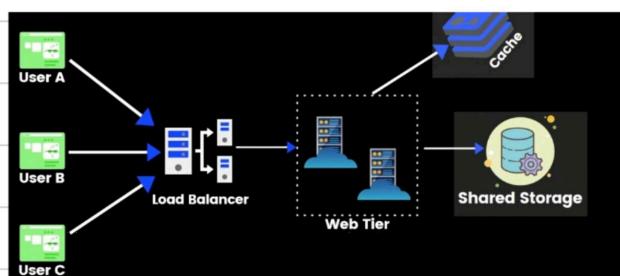
- * State of user is maintained in the local server, but not in the D.B. so now even if one server goes down, the user will lose all its session info.



- * Availability issues if one server crashes.
- * Scalability. (neither horizontally)

Stateless

- * State maintained in the shared storage or cache / DB.



Graph QL

- * Built on HTTP
- * Uses only 'post' request in where body, the query is included.
- * Query tells which fields do we want.
- * Solves problems like overfetching, underfetching.
- * Single Endpoint
 - Query
 - Mutation
- * Can't cache Graph QL

```
1 v {  
2 v   launchesPast(limit: 10) {  
3 v     mission_name  
4 v     launch_date_local  
5 v     launch_site {  
6 v       site_name_long  
7 v     }  
8 v     links {  
9 v       article_link  
10 v      video_link  
11 v    }  
12 v    ships {  
13 v      name  
14 v      home_port  
15 v      image  
16 v    }  
17 v    rocket {  
18 v      rocket_name  
19 v      first_stage {  
20 v        cores {  
21 v          flight  
22 v          core {  
23 v            reuse_count  
24 v            status  
25 v          }  
26 v        }  
27 v      }  
28 v      second_stage {  
29 v        payloads {  
30 v          payload_type  
31 v          payload_mass_kg  
32 v          payload_mass_lbs  
33 v        }  
34 v      }  
}
```

Sample
Graph QL

Query

- * Schema is req. to be followed

gRPC

* Built on HTTP 2.

* → Can't be directly used by the browser

→ proxy is needed.

like gRPC - web

→ because gRPC require HTTP 2 which can't be directly accessed by browser

* Info is sent in form of Protocol buffers.
Schema's.

* Data is sent in binary format and typically lesser packed size is req. to send same data, when compared to REST.

* Powerful

→ Do things like streaming (both one way and 2 - ways).

* Newer

but

less tools available

Harder

but takes less time.

* No status code but have error messages.

Defining A Message Type

First let's look at a very simple example. Let's say you want to define a search request message format, where each search request has a query string, the particular page of results you are interested in, and a number of results per page. Here's the .proto file you use to define the message type.



```
syntax = "proto3";  
  
message SearchRequest {  
    string query = 1;  
    int32 page_number = 2;  
    int32 result_per_page = 3;  
}
```



API-Design

Operation	Rest	RPC
Create Read Update Delete	Post Get Put Delete	Create Read update delete.

Example :-

version

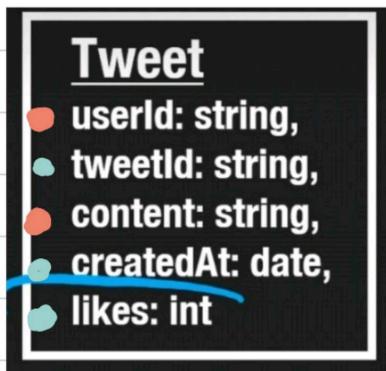
<https://api.twitter.com/v1.0/tweet> - POST

So, our API will take , these three as a parameter.

create Tweet [User ID, content, Parent ID=null]

Optional parameters are used in general to make the newer versions of an API, backward compatible

But sometimes, it is not possible to make the APIs backward compatible, in that case **versioning** comes into the picture.



- Required
- Can be created on the server directly.

To read the tweets, we may use :-

<https://api.twitter.com/v1.0/tweet/:id> - GET

This id helps us to get tweets one by one, else user's mobile won't be able to handle likes 100,000's of tweets simultaneously.

Moreover, many times **pagination**, is done in order to solve the above issue

<https://api.twitter.com/v1.0/users/:id/tweets?limit=10&offset=0>

user
ID

we want
tweets.

Query
parameters.

Get requests are not used to update something on the server, because they are cached by default. Therefore, every time the value will change.

Caching

CPU can read fastest from cache.

Objective

- Increase throughput
- lower the latency
- reduce the network cost.

Done by moving some data from primary storage to cache.

like suppose, I open a site very freq. and it is static in nature. So, instead of making network calls every time, the browser will cache the results of the older calls.

Header has

cache-control :- max-control - 3600

The content can be cached for at max 3600 sec.

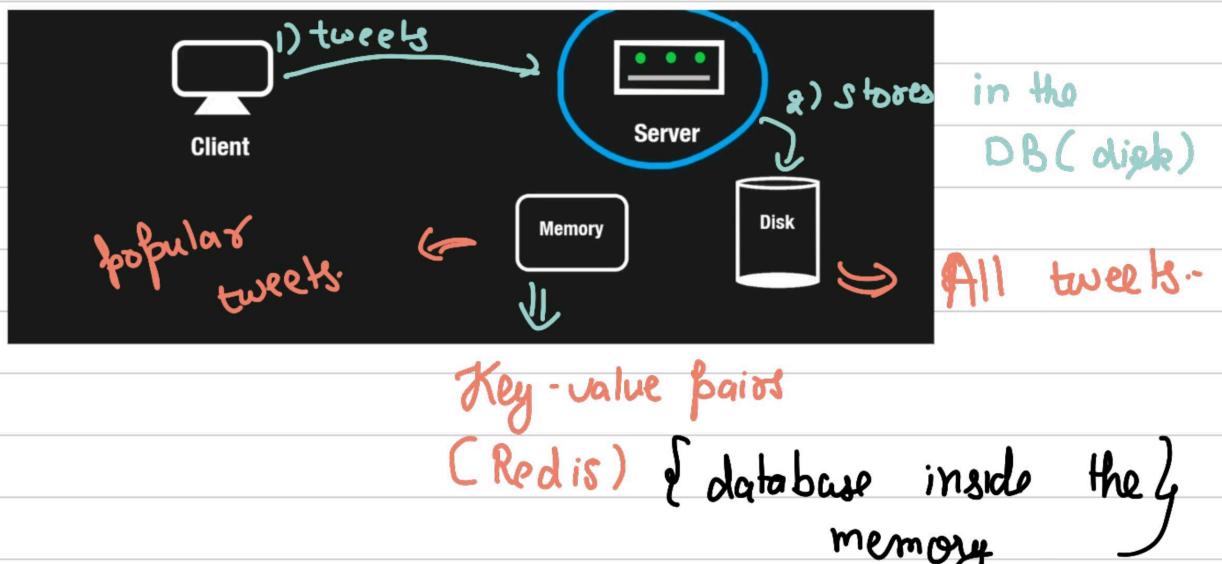
Data Disk >>> Data network.

Cache hit \Rightarrow if (cache[data] != null)

⇒ Cache miss ⇒ if (cache [data] == null or
cache [data] is expired)

⇒ Cache ratio :- $\frac{\text{hit}}{\text{hit} + \text{miss}}$

⇒ Caching from the server's perspective:-



Suppose, user wants to search for a tweet inside the D.B, now this get Tweet can be cached.

Example:- No one is going to read my tweets (so store in disk, because no. of req. ≈ 0)
But, for Narendra Modi, no. of reader will be high, so high no. of get requests. So, store in the memory as well as disk.

write-around cache :-

skip the cache entirely in the starting.

Disk
"

When we read, if cache

miss, then find it in the disk
and write it in the cache
and return to the user.

Write - through Cache :-

Write in memory / cache first
and then inside the disk.

Write - back Cache :-

might be used
for the case
of liking and
disliking tweets.

- * lazy way
- * faster when possible
- * But less reliable.

Problem

Soln

- * Store only into the memory not inside the disk.
fails in case of power supply
outage because ram are volatile
- * So, periodic cycles, the data is also written to the disk.

- * Used, where we can afford to loose some amount of data.

Eviction Policy

- * Cache limited in size
- * So, what to do when the cache is full?

- * Can do randomly or
- * Objective
 - * Maximize cache hit ratio.

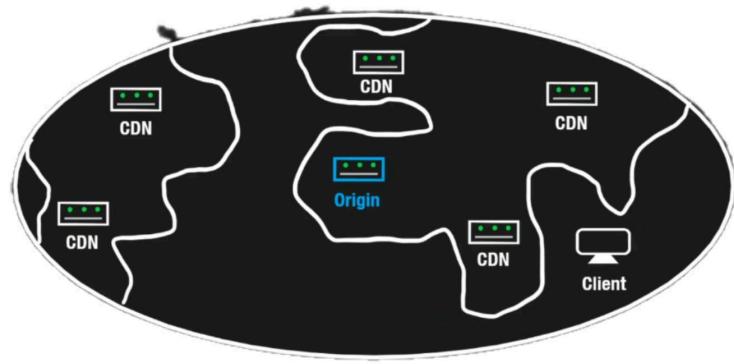
- * FIFO : first in first out.

More important

- * LRU : Least Recently used

* LRU :- Least frequently used.

CDN's



Content delivery Network.

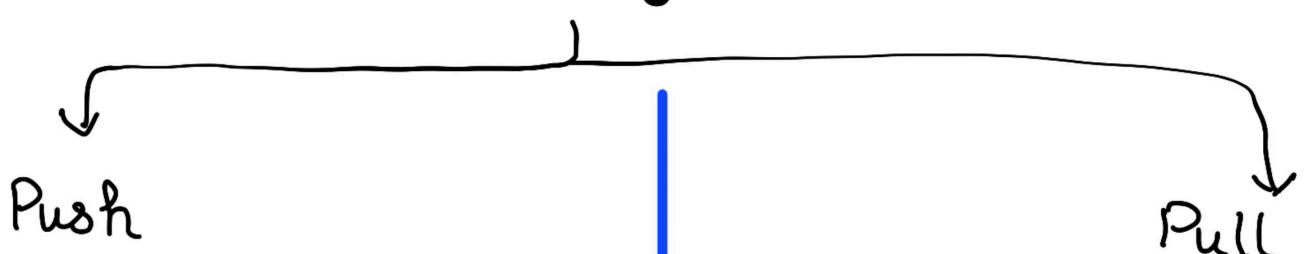
Used to reduce network latency. by making copies of ^(image/video) content / static code and placing them at various locations in the world.

Can't put everything on the CDN server

→ only static content can be stored here.

→ But these days Edge servers are getting introduced, which can run code on them.

Types of CDN:-



* Origin updates the data in the CDN.

* In case of a miss, CDN will call origin for data,

get it and cache it.

- * E.g
 - * User changes D.P.
 - * update DP at origin
 - * Origin pushes the same change to all the CDN's.
- * E.g
 - Match streaming (India vs Pak)
 - Now, only CDN near India will get the most req. So it will pull the same from Origin and store it with him.

* Cache control : public

The file is static in nature and can be cachable.

Proxies & Load Balancers

* Proxy (Forward)

- ⇒ A server which acts as a middle layer.
- ⇒ Hide client from the main server.
- ⇒ Proxy hides user's info.
- ⇒ Used to bypass network restriction.
like suppose a website is blocked,
then .

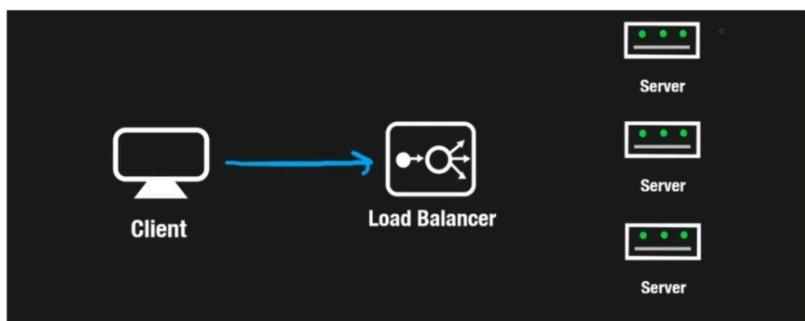
then we will access proxy. Proxy will access the website and return the content to us.

⇒ Proxies are also used to block some website like in the case of a school / corporate network.

Reverse Proxy :-

- * Hides server from the client.
- * But server will know about client.
- * Eg. CDN. and Load Balancers.

↓
Selects which server will serve our request.
↓



So, that every server will have eq. amount of load.
↓

weighted
←

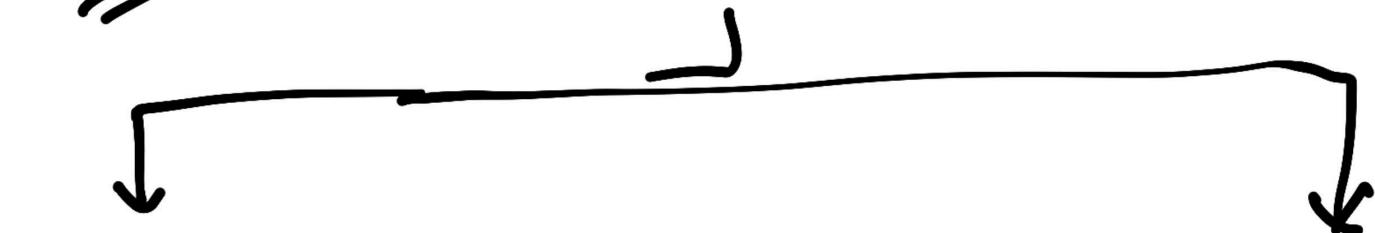
↓
Implement with REST

↓

Ways to implement the load Balancer.

Built using wto Round Robin. method
or
Server with least no . of connections
or
server nearest to the client.
or
using the Hashing.

Types of load Balancer

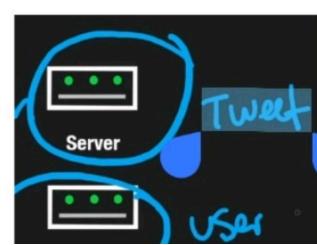


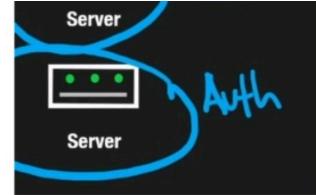
Layer 4 (TCP)

→ Location based
or
Round Robin.

Layer 7

(Application Layer)
→ Based on type of request.





→ More powerful but also expensive.

What if our one and only load Balancer goes down?

So have multiple / backup load balancer. But generally do not do not do much and hence won't go down.

Open Source Load Balancers :-

Nginx

Hashing.

#

Same as implementation of maps.

Suppose we have 'n' servers numbered from 0 to n-1. Now, take the MOD of the IP with N, and route the request accordingly.

So, a user always approaches a same server.

→ Benefit

↳ when the servers are stateful and then they won't req. to share their cache with one another.

→ loose

↳ It might be that user won't be going to the server with least amount of load.

↳ ex., one server crashed.

Now. 3 Requests $\Rightarrow \{9, 10, 11\}$

Initially (we had 3 servers)

Now, we have 2.
(Server 2 crashed)

$9 \% 3 \Rightarrow$ Gone to 0

$9 \% 2 \Rightarrow$ Goes to 1

1 Server
changed

$10 \% 3 \Rightarrow$ Gone to 1

$10 \% 2 \Rightarrow$ Goes to 0

$11 \% 3 \Rightarrow$ Gone to 2.

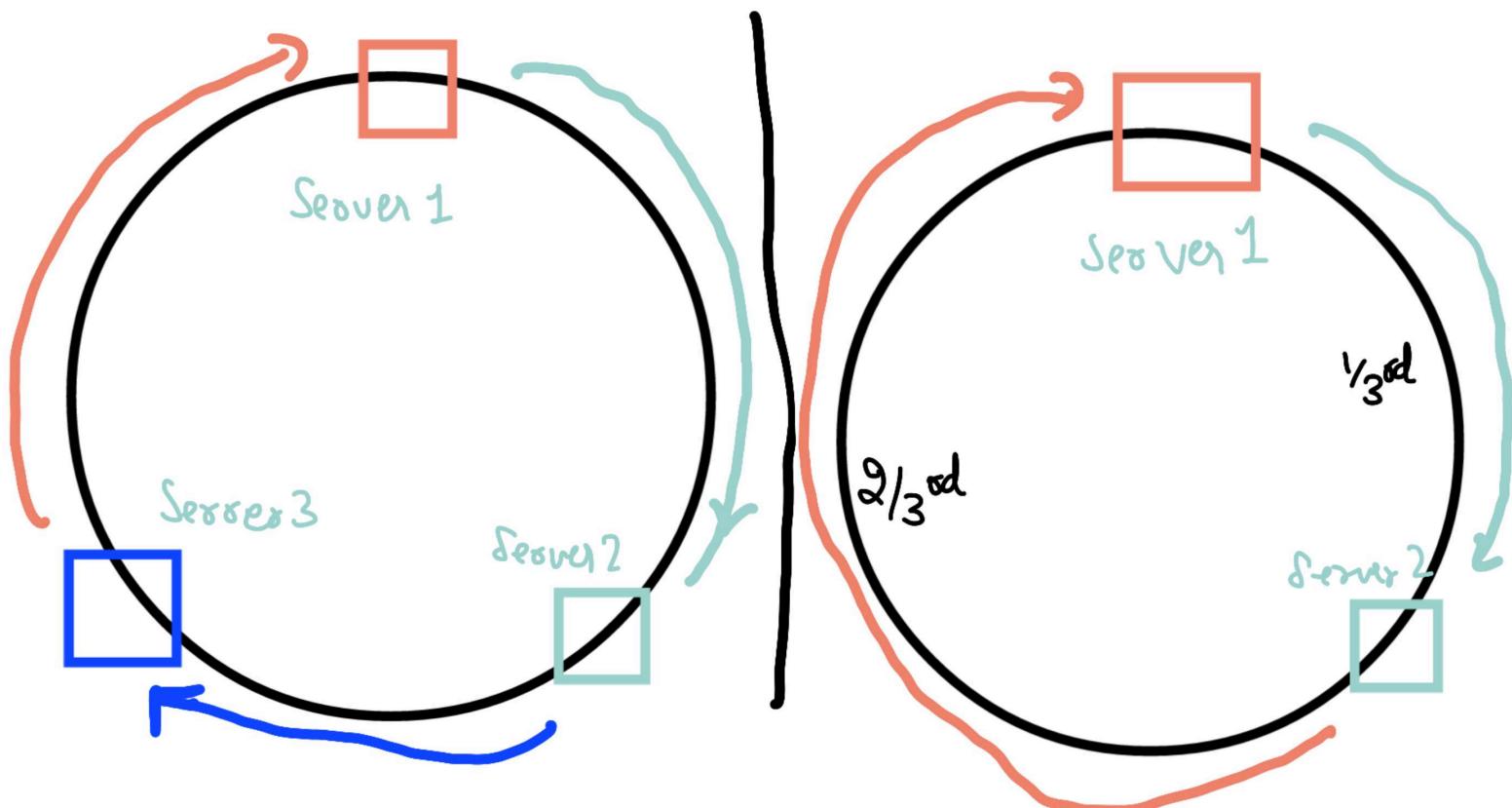
$11 \% 2 \Rightarrow$ Goes to 1

So, they lost
their data.

To solve above cache problem, consistent hashing is used.

Method of Consistent Hashing.

- ↳ Mark the servers on the circle at equal angles.
- Now, hash the incoming IP address.
- If the angle == angle of server then good, else move clockwise till you encounter a server.



Benefit :- Cache is preserved.

Downside :- Load Balancing do not occur in a good manner.

SQL

To improve consistent hashing, virtual servers are made. as explained by Gaurav Ben.

- ↳ Standard Query Language.
 - ↳ RDBMS
 - ↳ Stored on Disk.
 - ↳ want to read / write very efficiently
 - ↳ Generally B+ Trees are used for the same.
 - ↳ m-1 keys
 - ↳ stored in leaf node
 - ↳ sorted.
 - ↳ Schema
 - ↳ what is the field name
 - ↳ what is the data type.
 - ↳ Tables are related to each other with help of foreign key.
 - ↳ Every row should have a primary key.
 - ↳ Joins can be performed.
- These are ACID compatible.

ACID

A ⇒ Atomicity

C ⇒ Consistency

I ⇒ Isolation

D \Rightarrow Durability

for durability

- * Try to store things in disks in place of storage
 - * Speed of read/write reduces.
 - * Our data won't be lost in case of power outage.

Transaction

- ↳ A collection of queries.
- ↳ Start with the keyword 'Begin' and end with 'Commit'.

Atomicity

- ↳ A database transaction can't be splitted.
- ↳ A transaction should either be completed or should do nothing.
 - ↳ suppose, server crashes while transferring money, now the money is debited from A, but never able to reach B.

↳ More on atomicity

Isolation

↳ Isolation of one transaction from another.

elsewise

↳ Dirty Reads

↳ Reading inconsistent value (value which is not finalised yet).

↳ Phantom Reads

or will happen race condⁿ.

↳ Therefore, all the transactions should be serialised to avoid such problems.

↳ Think of B and C sending A (10 initially), 100 \Rightarrow . now final balance of A should be 210. But it might be 110 if both happened parallelly.

transaction B \Rightarrow Balance A (10) + 100 = 110

transaction C \Rightarrow Balance A (10) + 100 = 110

Consistency

↳ Like ^{name} column should not be null.

↳ Should have a primary key.

No SQL

- ↳ Not only SQL
- ↳ open ended.
- ↳ Non-relational databases.
 - ↳ So, Can't use SQL to query them.
- ↳ Solves problem of scaling up a d.B.

Eg of NoSQL D.B.

→ Key-Value :- Redis, Memcached, etc.
 |
 | Works like a hashmap.
 |
 | Adv :- very fast because uses RAM.

→ Document Store.

Document Store

```
{  
    "field1": {  
        "onetyp": [  
            {"id": 1, "name": "John Doe"},  
            {"id": 2, "name": "Don Joeh"}  
        ],  
        "othertyp": {"id": 2, "company": "ABC"}  
    },  
    "field2": {  
        "a list": [[1,42],[2,2]]  
    }  
}
```

A collection of Document
↓
A typical JSON object.
Example : MongoDB.
No Schema

→ Wide Column DB

- ↳ Handle massive massive scale
- ↳ mostly used, where a high no. of writes are going on.
 - ↳ like time-series D.B.
- ↳ And low no. of read and updates.

↳ Example

- ↳ Cassandra
- ↳ Big Table

→ Graph D.B

- ↳ Used in general in social media sites.
- ↳ Like A follow B, who follows C, who in turn follows D.
 - ↳ So, to do this in SQL, we have to join 4 tables(A, B, C, D)
- ↳ But graph GL can do it easily.
 - ↳ Uses directed graphs.

Scalability in NoSQL

Challenge 1.

- ↳ It is hard to scale a DB horizontally, because in that we won't know, which part of D.B. is our data located. Is it present in server 1 or is it present in the server 2?
- ↳ Or, suppose we are trying to join data from server 1 to data in server 2.
- ↳ So, it is why people try to prefer to vertically scale a database.

Challenge 2.

{ we have to maintain ACID properties in case of a SQL database. But that's not the case with NoSQL DB's. They sacrifice the acid property.

Acronym of NoSQL database.

Base

B - Basically Available

Users can access the database concurrently at all times, and one user doesn't need to wait for others to finish a transaction before updating a record.

S

S - Soft state

Without a guarantee of consistency, the state of the system can only be predicted with a certain probability after some time.

E

E - Eventually consistent

If the system is functional, the data will eventually get to a consistent state after enough time has passed.

Eventual Consistency

Suppose we have Server 1 as the main server and Server 2, which stores the copy of Server 1.

Now, what if some write happens in the Server 1? Because no update in Server 2 will make the D.B inconsistent.

So, it becomes the responsibility of the Leader (Main D.B.) to update the follower (Server 2)

Benefit

↳ Generally, the no. of reads are more than no. of writes. Therefore the above scenario comes into the play / picture.

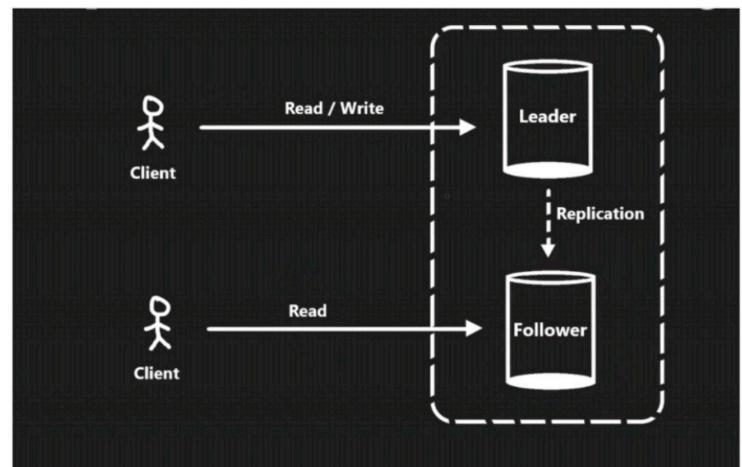
Replication (Leader - Subscriber) model

→ Creating copy of the D/B.

→ Used when ...

handle the no. of requests it is getting.
(Read/write)

→ In this case, we are allowed to only read from the Replica (follower).



Reason, because as we have read in eventual consistency, it is the responsibility of the leader to write in the follower.

→ Can be used to scale up the read requests.

→ Challenge

Asyc Data update to the follower.
Like data _{Leader} updated at 12:00 am.

Some one tries to read data from follower at 12:02 am?

Sync Data

- ↳ Replicate the write req. immediately.
- But this increases latency.

Steps

These extra steps increases the latency.

- i) Client initiates a transaction.
- ii) Update the leader.
- iii) Leader update the follower.
- iv) Get an ACK of the work done successfully.
- v) Leader tells the client that the write req. successful.

Benefit

- ↳ Increased availability and reliability.

Replication

Leader follower

Leader - Leader.

Downside

- ↳ Becomes out of sync easily.
- ↳ Do not remain consistent.

Pros

↳ Useful when we can divide the user base into two - or more sub groups and later assign D/B to them. Now if the D/B's get out of sync, that won't be an issue too.

Sharding

Partitions into smaller parts

~~use~~.com
An approach of horizontal scaling. Used when we have data on the scale of Peta Byte.

Queries are taking multiple seconds to run.

~~what we do~~

- ↳ we not only keeps the data on the multiple machines but also split the data into multiple smaller DB's and putting these smaller DB's into their each separate machine.

How to Shard?

* Range Based

- ↳ like if first name has first character in the range of A to L, then put it into Server A, otherwise put into server B.

Cons

↳ Running joins will become both slow as well as complicated.

- * Hash Based / Consistent Hashing.
- * Popular SQL databases such as MySQL or Postgres Do not support sharding by default. because we loose the consistency
- * Most NoSQL DB's have sharding by default.

CAP - Theorem.

C \Rightarrow Consistency } Relates to
A \Rightarrow Availability } Databases.
P \Rightarrow Partition Tolerance.

* It applies only to replicated Data base. Therefore

It won't be applicable to SQL d.b.

* CAP states that the 'P' is always applied and it remains our choice to choose from 'C' and 'A'.

* Partition :- It refers to the situation where a part of the D/B becomes disconnected from the rest of the D/B network.

Like due to some reason the connection b/w the leader and follower broke. The users are still able to connect to the leader and follower but, the leader is not able to connect to follower and vice versa.

* Partition Tolerance :- Our system will continue to work even if we have got a partition. Because, we won't want that due to a small partition, our whole system goes down.

* Consistency :- Data - consistency b/w the D/B's.
Every single read will get the most up to date data.

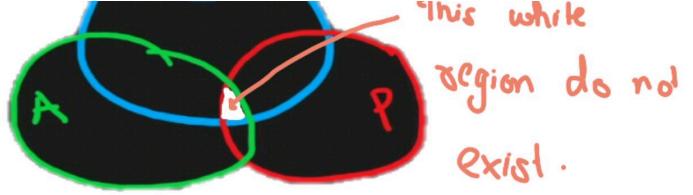
∴ In case of a partition, the updates b/w the leader and follower stops. Now, either we can take down the replica (Availability, ends) and make the DB consistent or let it be inconsistent but continue to serve.

* Availability ⇒ Availability of a Every node.
Na ki 99.9% availability.

∴ In case of a partition, the updates b/w the leader and follower stops. Now, either we can take down the replica (Availability, ends) and make the DB consistent or let it be inconsistent but continue to serve.

There are a very few D.B's which fits the 'CA' definition.

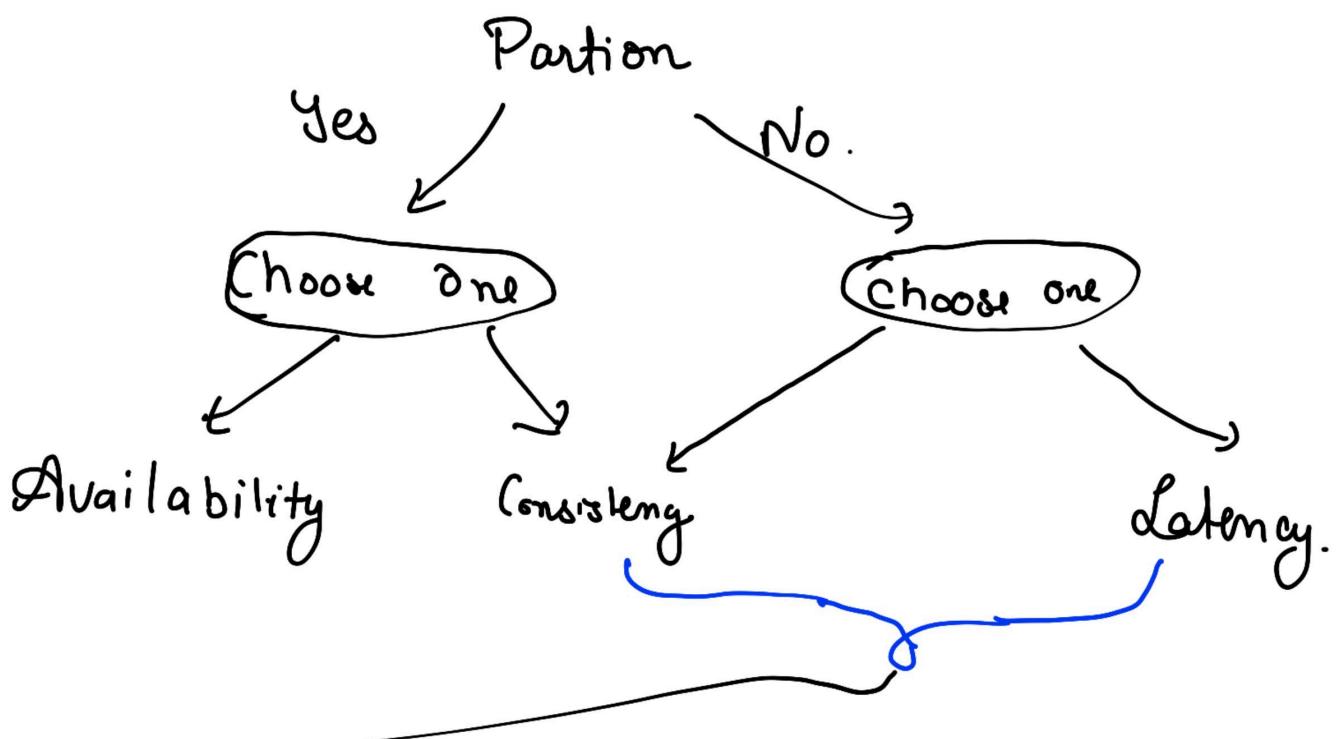
A - Availability
P - Partition



PACELC

Given P, choose A or C.

Else favour the Latency or Consistency.



→ Basically, when the partition is not there we might choose Consistency (by every write req. updating the replica and increasing the latency) or reduce the latency by choosing not to make the D.B.'s consistent.

OBJECT STORAGE.

- No hierarchy. → ∵ we can immediately find the objects.
- Data is stored in form of a flat way.
- Example :-
 - * AWS S3.
 - * Google Cloud Storage
 - * Azure.
- Can't update thing we add.
- Optimized to store large amount of data.
- Same as a Hashmap
 - * Unique key for every object.
 - * Therefore, can't store duplicate objects.

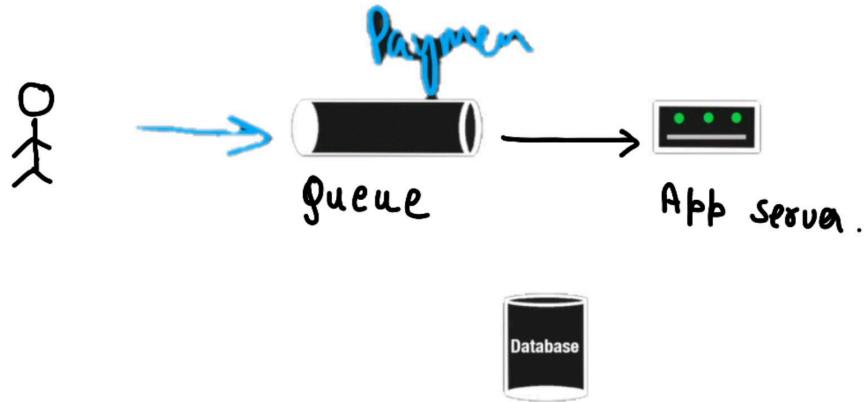
↳ Should have a globally unique names.

- * Videos and Images are stored in the Object storage as they are never going to be queried and will make the D.B. slow if stored.
- * Reading and writing is done using the HTTP. (but have own security problems.).

MESSAGE QUEUES.

↳ If we have a large amount of application events and they are going faster than what our app server can handle. And we do not want to scale up the server.
In this case, message queues comes into the picture.

Message Queues



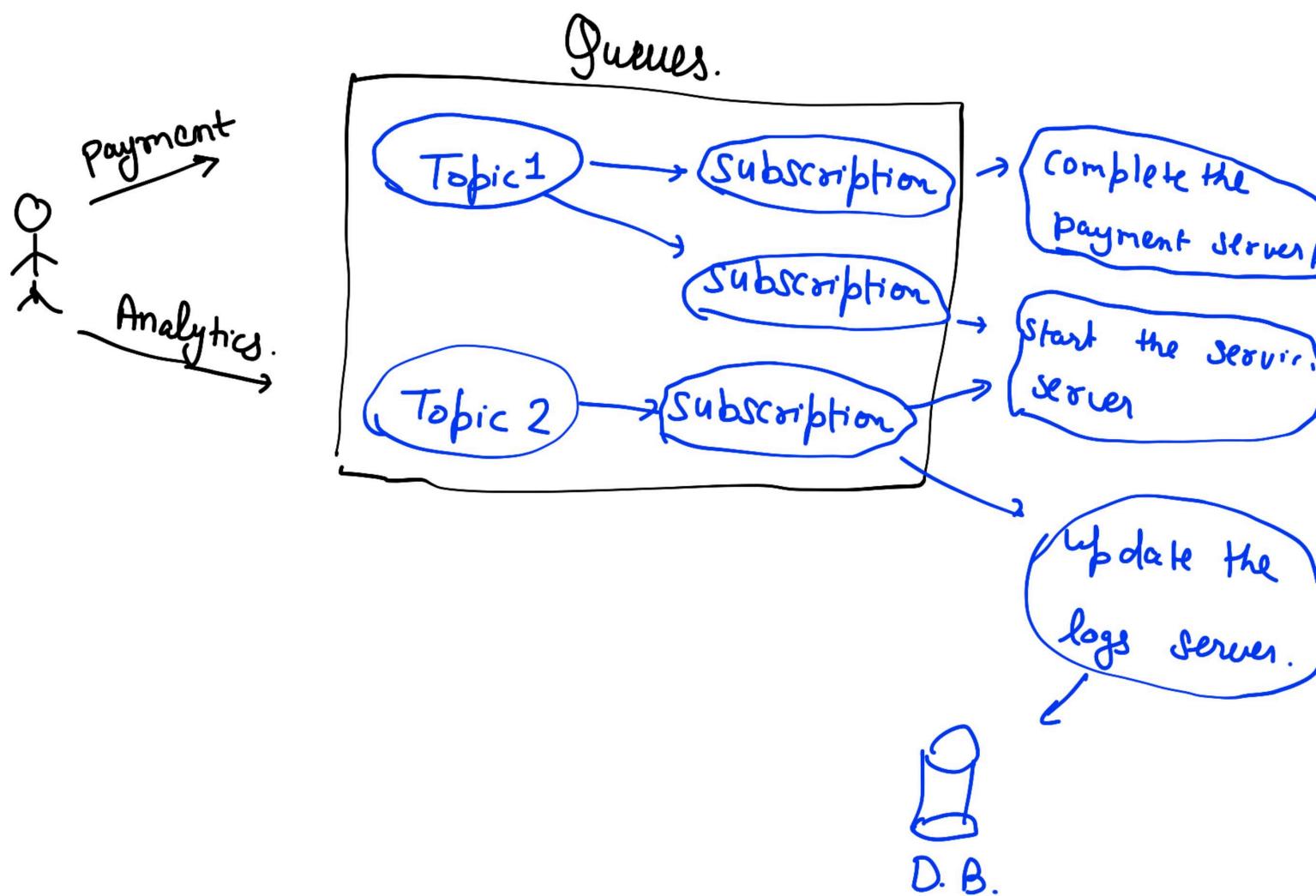
- * The queue is required to be durable. That is, if the system crashes, the message should still be there in the queue.

A Message Queue (MQ) allows applications to communicate and process operations asynchronously. This means that the sending application can send a message to the queue and then continue its own processing without waiting for the receiving application to process the message. The receiving application can retrieve and process the message at a later time, thus decoupling the sender from the receiver and enabling asynchronous communication. This approach is particularly useful for improving system scalability, reliability, and responsiveness. ▲

- * The queue will store the incoming requests into it, and when the server gets free it releases the requests. The order remains FIFO.
- * Queues uses 'Push.' to send the req. to the server. In case, server is not able to handle it, it will send a 'Pull' req. to the consumer.

In case of successful completion of a request, server sends a 'ACK' for the same. If no ACK, queue will try to resend the same req. to the server until server acknowledges or the max. trials for one req. runs out.

In case of many applications , sending req.'s to the queue, we get the Pub/Sub mechanism.



Sending same topic to many subscription is called. Fanning out.

Same req. can go to multiple applications.
Same Application can receive multiple type of requests.

Example.

