

~~system~~

ArgoExpress

content: 99 Fundamental Le

Amman Barnwal

Amman Barnwal

Resources: ① Grokking the system design (Educative.io) pdf

② System design video (clement)

③ Gaurav Sen

④ Tech Dummies

⑤ audicode (Yogita's channel) → System Design Primer course

⑥ GitHub link →
 ~~System Design Primer~~ →
 ~~System Design Interview~~ →
 ~~System Design Interview~~ →
 shashank88

Other Top YouTube channels for HLD:

① CodeKode

② Think Software

③ Udit Agarwal

④ Success in Tech

⑤ Jackson Gibbard

⑥ Coding Simplified

⑦ System Design Interview

Blogs and websites:

highscalability.com

Interviewbit → System Design

Gaurav Sen's Answer (Aurora)

Syllabus

Main Topics of System Design

Working system design (Educative.io)

- System Design Basics
- Key characteristics of Distributed systems
- Load Balancing
- Caching
- Data Partitioning
- Indexes
- Proxies
- Redundancy and Replication
- SQL vs NoSQL
- CAP Theorem
- Consistent Hashing
- Long Polling vs Websockets vs Server-Sent Events

System Expert (Element)

- Client-Server model
- Network protocols
- Storage
- Latency and Throughput
- Availability
- Hashing
- Relational Databases
- Key-value stores
- Specialized storage paradigms
- Replication & sharding
- Leader Election
- Peer-to-peer networks
- Polling & streaming
- Configuration
- Rate Limiting
- Logging and monitoring

copy
of this
copy

Windows Server 2016

Storage Spec. for Sys. ①

Data partitioning ②

Load balancing ③

Caching ④

Replication ⑤

Redundancy ⑥

SQL vs NoSQL ⑦

Consistent Hashing ⑧

Long Polling vs Websockets vs Server-Sent Events ⑨

Client-Server Model ⑩

Relational DB ⑪

Key-value stores ⑫

Specialized storage ⑬

Replication & sharding ⑭

Peer-to-peer networks ⑮

Polling & streaming ⑯

Configuration ⑰

Rate Limiting ⑱

Logging and monitoring ⑲

Windows Server 2016 ⑳

Storage Spec. ㉑

Data partitioning ㉒

Load balancing ㉓

Caching ㉔

Replication ㉕

SQL vs NoSQL ㉖

Consistent Hashing ㉗

Long Polling vs Websockets vs Server-Sent Events ㉘

System Design Problems

edabit.io

- ① step by step guide
- ② designing a URL shortening service like tinyURL
- ③ designing Pastebin
- ④ " Instagram
- ⑤ " Dropbox
- ⑥ " Facebook Messenger
- ⑦ " Twitter
- ⑧ " YouTube or Netflix
- ⑨ " Typeahead suggestion
- ⑩ " on API Rate Limit
- ⑪ " Twitter search
- ⑫ " a Web crawler
- ⑬ " Facebook's Newsfeed
- ⑭ " Yelp or Nearby Friends
- ⑮ " Uber Backend
- ⑯ " Ticketmaster
- ⑰ design AlgExpert
- ⑱ design a stockbroker
- ⑲ " Google Drive
- ⑳ " Reddit API
- ㉑ " Netflix
- ㉒ " Uber API
- ㉓ " slack
- ㉔ " Airbnb
- ㉕ " a code deployment system

Latency and Throughput:

Latency and Throughput are 2 most important measures of performance of the system.

Latency:

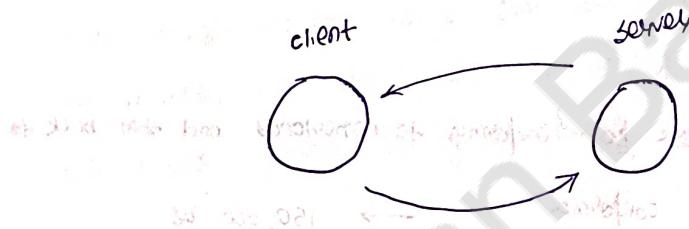
Latency is basically how long it takes for data to get from one point in the system to another point in the system.

Eg: Latency of a network request

(How long it takes for one request to go from a client)

and then back from the server to a client

Let's say, you are reading the data from either disk or memory then the time it takes to read data will also referred to as latency.



There is a tradeoff in the ways systems were built; bcs certain things will have higher latency but other things in the system will have lower latency.

Note:

① If you are reading 1 MB sequentially from memory, time it takes is 250 μs.
(i.e. It is a very fast operation).

② Reading 1 MB sequentially from SSD (solid state drive) → 1000 μs

i.e. Reading data from memory is a lot more faster than reading data from solid state drive (SSD).

- ③ sending 1 MB data over 1 Gbps network \rightarrow 10,000 us
($1 \text{ Gbps} = 1 \text{ Gigabit per sec}$)

(we are assuming that we are sending 1 MB to the machine that is right next to us i.e. we are not accounting any distance here)

Eg: Making a network request
" " API call

- ④ Reading 1 MB from HDD (hard drive) \rightarrow 20,000 us
(Hard drive)

Depending on the hardware that you have, sending something over the network (i.e. HDD vs SSD)
might actually be faster.

- ⑤ sending a packet over the network from California to Netherlands and then back to California \rightarrow 150,000 us

(lets say, packet here is much smaller than 1 MB)

Eg: 1000 to 1500 bytes typically

The time taken a lot more time than the remaining.

→ when you are designing the system, you are typically gonna optimize those systems by lowering the overall latency of the system.

some systems will not necessarily want to optimising for latency as much as others.

Eg: certain types of system might really care about low latency
(Eg: A High Trading Firm)

Video Game

i.e. If you find some lag in a videogame, that just means you as a player in that video game is having high latency. That might be bcz the server in which you are playing on might be located half way across the world from you and it just takes a while for a computer or client to make a network request to a video game server.

Even though 150,000 ms is fast, but if you are sending a lot more data over the network, then this 150,000 ms quickly add up and it will have a really bad user experience bcz of high latency.

→ on the other hand, certain websites might care less about latency. But what they care about more is accuracy or uptime of the website.

Throughput:

Throughput is how much work a machine could perform in a given period of time.

(we are referring really to how much data is transferred to one point in the system from another point in the system in a given amount of time).

→ Typically, we measure throughput in gigabits per sec i.e. Gbps

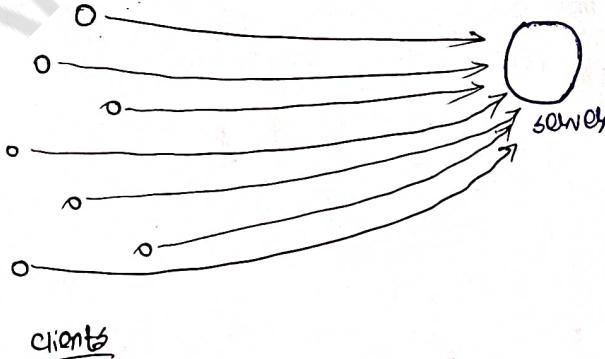
or we can also measure by kilobits per sec i.e. Kbps

or megabits per sec i.e. Mbps.

Eg: sending 2 MB data over 1 Gbps network. time $\rightarrow 10,000 \text{ us}$

This just means this network can support 1 Gigabit per sec

Eg: lets say, we have a bunch of clients and they all are issuing request to the same server.



Here, throughput is basically how many these requests can this server handle in a given amount of time.

or

how many bits can this server handle per sec?

How will I increase the throughput?

How will I optimize the system for throughput?

→ You just pay to ↑ the throughput.

Eg: The things that determine how many bits can go in and out of Algoexpert

server at any given time is determined by a cloud provider.

(Eg: Google cloud platform)

So, we will just pay these cloud providers to increase the throughput

i.e. we will just pay these cloud providers to increase the throughput

Note: Above is just a naive soln.

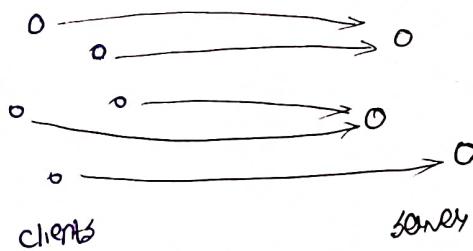
Let's say, a single server is handling multiple clients having a ton of requests.
It might be Google search or Facebook messenger that is expected to serve thousands of
requests or perhaps millions of request per sec.

so, just blindly increase the throughput on this network probably won't make sense.



so, a better way to fix this is to have multiple servers, handling these requests.

Eg:



Note:

Even though Latency and Throughput are very much related, they are not necessarily co-related.

Eg: You might have a system of small parts of that has very low latency i.e. they support really fast data transfers.

while you might have other part of the system which can have really low throughput.

i.e. ~~one~~ two parts of the system can have different latencies and throughputs based on one another.

You can't make assumptions on the latency and throughputs based on, one another.

Availability:

You can think of availability of a system as

availability → how resistant the system is to failure?

How fault tolerance is the system?

(Reported as system's fault Tolerance)

you can consider % of time in a given period of time like a month or year
or at least operational enough such that all of its primary functions
are satisfied.

→ Availability is a very important aspect to think about when you are evaluating a system.

→ Most systems almost have employed guarantee of availability.

Eg: when you purchase AlgoExpert, the main thing you purchase is access to
our content and courses.

so, availability matters a lot for any company.

There is a employed guarantee of availability when you purchase AlgoExpert.
but you expect the website of AlgoExpert to be up and fully operational.

so, availability matters a lot for any company.

→ There are varying degrees of availability that you can expect from different systems.

Eg: Let's say, we are dealing with a software that supports Airplane software.

The software that allows Airplane to function
properly when it's in flight.

if that software goes down, that would be fully unacceptable bcz it is directly attached
to lives of people flying in that aeroplane.

So, in this system, you would expect extremely high amount of availability.

Eg: 2nd example of this could be YouTube.

If YouTube ever goes down, then that's really bad. Bcz millions of users use YouTube everyday.

Eg: or take cloud providers for instance.

Let's say, if parts of GCP or AWS, then that will affect all the business and the customers that rely on their services for their own platforms.

How do you measure availability?

→ we measure typically availability as the percentage of system's uptime in a given year.

Eg: if the system is up and operational for half of an entire year, then you

would say that system has 50% availability.

which is really really bad for most services.

You could imagine if Facebook, Uber, etc is down for half of the time.

Nines

→ most systems or services often aim for really really high availability. so, we often end up measuring availability not exactly in percentages but rather in nines.

Nines are effectively percentages but they are specifically percentages with the number 9.

Eg: if you have systems having 99% of availability, meaning that it is up 99% of the time during a year then we say that your system has 2 nines of availability

(bcz the number 9 comes twice in this percentage).

Eg: (i) If your system has 99.9% of availability then we say that your system has 3 nines of availability.

(ii) " " has 99.99% "

- 4 nines of availability

Note:

% of Availability

Nines of availability

Downtime per year

Similarly,

90%

→ 1 nine

→

36.53 days

99%

→ 2 nines

→

3.65 days

99.9%

→ 3 nines

→

8.77 hours

99.99%

→ 4 nines

→

52.6 minutes

99.999%

→ 5 nines

→

5.26 min

99.9999%

→ 6 nines

→

31.56 sec

99.99999%

→ 7 nines

→

3.16 sec

99.999999%

→ 8 nines

→

315.58 ms

99.9999999%

→ 9 nines

→

31.56 ms

→ calling a system highly available is an actual recognized term in system design.

also and it's often abbreviated as "HA"

means for Highly Available

so, a lot of system aims to be highly available system bcz they really care for availability.

→ Availability matters a lot - both for system designers and end user of the system.

→ In fact, it matters so much that for certain systems, you don't have an employed guarantee of availability, you have an explicit guarantee of availability.

eg: Amazon guarantees 99.99% availability for its services.

SLA :

→ Many service providers out there have an SLA .

↓
Service Level Agreement

SLA is an agreement b/w service providers and user of this service on that system's availability on that other things.

so, many service providers have explicit written SLA that basically tells customer that "Hey, we guarantee this amount of availability at this % of uptime of system".

SLO : It stands for "service level objective".

You can think of SLO as the components of SLA.

SLA is made up of SLO.

Ex:

→ If I provide service for you and if I guarantee some percentage of uptime for that service, that % of uptime guarantee would be SLO.

If I guarantee you that you'll get only a no. of errors when you use my service, then that is another SLO.

Note :

All cloud providers like GCP or AWS have the very clear cut SLA's.

Infact, you could see them on their product pages.

These are very important when you are considering buying their services.

SLA is very serious and availability is taken very seriously by any firm.

→ Availability is not always super super important i.e. it's not always the case that you must have 5 nines of availability for all systems.

bcz having high availability comes with trade offs, such as higher latency and low throughput.
it's difficult to achieve high availability

so, you as a system designer have to think long and hard about whether or not your system should be highly available.

what part of system require higher availability and what part of the system would be okay to have lower availability.

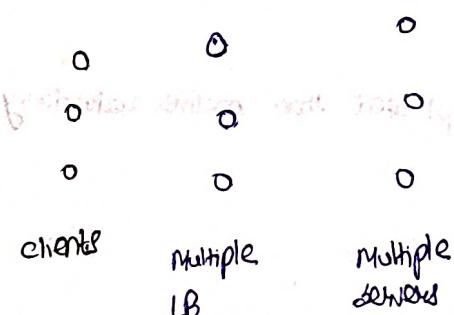
How do you actually improve the availability of the system?

→ You want to make sure first and foremost that your system does not have a single point of failure.
& How do you eliminate single point of failure?

↓
with Redundancy

Redundancy: It is basically the act of duplicating or triplicating or multiplying even more certain parts of the system.

Eg: Adding servers using multiple load balancers to avoid SPOF.



- You can introduce redundancy in a lot of different parts of your system just by literally adding machines to those parts of the system.

①

Passive Redundancy, when you have multiple components at a given level layer in your system. (using multiple servers and multiple load balancers, etc).

→ If at any point, one of these components fails, nothing really will happen bcs other components can handle the requests smoothly.

Eg of a passive redundancy →

An airplane engine

If one engine fails, the airplane can very smoothly function using other engine.

② Active Redundancy:

It is when you have multiple machines that work together in such a way that only one or few of the machines are typically handling the traffic.

→ If the machine which is handling traffic fails, the other machine are gonna somehow know that the other one failed and they are gonna take over.

→ Active redundancy is a little more complicated than passive redundancy.

Caching :

We use caching in algorithms to avoid redoing the same operation especially computationally complex operations, which might take a lot of time.

↓
So, we use caching to ↑ the time complexity & do speed up our algorithms.

→ In system design, caching is used to speedup the system. It is to reduce the cost of computation & improve the latency of the system.

You can use caching in a bunch of different places in a system.

Eg: You can cache at the client level. ie. caching can be somewhere at the client level such that it no longer has to go to the server to retrieve it.



similarly, you could cache at the server level. May be you always need server to interact with the client but may be you do not always need database to go to database to retrieve data.

→ You can also have cache in between two components of the system.

→ You can have a cache between a server and a database.

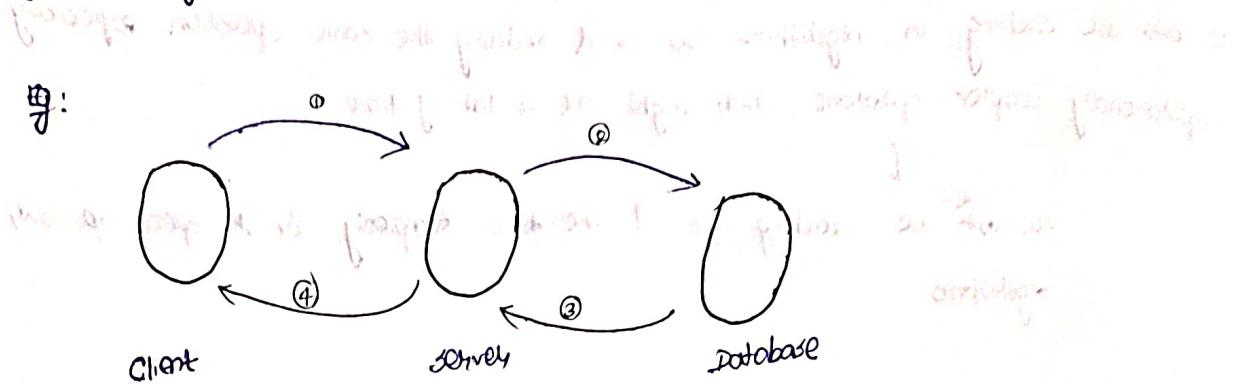
Maybe, you could have a cache between a client and a server.

→ In fact, you can even have caching at the hardware level. in modern day computers.

Eg: There are CPU caches that basically make it faster to retrieve data from memory.

SL

→ Caching is really helpful if you are doing a lot of network requests.



You can imagine here, the client might have made request to the server. The server might then make request to the database. The database is gonna return data to the server. The server is gonna return data to the client.

Perhaps here you might want to speed up this operation by caching the results of this network request to minimize as many requests.

Eg → Another eg in which caching is needed is when you are doing some very computationally long operation.



Eg: Imagine you have multiple servers instead of just 1 server. All of these servers are hitting the database with the same network request. May be there are bunch of clients, they may be there are bunch of clients, they also are making individual requests to the different servers to get a certain Instagram profile, for instance. (Eg: A popular celebrity)

Here, you can use caching, not to increase the speed of each network request when it's trying to get the profile. That's not something that you need to optimize on. Individual network requests would be very quick.

But you don't want to read from dB a million times so you can use caching so that you do not need data from dB each and every time.

Read operation from dB ↓

Note:

- Caching is used to either increase the speed of network requests or to either to speed up an operation that involves network requests by avoiding having to do network request altogether or atleast avoiding to do it multiple times.

either to speed up an operation that involves network requests by avoiding having to do network request altogether or atleast avoiding to do it multiple times.

May be you want to speed an operation bcz it's computationally long.

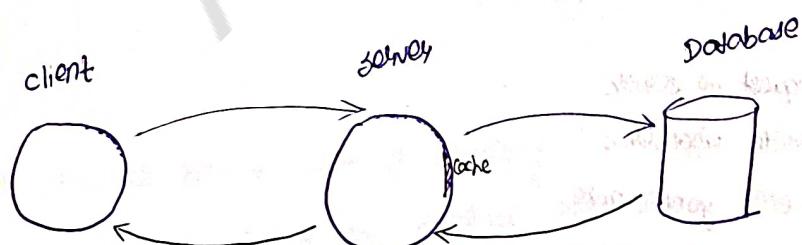
or may be, you have an operations that is being done tons of time, you don't necessarily need to speed them up individually. You just use caching here.

Redis → It is a very popular in-memory database (key-value store)

Eg:
→ Let's imagine, we are designing a web application where users can read or write posts.

And they can also edit their post.

Eg: writing facebook or LinkedIn post



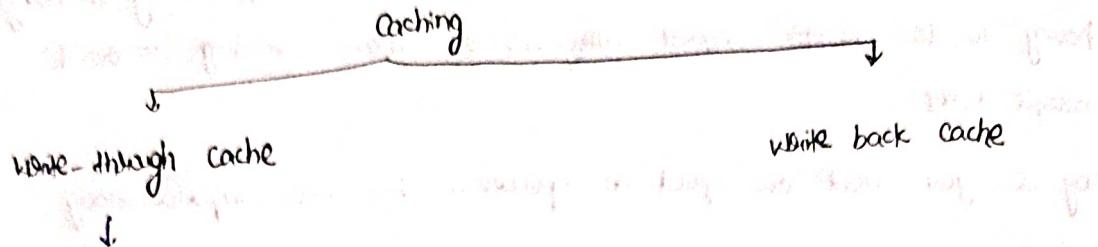
Let's say, you are caching the posts at the server level in memory; Now, you have 2 sources of truth. Your post is stored both in database and also in cache stored in your server.

Let's say, now you're editing an post that you wrote.

How do you deal with these 2 sources of truth?

and how do you know that when to write to the cache and when to write

to the database?



→ it is a type of caching where when

if you make an edit to a piece of data or when you write a piece of data,

your system will write that piece of data

to both in the cache and in main source

of truth at the same time or literally

at the same operation.

Eg:

Let's say, your post is stored both in db.

and in cache in the server.

And as a user, you want to edit your post.

You're gonna make a network request to server.

And the server is gonna override whatever is in the cache. And then also, it is gonna make a request to the database and overrides what's in the database.

On this way, the cache and db is always in sync.

But the downside here is you still end up going to the db.

Note: with write-through caching, every single time you are gonna override something in the cache and in the database, you are doing 2 things and you still end up going to the db.

Write-back cache:

Now back to the write-back cache. So here we are gonna do database writes to the client and then when the client does some update, it's gonna update the database and then the database will update the client.

Now if you do a write-back cache, the db info that you have can be inconsistent because there is no actual update to the database. Eg: let's say, you are a user and you are editing your post. You are gonna make a network request to the server and the server is gonna update the cache, i.e. only the cache. And then it will immediately go back to the client.

so here, your cache is not in sync with database.

Your system will asynchronously update the database with the value that are stored in the cache. It can be done either at every interval ie. every 5 sec or every 5 min, etc.

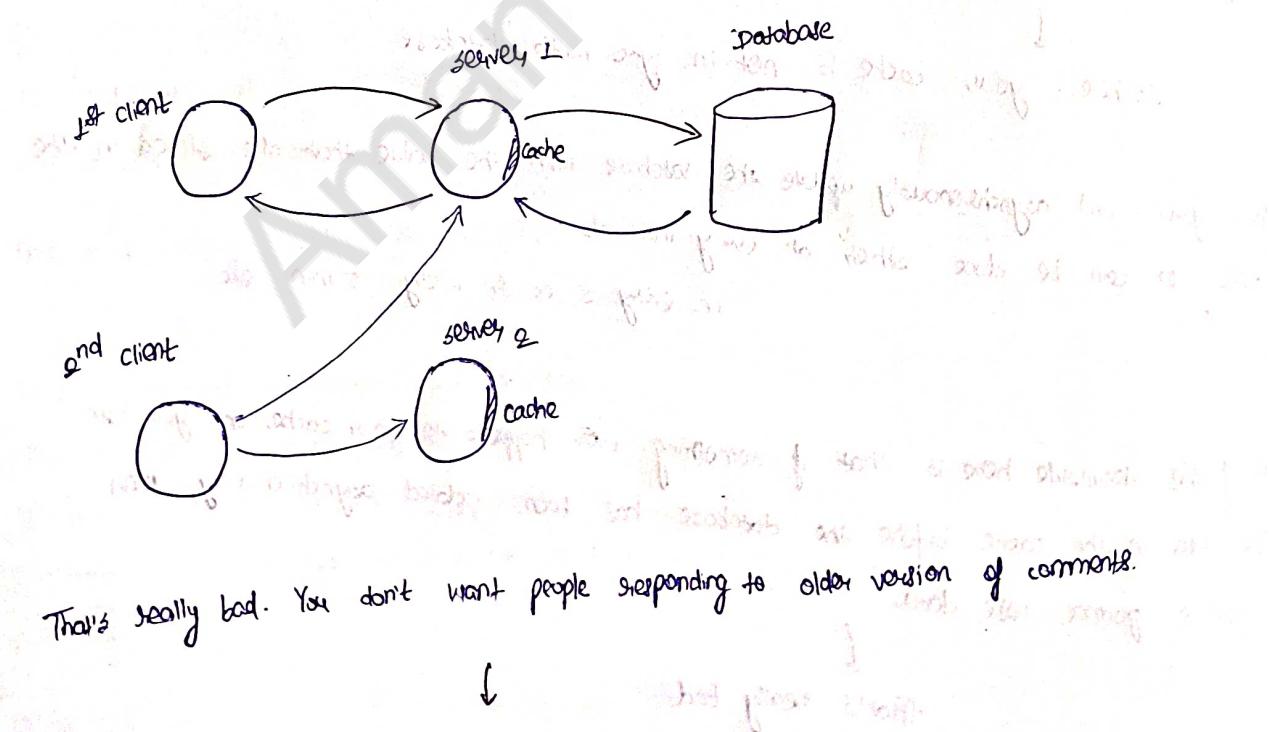
One of the downside here is that if something ever happens to your cache, and you lose the data in the cache before the database has been updated asynchronously, then you're gonna lose data.

That's really bad.

Ex: Let's say, we have to design YouTube comment section system. You have a bunch of servers. Every server caches in memory the comments on a single video. We have clients which are communicating to these respective servers. When they are reading comments, they just read from cache in their respective servers.

Now, let's say, 1st client posted a comment on the video. And our 2nd client goes to the video. The server goes to the dB to fetch all the comments and stores them in the cache.

Then 1st client edits that comment. And the 2nd client if goes to that comment sees the old version of 1st client's comment, not a newly edited one. (Bcz it doesn't go to dB. It fetched data from cache in the server). And 2nd client responds to the old comment, not the newly edited one.



That's really bad. You don't want people responding to older version of comments.

Result in staleness of caches.
Here, our client will often deal with stale caches, i.e. the caches that hold stale data.

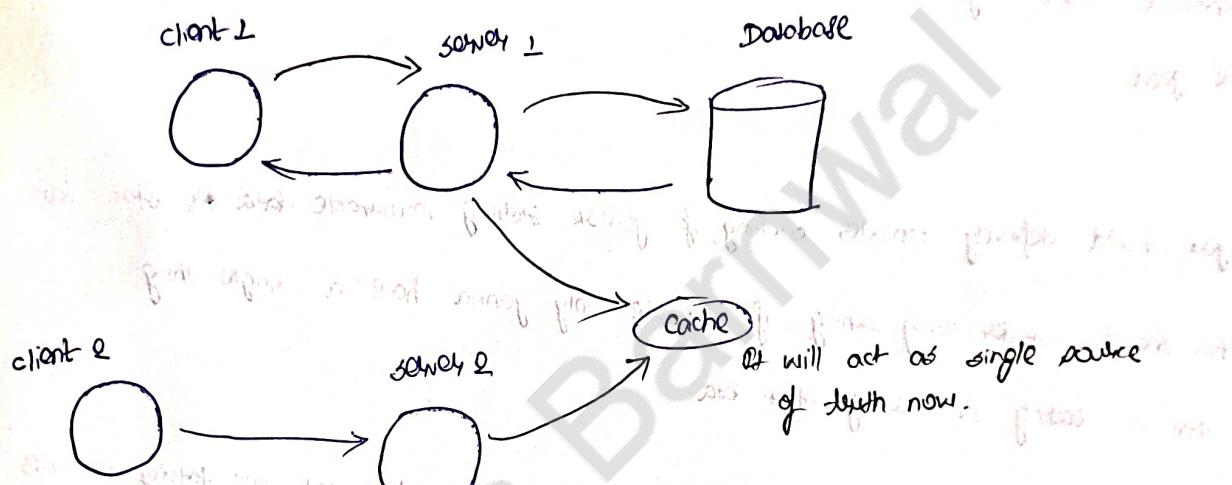
Note:

Caches can become stale if they have not been updated properly.

on the above ex, one of the soln is to move our cache out of the servers to and to put a single cache. And all servers hit this cache and then we'll have a (let's say, redis)

single source of truth for the caching mechanism.

IP.



→ For some part of the system, we might not care much about staleness or non-staleness of data in our caches.

Eg: Let's take viewcount on your YouTube video.

Viewcount is not necessarily the most important info on the YouTube video.

And if one user sees a slightly stale version of a viewcount on a video, that's not trouble much or that's not a big concern.

So, these things can be discussed with the interviewer in the system design interview that do we care about the accuracy of data that much?

Summary:

Caching has a lot of pitfalls.

- In general, if the data that you're dealing with is static data or immutable data, then caching is just beautiful.
- But if you're dealing with data that is mutable, then you have 2 or more locations where the data exists. You have to make sure that those locations are in sync. Otherwise data might become stale and depending on your use-case, that might not be good.

- ① So, you should definitely consider caching, if you're storing immutable data or static data.
- ② You should consider using caching, if you are only gonna have a single thing that is reading or writing that data.
- ③ If you don't care about consistency or staleness of data, then you can totally consider caching.
- ④ If you are able to design a system such that you can properly validate or get rid of stale data in your caches, then you can consider caching.