

O'REILLY

# Designing Data-Intensive Applications

THE TOP DATA INTEGRATION PLANS  
AND THE PRACTICAL APPROACHES  
TO BUILDING HIGHLY SCALABLE  
AND RELIABLE SYSTEMS

SECOND EDITION

# SYSTEM DESIGN INTERVIEW



By Alex Xu

Alex Xu

# TIPS for System Design Interview

Interview Questions and Answers  
System design and Expectations, i.e. functional, non-functional requirements

## ① Don't go into details prematurely

In fact, when it comes to DB design, ask if they want you to show ER diagram right away (initially, we also want to do that)

## ② Don't have an solid Architecture ready

i.e. not having global requirements and design as you've already read the question

## ③ keep it simple, stupid (KISS)

## ④ Don't make statements without justification

e.g. saying Join use in MySQL DB unless any relevant statement to back the design

## ⑤ Be aware of current technologies

(e.g. say you're use a CDN, it might be nice to mention a cool company like Akamai  
(or) cloud front

(e.g.)  
for a message system, say you're use maybe  
pusher.io / rabbitMQ /

② Don't do capacity estimation without any reason  
↳ magic speak out with interviewer which you  
will notice or presentation service needs to know  
what's your capacity when you're looking some-  
thing

③ Mention what protocols like communication of the  
systems in distributed DB (mention pros & cons of  
using one over other (ex. MySQL vs MongoDB))  
④ Speak out internal of system (working on MySQL  
DB like Cassandra over MySQL what does Consistency  
bring to the application other DB's connect)

## ⑤ TRADE OFFS (1)

# Terms & Analogy

Saturday, 25 December 2021 3:07 PM

① Scaling → Horizontal  
→ Vertical

② Load Balancing

③ Backups (Avoid single point of failure)

④ Microservices Architecture (Isolate Responsibilities)

⑤ Distributed System (Partitioning)

⑥ Decoupling

⑦ Logging & Metrics

⑧ Caching

⑨ API (Application programmable interface)

⑩ CAP Theorem (Consistency, Availability, Partition Tolerance)

⑪ CDN (Content Delivery Network)

⑫ DNS (Domain name system)

⑬ Ports (ex. HTTP - 80, HTTPS - 443)

⑭ Network protocols (IP, HTTP, TCP)

# Snippets

Friday, 14 January 2022 5:23 PM

Read heavy?

↳ Throw a Cache

Write heavy?

↳ Push write reqs to a message Queue.

Search Data?

↳ Put a search index

More users? Need to scale Database?

↳ Need ACID properties?

↳ Shard the RDBMS Database

↳ ACID Properties not mandatory?

↳ use NoSQL Database

High throughput?

↳ Scale application servers & put load balancer with  
consistent hashing

# Network Protocols

Wednesday, 12 January 2022

10:39 PM

IP  
TCP  
HTTP

Just a set of rules to interact b/w 2 things.

## IP (Internet Protocol)

$2^{16}$  bytes  $\approx 0.06$  MB

2 main things in IP packet are

① Header (src, target, total size of packet, IP version etc.)  
↓  
very small in general

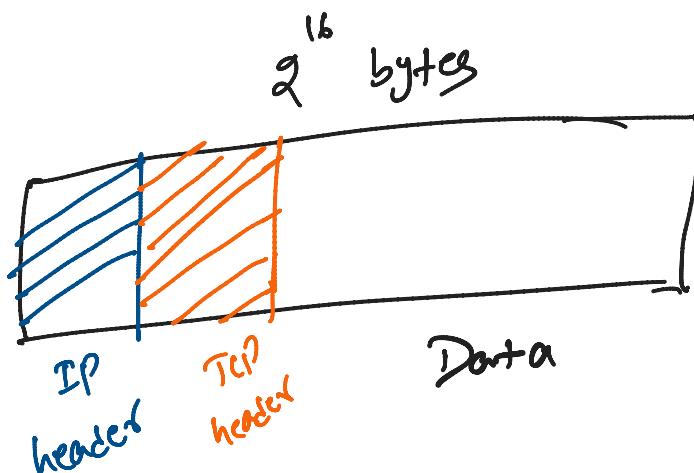
IPv4 / IPv6

② Data

## TCP (Transmission Control Protocol)

↳ Built on top of Internet Protocol

→ sends IP packets in order, reliable, error free way



## Hyper Text Transfer Protocol

## HTTP (Hyper Text Transfer Protocol)

↳ Built on top of TCP

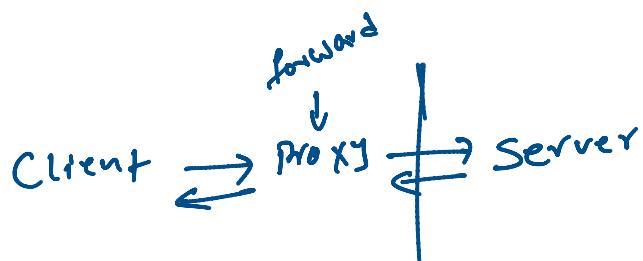
↳ This abstraction is called request-response paradigm.

# Proxy (on Behalf)

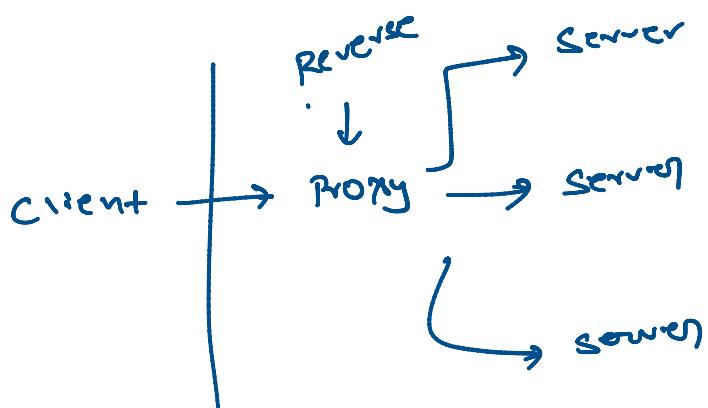
Friday, 31 December 2021 8:07 PM

Forward proxy → Server doesn't know which Client it is connecting from (VPN)

Reverse proxy → Client doesn't know which Server it is connecting to (balancing)



Forward proxy :- proxy sits b/w Client & Server and  
disguises client IP address, block malicious  
traffic from clients etc.



Reverse proxy :- It's basically a server side proxy.  
↳ can be used for traffic control, Load Balancing,  
caching response from servers etc...



## Expt. 7 Random

Random number generation

(↳ consistency)

(↳ availability)

(↳ partition tolerance)

Any real number system (continuous infinity) cannot  
have an integer modulus function. Therefore, the result can  
never be a random (uniformly distributed) function set of them.

Implementation:  
a) Naive Implementation: If we write  $\text{rand} \in [0, 1]$  and multiply it with  $n$ , we get a uniform distribution of numbers from  $[0, n]$ .  
However, this is not uniformly distributed because the floating point representation has a limited range.

Implementation: Every number  $x \in [0, 1]$  is system specific. In the case of floating point, if we request, then the guarantee of uniformity is lost.

Implementation: System will be responding to an input  $R$  which represents  
number of partitions. System will be responding to an input  $R$  which represents  
number of partitions (or subintervals) or subintervals (or subintervals).

3. Growth in performance many happens through reuse strategy. If we have distributed  
system, this is almost unavoidable. (i.e., we cannot give up on this property)  
System, this is almost unavoidable. (i.e., we cannot give up on this property)

7. Our design & system with cluster

(i) Partition tolerance + Availability (P)

(ii)

(iii) Partition tolerance + consistency (C)

## Design of Consistency & Availability

Ques: We know that we cannot provide both consistency and  
partition tolerance. How can we provide both consistency  
and availability with partition tolerance? In practice,  
we can provide some level of consistency and availability.  
In fact we can provide some level of consistency and availability  
but not both.

↳ This obviously depends on whether if we can tolerate  
it or not!  
Ques: If we tolerate no way we can be inconsistent, then what  
can we do? Because no way we can be inconsistent, it means  
we can't implement either (i) partition tolerance or  
consistency.

"In ATM machine, no way we can do --  
have to be 100% consistent alright. (or) Booking Tickets System.

2) Let's say, this time we're talking about No. of views  
for a youtube video. Here, you can forfeit the 100%  
consistency (The nodes will eventually reconcile & update  
anyways).

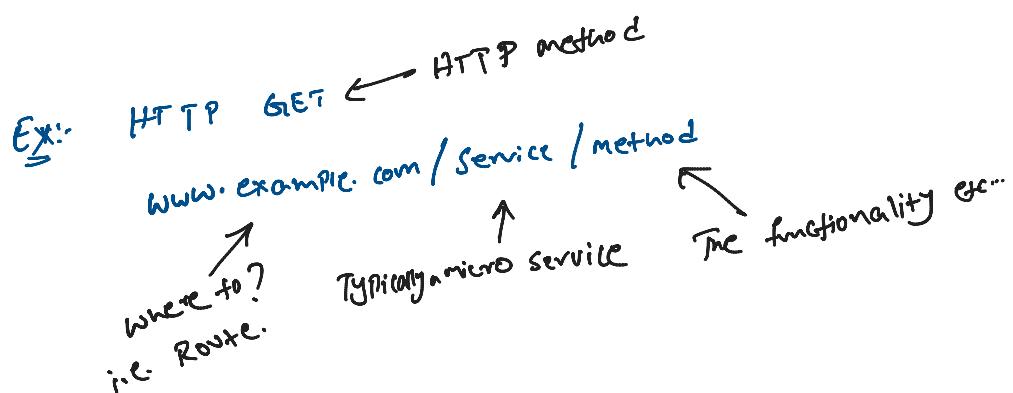
# API (It's only a contract!!!)

Tuesday, 28 December 2021 7:44 PM

- ↳ A documented way in which external consumers know how to interact with your code (they don't know how CCP works, they only know how to interact with it).

## Things to Remember when Designing API

- ① Naming is important
- ② Don't take additional parameters unless absolutely necessary. (to optimise the call.)
- ③ Don't give more info. than caller needs
- ④ Error handling → Response codes (HTTP) → Respond with correct HTTP codes.  
(Don't overdo it)  
Response messages → Don't give out sensitive info here!!



# Content Delivery N/W

Wednesday, 29 December 2021

8:47 PM

↳ example Akamai, S3

CDN does: (Primarily focuses on nearly fast data delivery by having multiple data centers geographically distributed !!)

- ① Host the data close to users
- ② Follow Regulations
- ③ Provides some interface to update/remove? data from the data centers?  
↳ The CDN basically?
- ④ can configure how much time to store some data in this CDN.

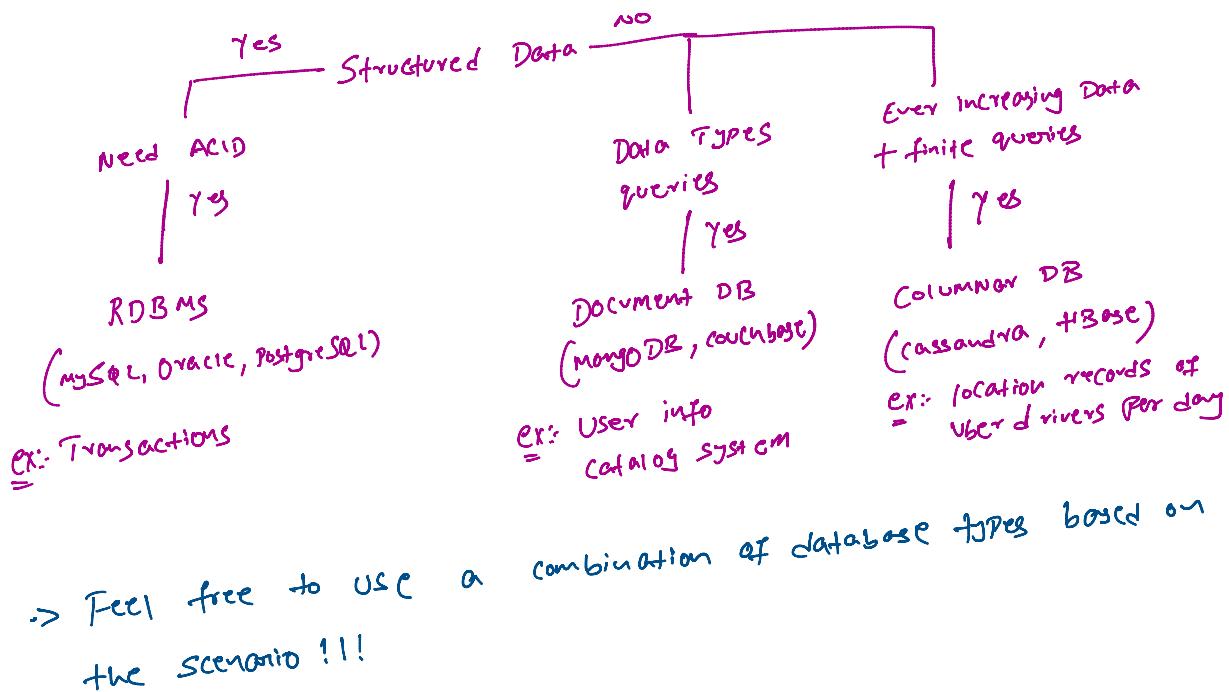
## Cloudfront

- ↳ super cheap
- ↳ reliable & easy to use

## SQL vs NoSQL

Monday, 24 January 2022 1:35 AM

- ① Always think of caching to reduce latency.
  - ↳ the server API can cache
  - ↳ Redis, Memcached (key value stores) → non relational DBs
- ② Storing Images, Videos? These are blob storages & not queried upon.
  - ↳ So use Amazon S3 to store
  - ↳ use a CDN in front of S3 to distribute same data geographically! (faster way)
- ③ Text searching feature (ex: Amazon Product etc Search / Netflix show title search etc..)
  - ↳ Elasticsearch / Solr



### ① Relational Databases

- ↳ tables & rows with a fixed schema
- ↳ ACID properties
- ↳ use case → Entirely structured data / need acid properties
- ↳ ex: Employee Table, Transactions etc..
- ↳ vertical scaling is easy (add more RAM, CPU etc.)
- ↳ horizontal scaling is tough here (i.e. sharding)

... databases (Horizontal Partitioning inbuilt mostly)  
... imp.)

↳ Built for scale

- ② non relational Databases (Horizontal Partitioning inbuilt mostly)
- ↳ key-value store (Redis, Memcached)
    - ↳ fast & provide quick access & they are in memory
    - ↳ used in caching solutions
  - ↳ Document Based (MongoDB)
    - ↳ used when not sure about structure of data & would evolve over time (i.e. No fixed schema, and can be changed easily)
    - ↳ supports heavy reads & writes, provides sharding capabilities
    - ↳ They have collections (like tables) and documents (like rows)
    - ↳ usecase → Product details of Amazon (since we want all data of one product at one place, no joins etc..)
  - ↳ Cons
    - ↳ NO ACID properties though we can handle it through application code, main loss is strong consistency.
    - ↳ Joins are tough to achieve.
  - ↳ column DB (Cassandra, HBase, Sybase)
    - ↳ sort of mid way of relational & document DBs
    - ↳ doesn't support ACID
    - ↳ usecase → streaming data / event data for analytics
    - ↳ distributed databases in general
  - ↳ Search Databases (ElasticSearch, Solr)
    - ↳ Full text Search
    - ↳ Fuzzy Search
    - ↳ Keep in mind, this is NOT PRIMARY DATABASE!

# For Reference

Tuesday, 28 December 2021 9:46 PM

The screenshot shows a mobile application interface with a title "Numbers Everyone Should Know". Below the title is a list of memory-related facts, each with a value and a unit. There are two red arrows drawn on the screen: one pointing to the first fact ("10 bytes reference") and another pointing to the last fact ("11m,930,000 ns").

10 bytes reference	0.5 ms
Branch mispredict	2 ms
12 cache references	7 ns
Matrix multiplication	163 ns
Memory refresh rate	100 ms
Processor L1 latency with 2-way	10,000 ns
Band of bytes from L1 cache to L2	10,000 ns
Round-trip memory access latency	250,000 ns
Memory copy within same processor	500,000 ns
Disk seek	
Read - 100 sequentially from memory	10,000,000 ns
Read - 100 sequentially from disk	10,000,000 ns
Read - 100 sequentially from memory (1GB)	11m,930,000 ns

# How to Avoid Cascading failures

Friday, 31 December 2021 1:47 PM

## Problems

- ① Cascading failures → Rate limiting
- ② Going viral → Pre scale
- ③ Predictable load increase → Auto scale
- ④ Bulk job scheduling → Batch processing
- ⑤ Popular posts → Approximate statistics.

## Database Searching

Section 10 Database Searching

Problem → How do query work when there is?

- ① Multiple queries (but not of same, even with different queries, it's not good enough)
- ② Queries in same (func. but not care of others...)

## ③ Searching

### 3. Horizontal Partitioning

→ Horizontal partitioning → breaks up data into different sets to take a key and partitions the data to different sets to (this key is generated on basis of data we're saving)

→ Practice of separating one function into multiple different tables, however by partitioning → searching

## Problems

- ① joins across multiple servers

- ② types of servers

- ③ types of horizontal partitioning

↳ the conventional partitioning → lets say each server handles → one shard is handle → in which shard we can again → multiple

one Shard :-  
Now, in this shard we can again  
shard / Partition the big shard into multiple  
shards.

② What happens if a shard goes down

↳ use master/slave architecture

↳ master is the most updated one &  
slaves keep reading from master for  
updates

↳ All writes go to master, reads are  
distributed to slaves

↳ If master fails, one slave becomes  
master.

## Distributed Caching

Version: 20 December 2021 - 10:00 AM

- ▷ Querying DB for company info about Louise Colee? (In memory cache)

Yields no serving info cache

- ▷ ~~Partial information given to client~~  
DB returns all company info every time  
↳ Client has to compare all data &  
then querying DB again once it's  
known in cache.

- ▷ This will avoid bad DB

comes generally from an SCD (batch data flow)  
comes generally from an SCD (batch data flow)  
↳ like job engine!!!

- ① When to load data into cache
- ② When to refresh data from DB

Cache policy!

- ▷ Cache policy determines how often performance is checked
  - ↳ O-LRU (least recently used) → mostly used
  - ↳ LRU (least frequently used) → never used

## Threading

- ↳ (currently) reading a memory zone from cache without ever using the registers!

## Problem

Given two dates. Is one date in future?

Given two dates. Is one date earlier than the other?

Given two dates. Is one date later than the other?

## Simple Rule Cache

Given two dates. Is one date (a) and date (b)

Is one date (a) later than date (b)?

For update to work in cache, then update (a)

↓  
This friend will work  
for date (a) because it  
will always be later than  
date (b).  
For update to work in cache,  
then update (b) because it  
will always be later than  
date (a).

## Simple Rule Cache

Given two dates. Is one date in future? (Simple rule)

Given two dates. Is one date earlier than the other?

Given two dates. Is one date later than the other?

Given two dates. Is one date later than the other?

## Geometrische Winkel

### Winkel Summung

In  $\triangle ABC$  ist  $\angle A + \angle B + \angle C = 180^\circ$   
 und ist  $\angle A + \angle B + \angle C + \angle D = 360^\circ$   
 dann  $\angle D = 180^\circ$

→ **Winkelsummenpostulat** mit der Winkelsumme  $180^\circ$

$$\text{Gesucht: } \angle C = ?$$

$$\begin{array}{l} \text{Gegeben: } \angle A = 60^\circ \\ \angle B = 70^\circ \end{array}$$

Winkelsumme ist  $180^\circ$

zu  $\angle C$  fehlt  $180^\circ - 60^\circ - 70^\circ$   
 also  $\angle C = 50^\circ$

und  $\angle C$  kann mit dem Winkelmaß  
 nach  $50^\circ$  abgelesen werden

→  $\angle C = 50^\circ$  ist die gesuchte Winkelsumme

### Winkelhalbierer

Nach Zeichnung eines Winkels (Winkelwinkelform)

die Winkelhalbierer sind

z.B. In  $\angle ABC$ , die Winkelhalbierer ist  $\angle B$   
 (durch  $B$  verläuft ein Strahl, der  $\angle A$  und  $\angle C$   
 gleich teilt)

→ Die Winkelhalbierer teilt den Winkel in zwei gleiche Teile  
 Beispiel: Abstand von zwei Sesseln kann  
 bestimmt werden, wenn die Winkelhalbierer  
 zwischen den beiden Sesseln auf einer  
 Linie steht.

(Dann steht es in einem Winkel von  $90^\circ$ )

Winkelwinkelform kann in einem Winkel  
 mit Winkelmaß gezeichnet werden  
 und die Winkelhalbierer ist gleichzeitig  
 die Stelle, an der diese Winkelmaße

### Satz (Von Thales) (Thales-Satz)

Winkelmaß  
 $\rightarrow 180^\circ$

Solution (use Consistent Hashing)

① Let's say our hash function outputs in  $[0, M-1]$  range

② Servers themselves have IDs also, we hash these IDs also and take remainder with M.

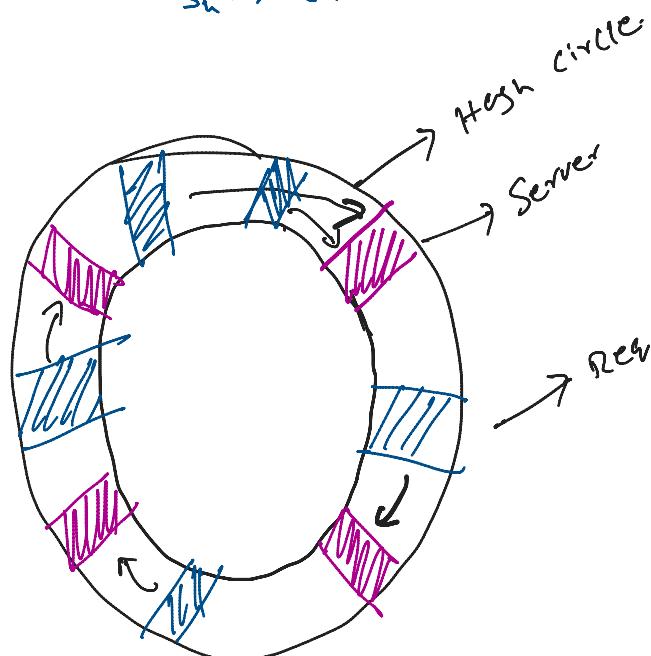
$$\text{i.e. } S_0 \rightarrow h(0) \% M$$

$$S_1 \rightarrow h(1) \% M$$

$$S_2 \rightarrow h(2) \% M$$

$$S_3 \rightarrow h(3) \% M$$

$$S_n \rightarrow h(n) \% M$$



⇒ The request is served by the immediate server on the clockwise direction.

⇒ If a new server is added, it goes into the hash ring.

⇒ If a server goes down, all requests that would go to this server will now be served by the next server on clockwise direction.

Problem

⇒ If there are 5 requests that have Server "A" on clockwise direction. If this server "A" goes down, then what happens?

- If the 5 requests go to next server on clockwise direction. If this server "A" goes down, then all these 5 requests go to next server on clockwise direction, lets say server "B". Due to this server "B" load increased a lot!!

SOLUTION:-

use multiple hash functions for servers

$$\text{I.C. } \begin{aligned} s_0 &\rightarrow h_1(0) & h_2(0) \rightarrow \dots \\ s_1 &\rightarrow h_1(1) & h_2(0) \rightarrow \dots \end{aligned}$$

- If we have "K" hash functions, we can have K points for each server.
- Likelihood of one server getting lot of load is now very less.
- If a server is removed, then remove the K points of that server.

# Message Queue

Sunday, 26 December 2021 3:34 PM

- Takes tasks, persists the tasks, assign tasks to right servers,  
waits for task completion, if taking long time/server dead  
assigns to another server

All of these are done by a  
"Message Queue"

Ex:- RabbitMQ, JMS (Java messaging service)

## Horizontal vs Vertical Scaling

Horizontal vs Vertical Scaling

→ Tech services are good (fixed cost per unit)



Can be your product/service / anything

Problem → You do benefit more replicability  
i.e. Scalability

- ① Buy bigger machine → horizontal scaling
- ② Buy more machines → horizontal scaling

## Drawbacks

Horizontal	Vertical
① need load balancing	① N/A
② redundant	② single point failure
③ major cell (IP) between servers	④ longer process (communication)
④ Data consistency issue	⑤ consistency
⑤ Starts well	⑥ maintains units

Solution

Horizontal scaling but each server is powerful enough