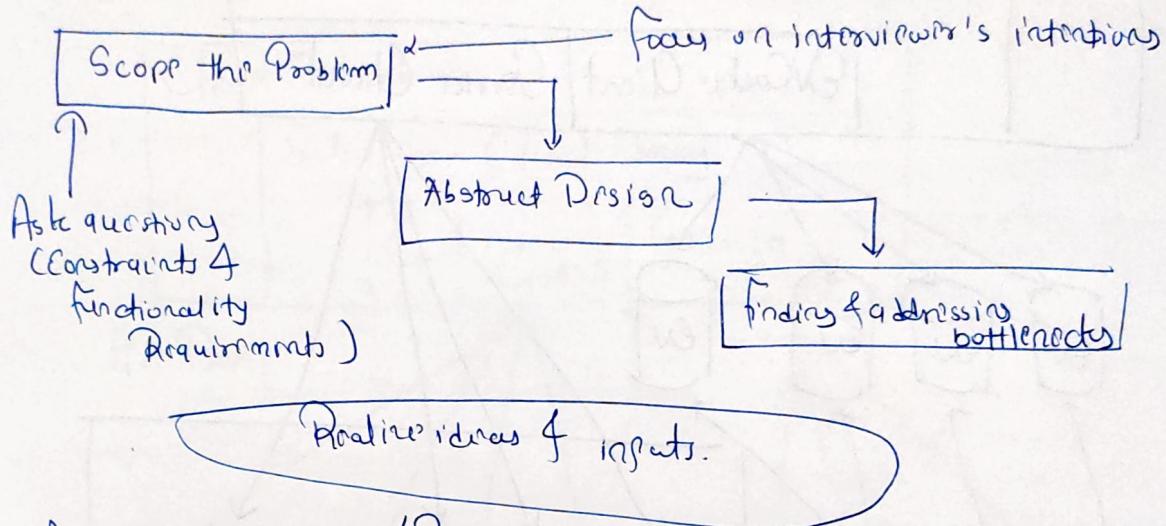


System Design Basics Amazon SDE-2

- ↳ Try to Break the Problem into simpler modules (Top down approach)
- ↳ Talk about trade-offs
(No solution is perfect)
Calculate the impact on system based on all the constraints and the end test cases.



- ↳ Architectural choices / Resources available
- ↳ How those Resources work together
- ↳ Utilization of Trade-offs

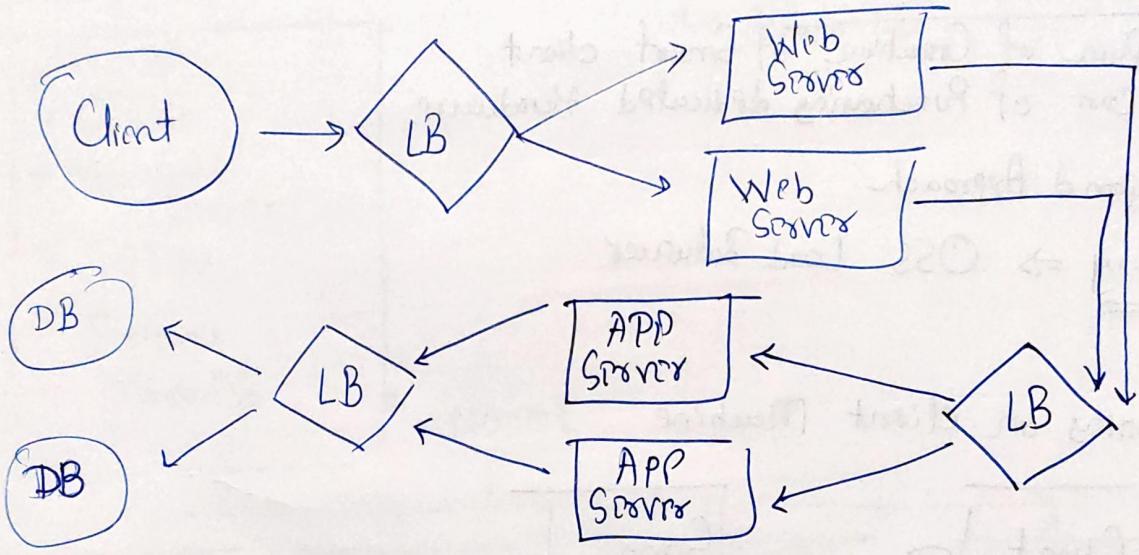
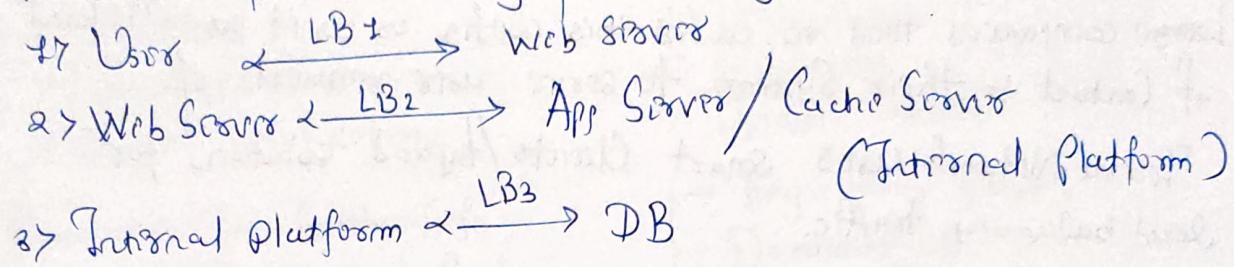
- Consistent Hashing
- CAP theorem
- Load balancing
- Queues
- Caching
- Replication
 - SQL vs NoSQL
 - In-drawer
 - Proxy
- Data Partitioning

Load Balancing (Distributed System)

Types of distribution

- Random
- Round-Robin
- Random (Weights for memory & CPU cycles)

To utilize full Scalability & Redundancy, add 3 LB.



Smart Clients

Takes a pool of Service hosts & balances load

↳ detects host that are not responsive

↳ Recovered Host

↳ Addition of new hosts

Load balancing functionality to DB (Cache, Service)

★ Attractive solution for developers.
(Small scale Systems)

As system grows → LBs (Standalone Server)

Hardware Load Balancers:

Expensive but high Performance.

e.g. Citrix NetScaler

Not trivial to Configur.

Large companies tend to avoid this config. as we it as 1st point of contact to their System to serve user requests & Intra Network uses Smart Clients/Hybrid solution for load balancing traffic.

Software Load Balancers:

No Pain of Creation of Smart client.

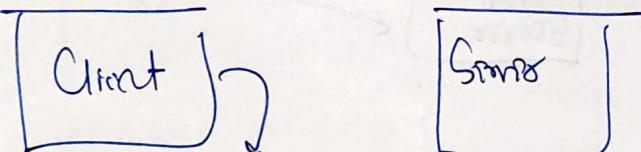
No cost of Purchasing dedicated Hardware.

↳ Hybrid Approach

HAProxy ⇒ OSS Load Balancer



⇒ Running on Client Machine



(locally bound port)

e.g. localhost:9000

↳ Managed by HAProxy
(With efficient management of Requests on the Port)

⇒ Running on intermediate Server: Proxies running b/w different Server side Components

HAProxy

- Manages health checks
- Removal & addition of Machines
- Balances requests a/c pools.

World of Databases

[SQL vs NoSQL]

Relational Database

- ↳ Structured
- ↳ Predefined Schema
- ↳ Data in rows & columns

Row → One Entity Info
 Column → Separate data points

- MySQL
- Oracle
- MS SQL Server
- SQLite
- PostgreSQL
- MariaDB

Non-relational Database

- ↳ Unstructured
- ↳ Distributed
- ↳ Dynamic Schema

- key-value store
- Document DB
- Wide column DB
- Graph DB

NoSQL

Key-value store

Data → Array of key-value pairs
 key → Attribute
 Value → to distinct

Redis
 Voldemort
 Dynamo

Document DB

Data → Documents
 ↴ grouped into
 Collections
 Each doc can be different.

CouchDB
 MongoDB

Wide-column DB

→ Instead of tables, column families
 ↴ Contains rows
 Non need of knowing all columns upfront.
 Each Row → different no of columns.
 Analysis of large datasets, etc.

Cassandra
 HBase

Graph DB

↓
 Data where Relations are best represented in Graphs
 → Nodes (Entities)
 → Properties (Information of Entities)
 → Links (Connection b/w Entities)

Neo 4J
 Infinite graph

High level differences b/w SQL & NOSQL

Property	SQL	NOSQL
Storage	Tables (Row → Entity, Column → Data Point) e.g. Student (Branch, ID, Name)	Dif. Data storage models (key value, Document, graph, columnar)
Schema	Fixed Schema (Columns must be decided f chosen before data entry) (can be altered ⇒ Modify whole Database (need to go offline))	Dynamic Schemas (columns addition on the fly) Not mandatory for each row to contain data.
Querying	SQL	UnQL (Unstructured Query Language) Queries focused on Collection of Documents. Diff. DB ⇒ Diff. UnQL
Scalability	Virtually Scalable (+ horsepower of hardware) Expensive Possible to Scale across multiple servers ↳ Challenging & time-consuming	Horizontally Scalable. Easy addition of Servers. Hosted on Cloud or cheap commodity hardware. → Cost Effective
Reliability or ACID Compliancy	ACID Compliant → Data Reliability → Guarantee of transaction → Still a better bet.	Sacrifice ACID Compliance for Scalability & Performance. (ACID) - Atomicity, Consistency, Isolation, Durability

Reasons to Use SQL DB

⇒ You need to ensure ACID Compliance:

ACID Compliance

- Resists anomalies
- Protects integrity of the Database.

for many E-commerce & financial Applications

↳ ACID compliant DB is the first choice.

⇒ Your data is structured & unchanging

If your business is not experiencing rapid growth or sudden changes

↳ no requirement of more servers

→ Data is consistent

then there's no reason to use System Design to support variety of Data & high ~~processing~~ traffic.

Reasons to use NoSQL DB

When all other components of system are fast

↳ Querying & Searching for data → bottleneck

NoSQL Document data format brings bottleneck.

Big Data → Large success for NoSQL.

↳ To store large volumes of Data (little/no structure)

No limit on type of Data.

Document DB → Stores all Data in one Place
(No need of type of Data)

⇒ Using cloud & storage to fullest.

Excellent cost saving solution. (Easy spread of Data across multiple Servers to Scale up)

⇒ Commodity hardware on site (affordable, smaller)

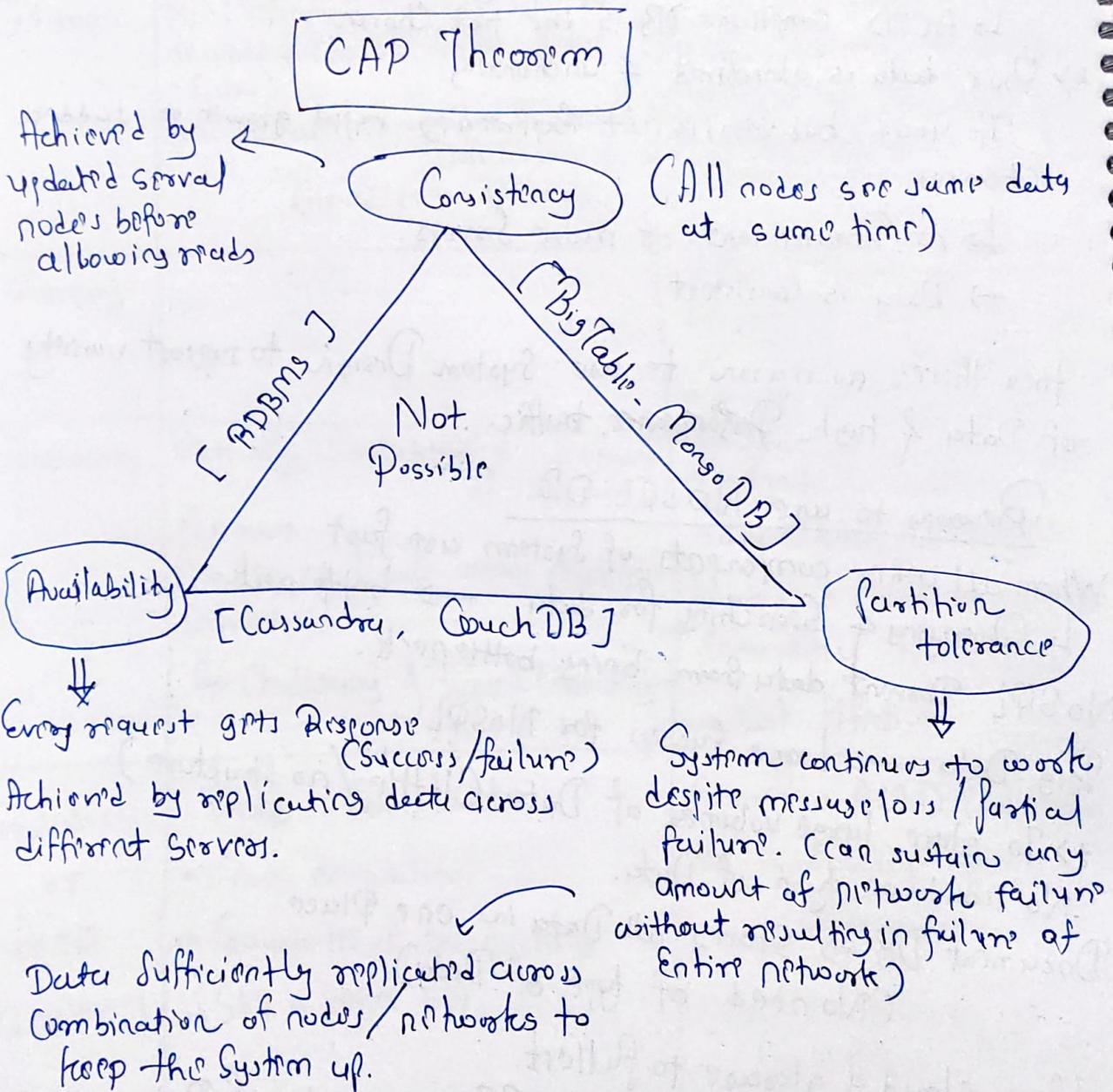
→ No headache of additional software

↳ NOSQL DBs like Cassandra → Designed to scale across multiple data centers out of the box.

3) Useful for rapid / agile Development

If you're making Quick iterations on schema

→ SQL will slow you down.



It is impossible for a distributed system to simultaneously provide more than two of three of the above guarantees.

We cannot build a distributed system which is:

- 1> Continually Available
- 2> Sequentially Consistent
- 3> Partition failure tolerant.

Because
To be consistent \rightarrow All nodes should see the same set of updates in the same order.

But if network suffers partition,
Updates in one Partition might not make it to other Partition.
 \hookrightarrow Client reads data from out-of-date Partition after having read from up-to-date Partition

Solution : Stop serving Requests from - out-of-date Partition.

\hookrightarrow Service is no longer 100% available.

Redundancy & Replication

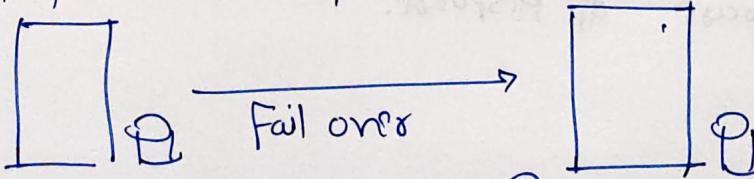
\rightarrow Duplication of Critical Data & Services.

\hookrightarrow increasing Reliability of System.

For critical Services & Data \Rightarrow Ensure that Multiple Copies / Versions are running simultaneously on different Services / databases.

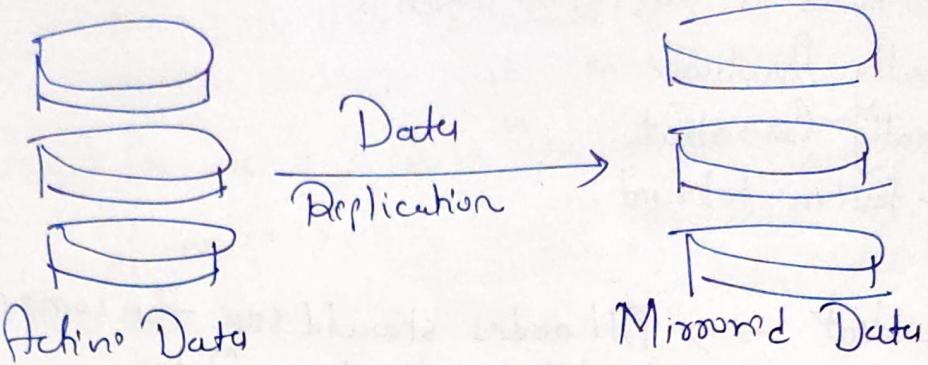
\hookrightarrow Secure against single node failures.

\rightarrow Provides backups if needed in Crisis.



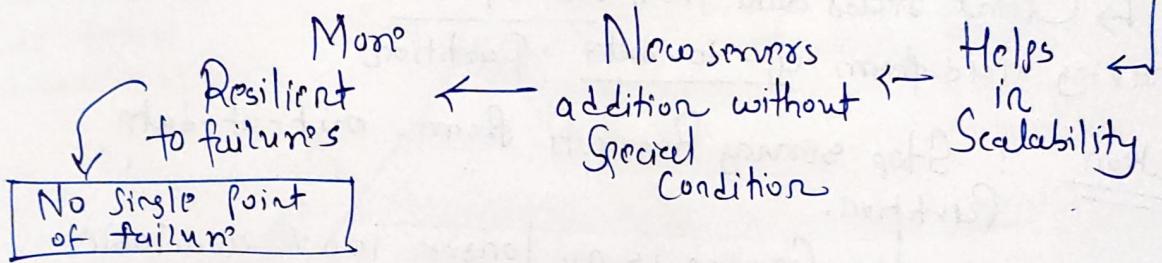
Primary Server

Secondary Server



Service Redundancy: Shared-nothing Architecture.

Every node → independent. No Central Service managing State.



Caching

Load balancing ⇒ Scales Horizontally

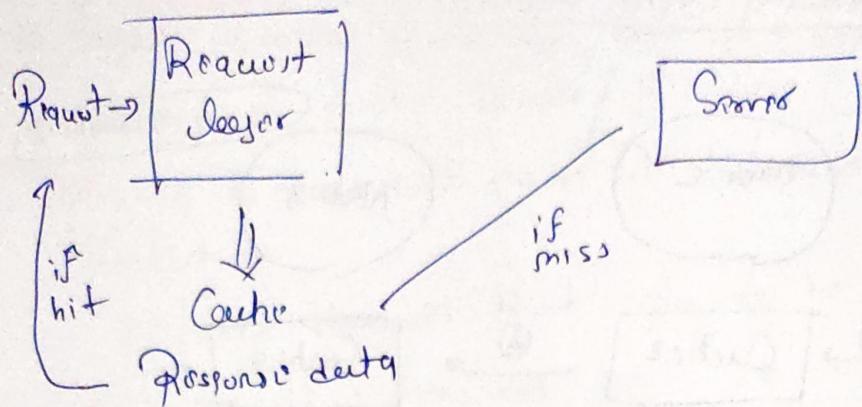
Caching: Locality of Reference Principle.

↳ Used in almost Every layer of Computing

⇒ Application Server Cache:

Placing a cache directly on a request layer node.

↳ Local Storage of Response.



Cache on one Request layer node

Can be located

Memory
(Very fast)

Node's local disk

(faster than going to network storage)

Bottleneck: If Load Balancing Distributes Request Randomly.

↳ Same Request \Rightarrow different Node

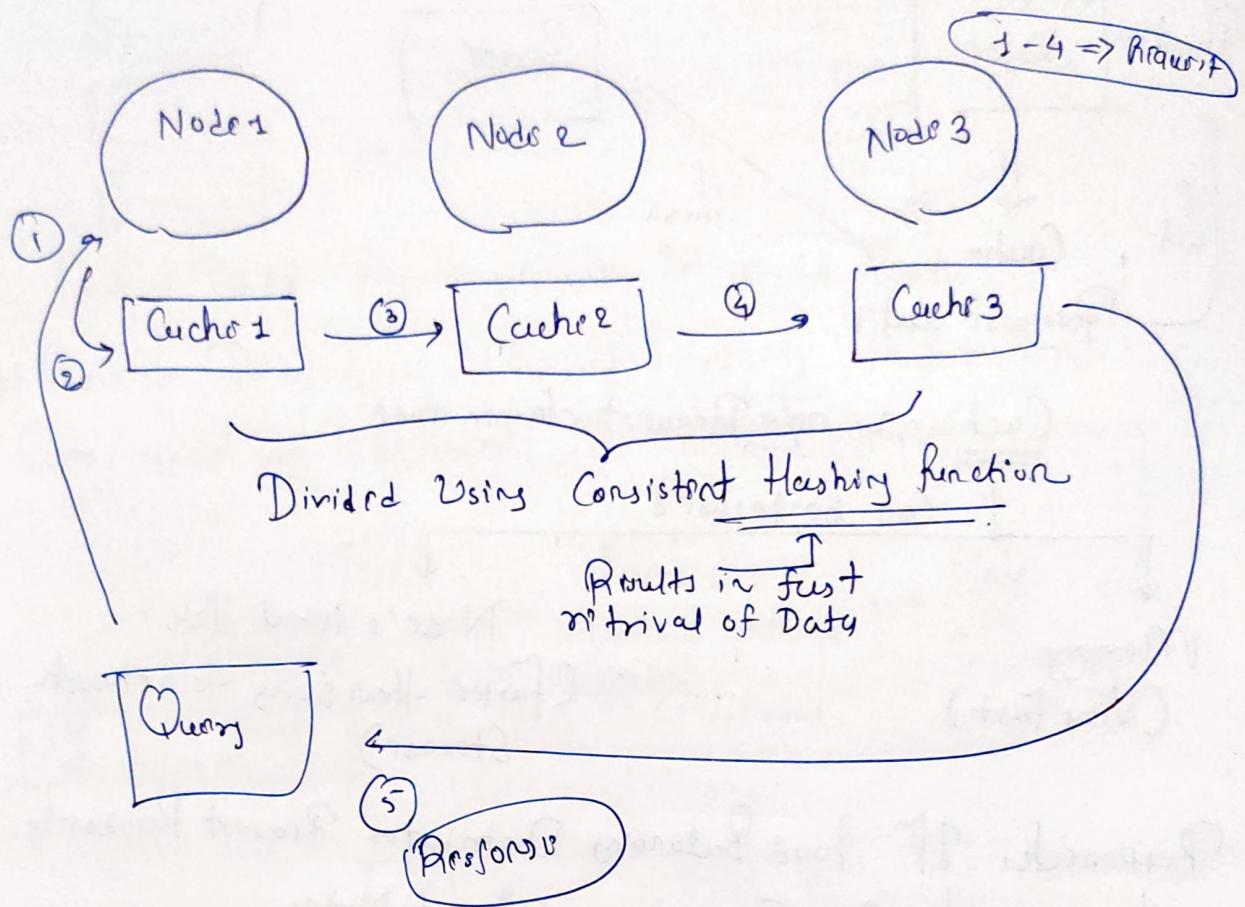
More cache Miss

\Rightarrow Global Caches

\Rightarrow Distributed Caches

Can be
overcome
by

Distributed Cache



Easy to increase cache by Space by adding more nodes

Disadvantage : Resolving a missing node

Storing multiple copies of data on different nodes) can be handled by

↳ When making it more Complicated .

Even if node disappears →

Processor can pull data from Origin

Single cache issues for all the nodes.

↳ Adding a cache server/file store (faster than original store)

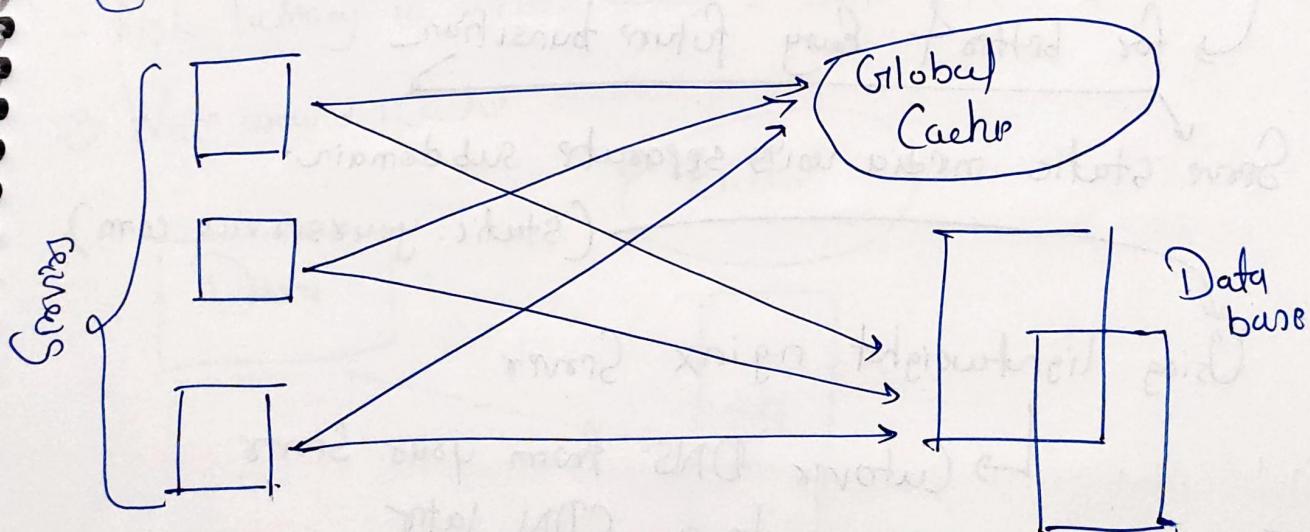
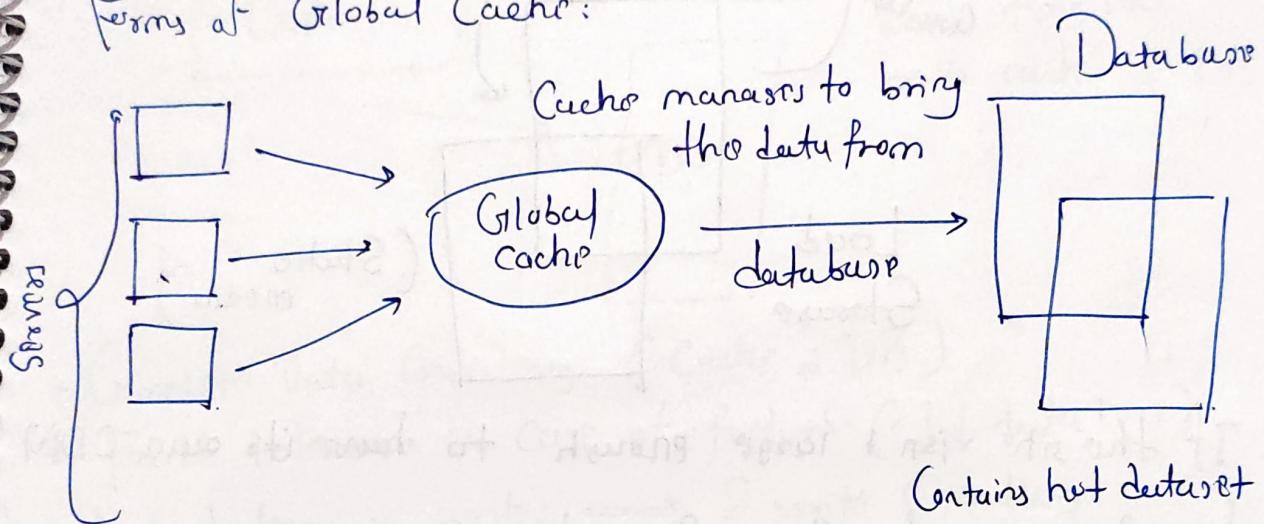
Difficult to manage if no of clients/requests increases

Effective if

 > Fixed dataset that needs to be cached

 > Special H/W \Rightarrow Fast I/O

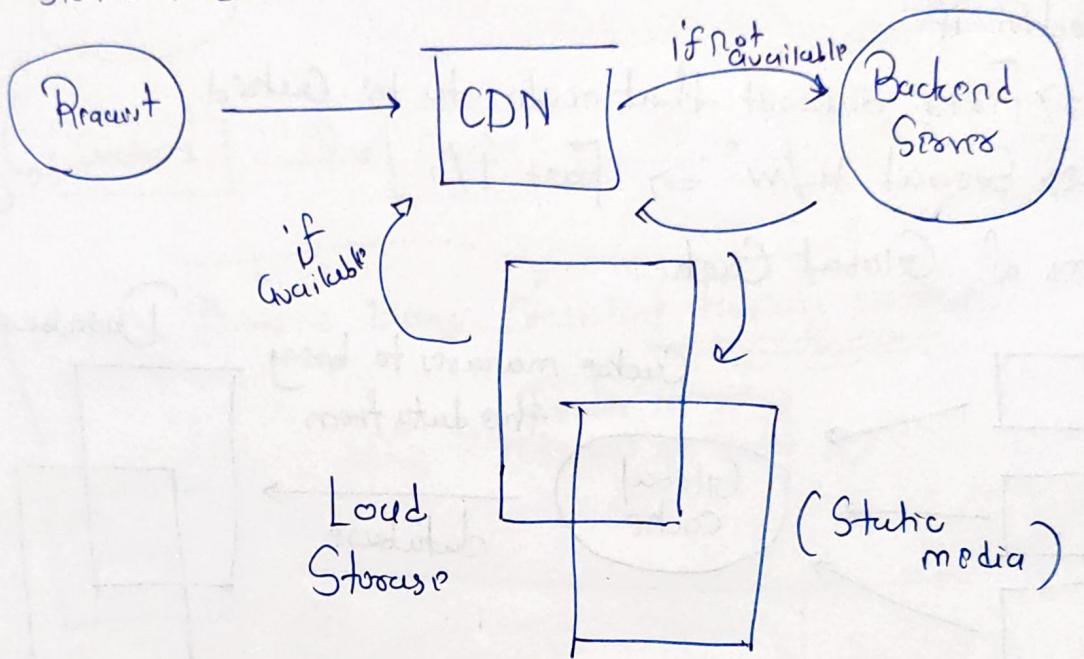
Forms of Global Cache:



Application logic understands the eviction strategy/hotspot better than cache.

CDN: Content Distribution Network

↳ Cache store for sites that stores large amount of static media.



If the site isn't large enough to have its own CDN.

↳ for better & every future transition

Some static media using separate sub domain

(static.yourservice.com)

Using lightweight nginx Server

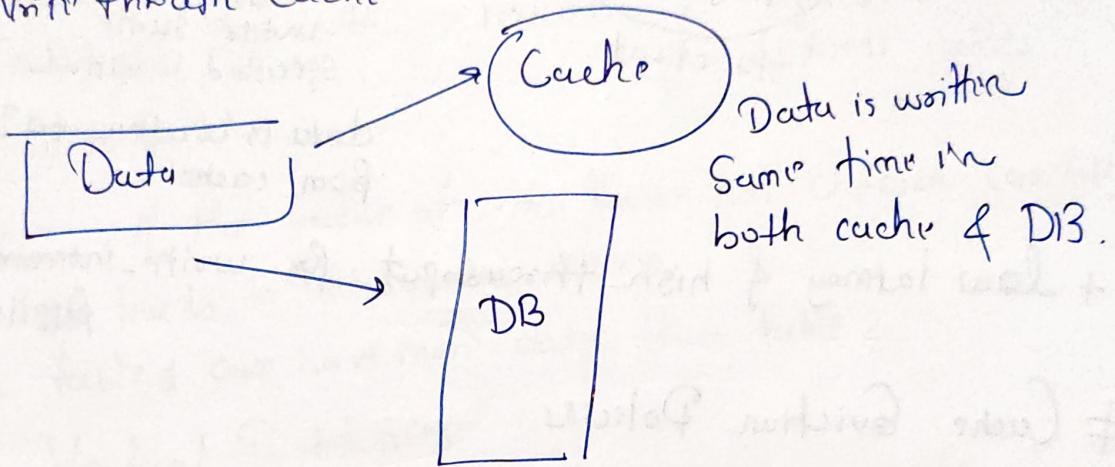
↳ Cutomize DNS from your Server
to a CDN later

Cache Invalidation

Cache \rightarrow needs to be coherent with this database if data in DB modified \Rightarrow Invalidates this cache data.

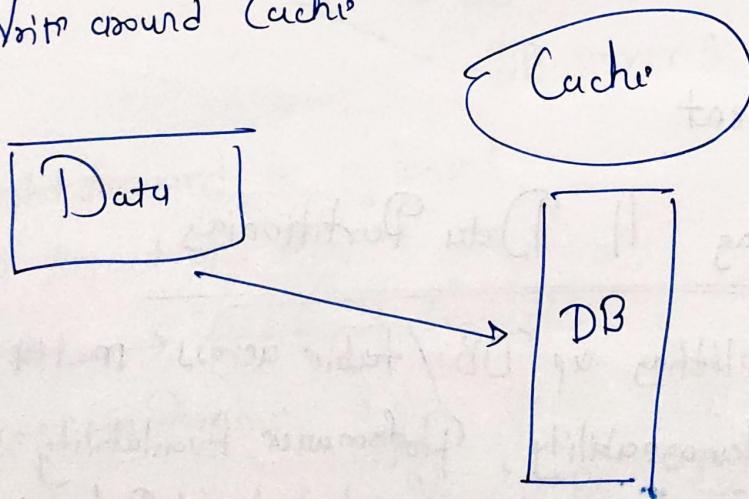
3 Schemes

1) Write through Cache



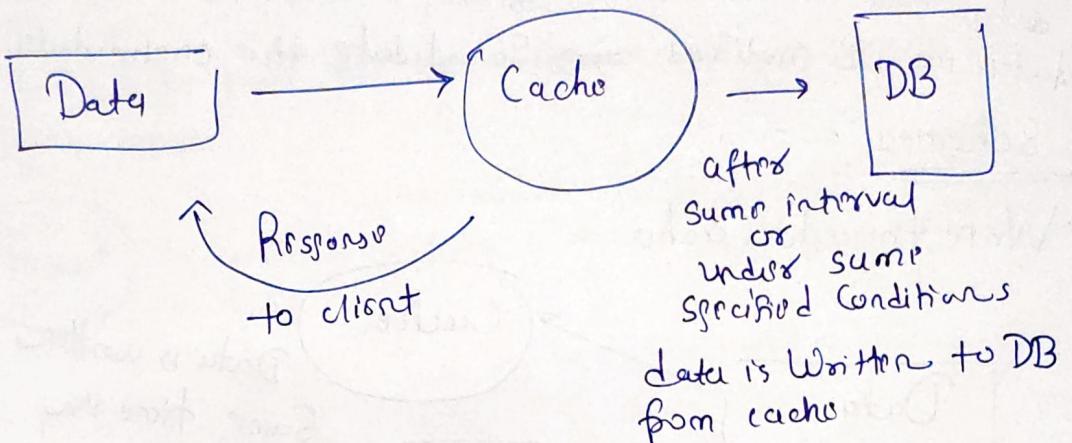
- + Complete Data Consistency ($\text{Cache} = \text{DB}$)
- + Fault tolerance in case of failure ($\downarrow\downarrow$ data loss)
- high latency in writes \Rightarrow 2 write Operations

2) Write around Cache



- + no cache flooding for writers
- Read Request for newly written data \Rightarrow miss
higher latency

3) Write back cache:



+ low latency & high throughput for write-intensive Application

Cache Eviction Policies

- 1) FIFO
- 2) LIFO or FILO
- 3) LRU
- 4) MRU
- 5) LFU
- 6) Random Replacement

Sharding || Data Partitioning

Data Partitioning: Splitting up DB / table across multiple machines \Rightarrow Manageability, Performance, Availability & Load Balancing

After a certain scale point, it is cheaper and more feasible to scale horizontally by adding more machines instead of vertical scaling by adding bigger servers.

Methods of Partitioning:

1) Horizontal Partitioning: Different rows into different tables
dense based Sharding.

E.g. Storing locations by Zip

Table 1: Zips with < 100000

Table 2: Zips with > 100000

different servers in different tables

and so on

Cons: If the value of the Range not chosen carefully

leads to unbalanced servers.

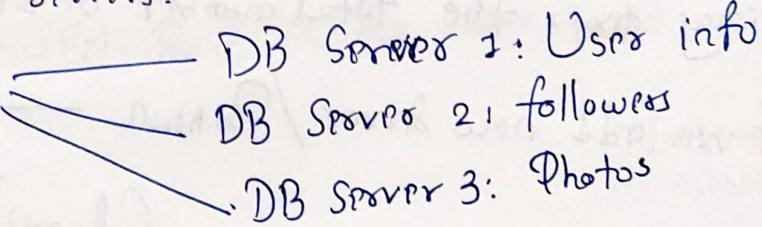
E.g. table 1 can have more data than table 2.

Vertical Partitioning

Feature wise distribution of data

in different servers.

E.g. Instagram



Straight forward to implement

low impact on Apps.

if app \rightarrow additional Growth

need to Partition feature specific DB across various servers.

i.e., if would not be Possible for a Single server to handle all meta data queries for 10 million Photos by 140 million users.

Directory Based Partitioning

Directory Based Partitioning

↳ A loosely coupled approach to work around issues mentioned in above two partitionings.

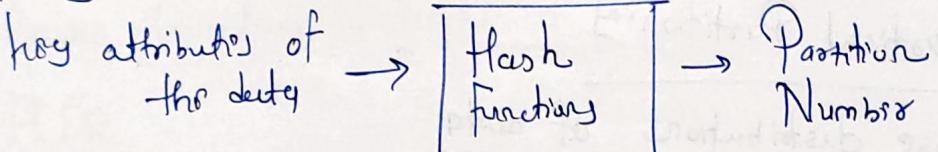
Create lookup Service \rightarrow Current Partitioning Scheme 4
Abstract it away from the DB access code.

Mapping (tuple key \rightarrow DB Server)

Easy to add DB servers or change Partitioning Scheme.

Partitioning Criteria

↳ Key or Hash based Partitioning:



Efficiency fixes the total number of servers / partitions

So if we add new servers / partition \downarrow

Change in hash function

downtime because of Redistribution

Solution:
Consistent Hashing

2) List Partitioning:

Each Partition is assigned a list of Values.

new record \rightarrow Lookup for key \rightarrow Store the Record

(Partition based on the key)

3) Round Robin Partitioning:
Uniform Data Distribution

With 'n' partitions

→ the $(i,)$ tuple is assigned to Partition
 $(i \bmod n)$

4) Composite Partitioning:

Combination of above Partitioning Schemes

Hashing + List → Consistent Hashing

↓
Hash reduces the key space to a size
that can be listed

Common Problems of Sharding

Sharded DB: Extra Constraints on the different DB Operations

↓

Operations across multiple tables or
multiple rows in the same table

↓

No longer running
in single server.

⇒ Joins & Denormalization:

Join on table on single Server → Straight forward

not feasible to perform joins on sharded tables

↳ less efficient (data needs to be compiled from
multiple servers)

Workaround → Denormalize the DB.

So that the queries that previously required joins can be performed from a single table.

Cons: Penalties of Denormalization

↳ Data inconsistency

Referential Integrity: Foreign keys on Sharded DB
↳ Difficult

Most of the PDBMS does not support foreign keys on sharded DB.

If an application demands Referential integrity on sharded DB

↳ Enforce it in application code (SQL joins
to clean up dangling References)

3) Rebalancing:

- ↳ Reasons to change sharding Schema:
 - a) Non-Uniform distribution (Data wise)
 - b) Non-Uniform load balancing (Request Wise)

Workaround: ↳ Add new DB

↳ Rebalance.

↳ Change in Partitioning Schema

- ↳ Data movement
- ↳ Downtime

We can use directory-based Partitioning

↳ Highly Complex

↳ Single point of failure
(Lookup service/table)

Indexes

→ Well known because of databases.

→ Improves speed of Retrieval.

→ Increased storage overhead

- Slower writes

↳ Write thru cache

↳ Update the index

→ Can be created using one or more Columns.

* Rapid Random lookups.

& Efficient access of ordered Records.

Data Structure

Column → Pointers to Whole Rows

→ Create different views of the same data.

↳ Very good for filtering / Sorting of large data sets.

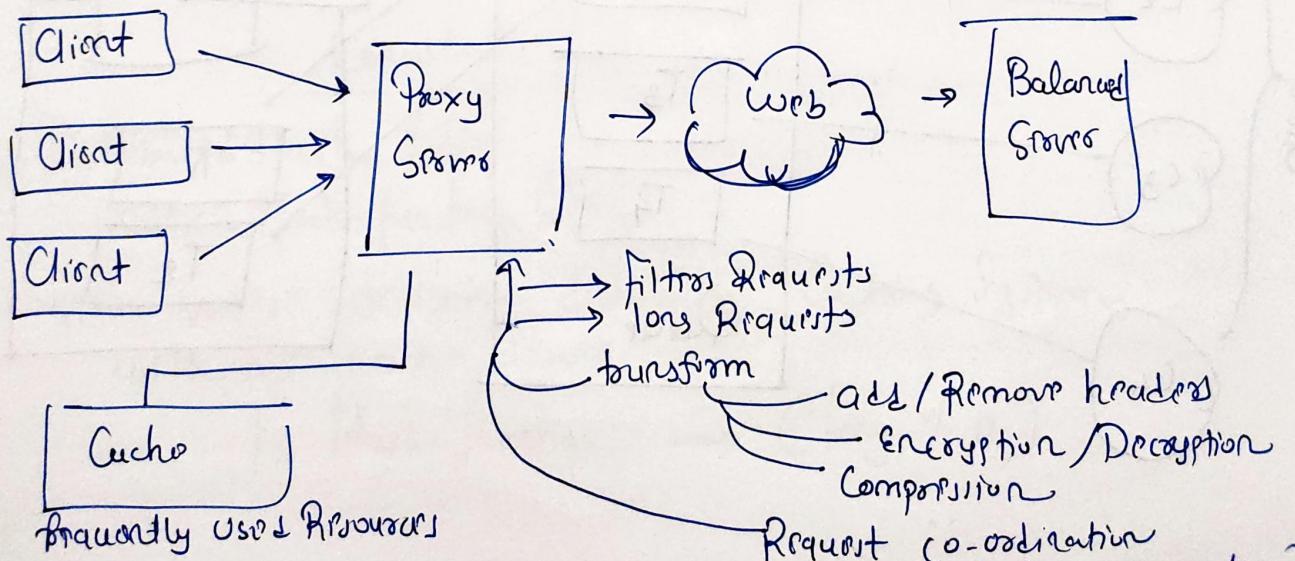
↳ no need to create additional copies.

Used for datasets (TB in size) of Small Payload (KB)

Spreading
Several Physical
Devices

→ We need some way to find the correct
Physical location. i.e. **Indexes**

Optimal under high load situations
if we have limited Caching
↳ batches several Requests into one



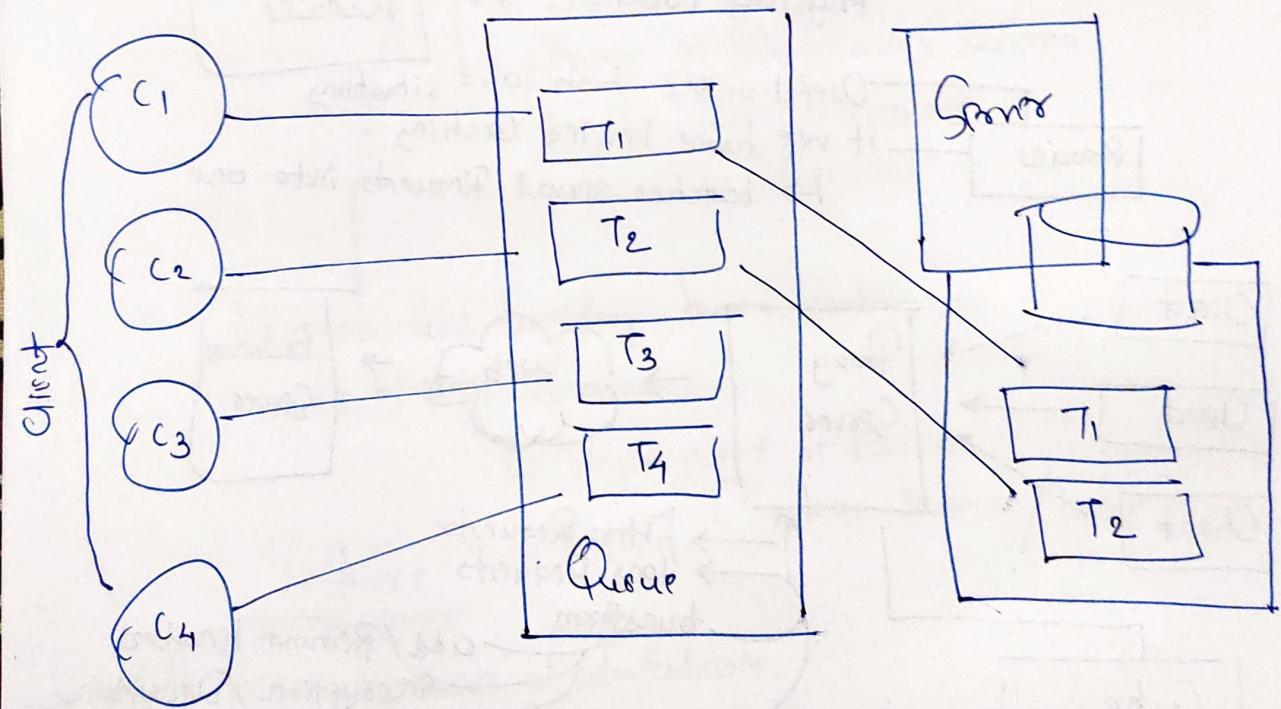
We can also use spatial locality

↳ Collapsing Requests for data
that is spatially close

Queues

Effectively manages Requests in Large-scale distributed System.

- In Small Systems → Writings can fail.
 - In complex Systems → high incoming load
if individual writes take more time.
 - * To achieve high performance & availability
 - ↳ System needs to be Asynchronous
 - ↳ Queues
- # Synchronous behaviour → Degrades Performance
- ↓
Can improve load balancing
- Difficult for
fair & balanced distribution



Queues: Asynchronous communication Protocol.

↳ Client sends task

↳ Gets ACK from queue (Receipt)

It stores as reference for the result in future

↳ Client continues its work.

limit on the size of Request

& Number of Request in Queue.

Queues: Provides fault-tolerance

↳ Protection from Service outage / Failure.

↳ Highly Robust

↳ Retry Failed Service Request

↳ Enforces Quality of Service guarantee

(Does not expose clients to outages)

Queues: Distributed Communication

↳ Open source implementation

↳ RabbitMQ, ZooMQ, ActiveMQ, BeanstalkD

Consistent Hashing

Distributed hash table.

index = Hash - Function (key)

Suppose we're designing distributed Caching System
with n cache servers

↳ Hash-function \rightarrow (key % n)

Drawbacks:

- 1) Not Horizontally Scalable
 - ↳ addition of new servers results in
↳ need to change all Existing Mapping
(downtime of system)
- 2) Not load Balanced
 - (Because of non-uniform distribution of data)

Some Caches: Hot & Saturated
Other Caches: Idle & Empty

How to tackle above Problems?

↳ Consistent Hashing

What is consistent hashing?

- Very useful strategy for distributed Caching & DHJs.
- Minimizes Data Reorganization in Scaling up / Down.
- Only $\lceil k/n \rceil$ keys needs to be remapped.

$k \rightarrow$ total number of keys

$n \rightarrow$ Number of Servers

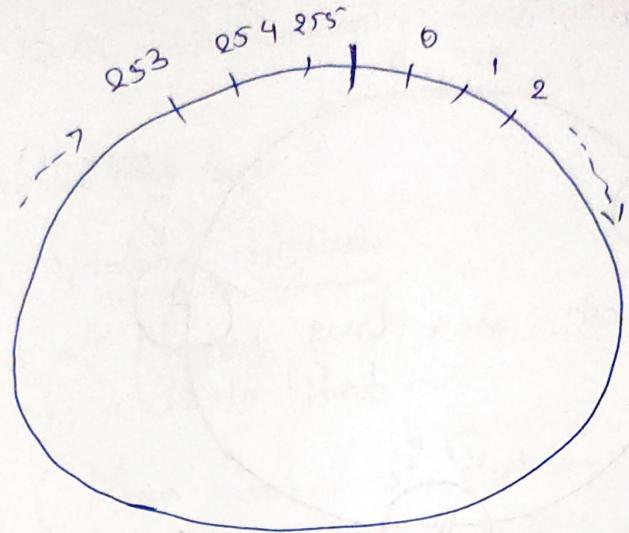
How it Works?

Typical hash function suppose outputs in $[0, 256]$

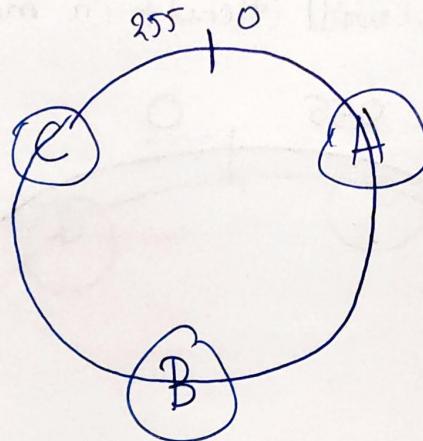
In Consistent Hashing.

Imagine all of these integers are placed on a Ring.

If we have 3 Servers: A, B & C



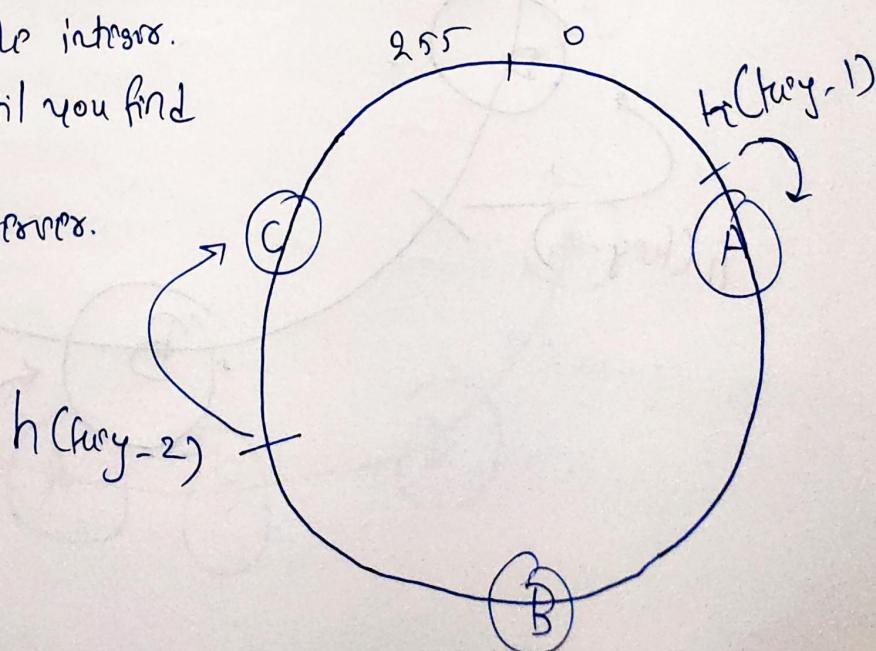
1> Given a list of Servers, hash them to integers in the Range.



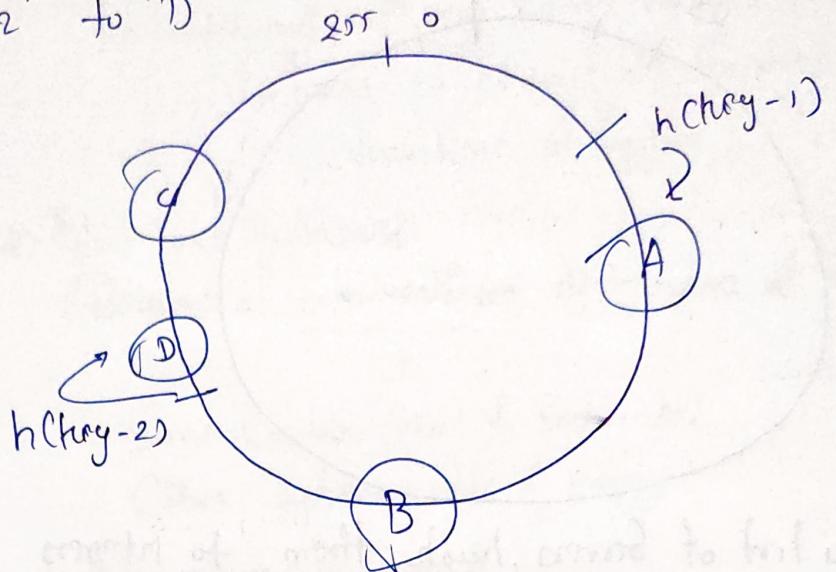
2> Map key to a server :

- Hash it to Single integer.
- Move CLK wise until you find Server.

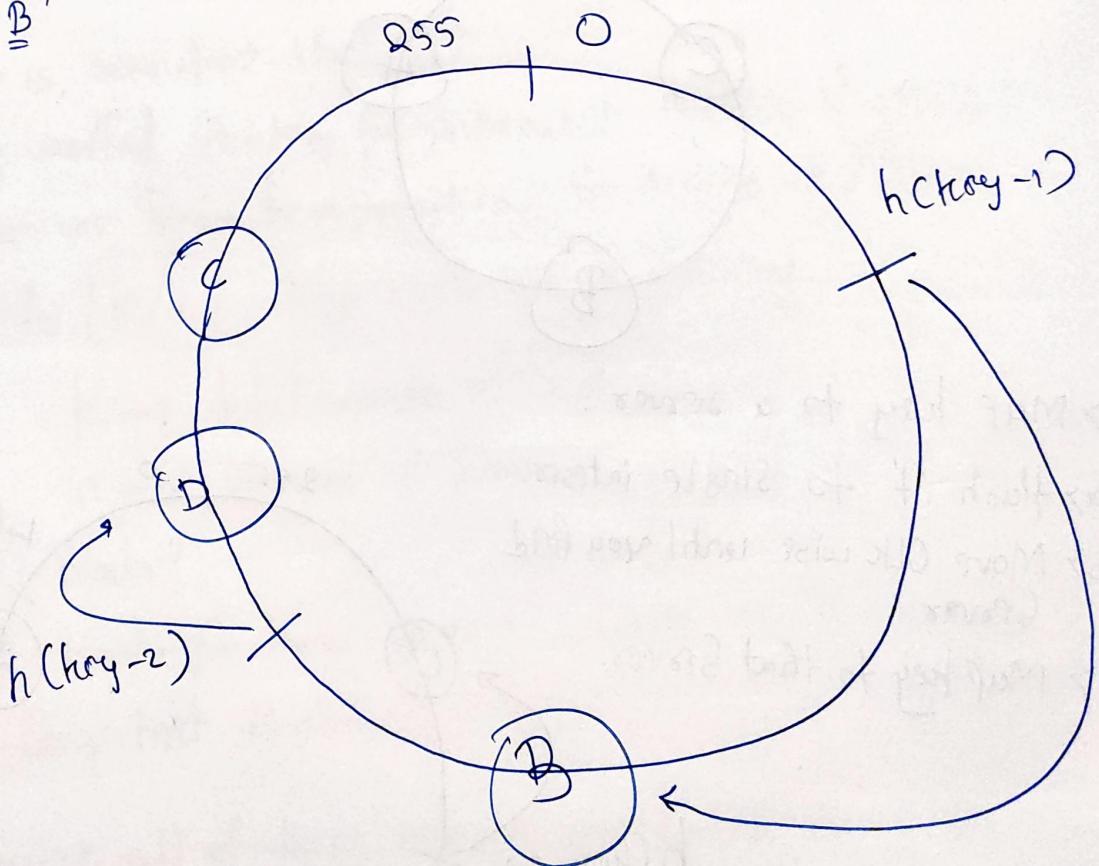
c> Map key to that Server.



Adding a new Scores 'D', will Result in moving the 'key-2' to 'D'



Removing Scores 'A', will Result in moving the 'key-1' to 'B'

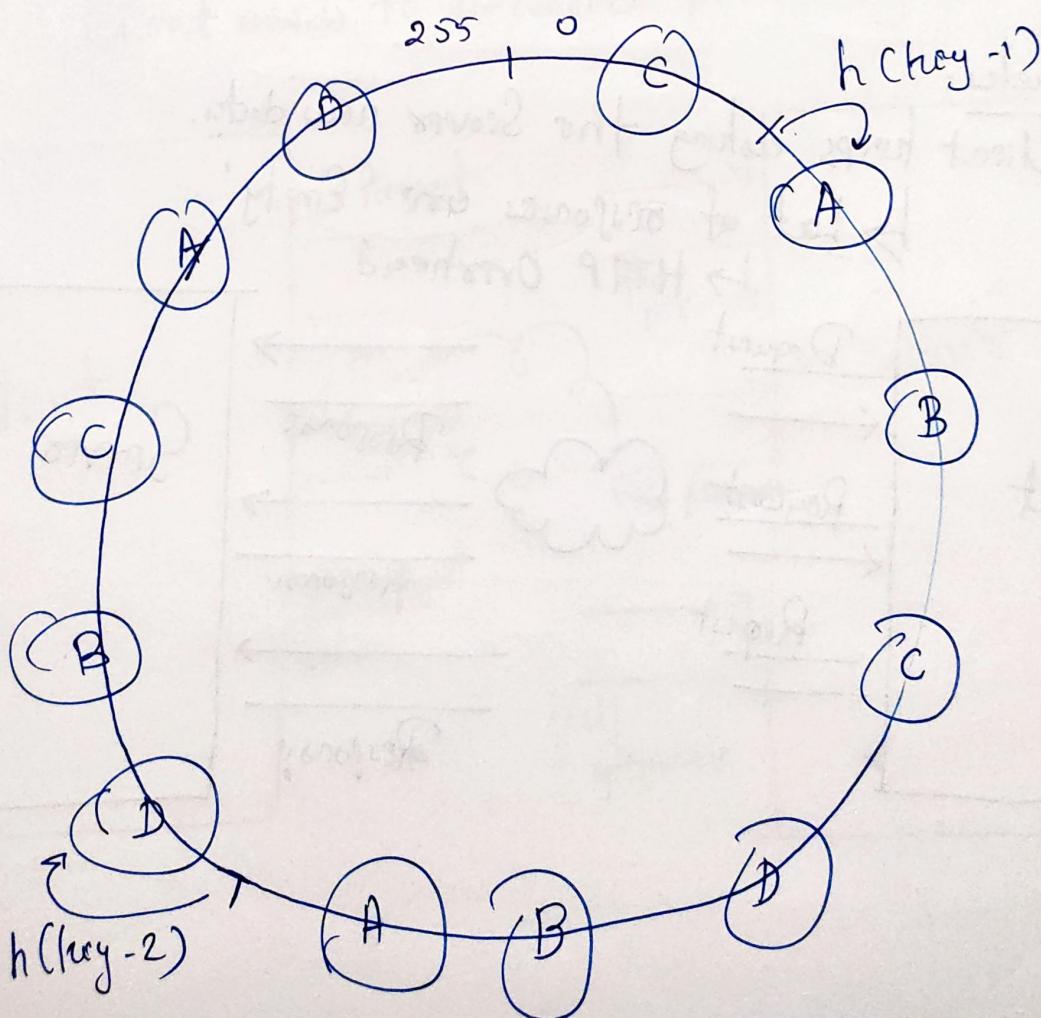


Consider Real World Scenario
data \rightarrow Randomly distributed
 \hookrightarrow Unbalanced Cache.

How to handle this issue?

Virtual Replicas

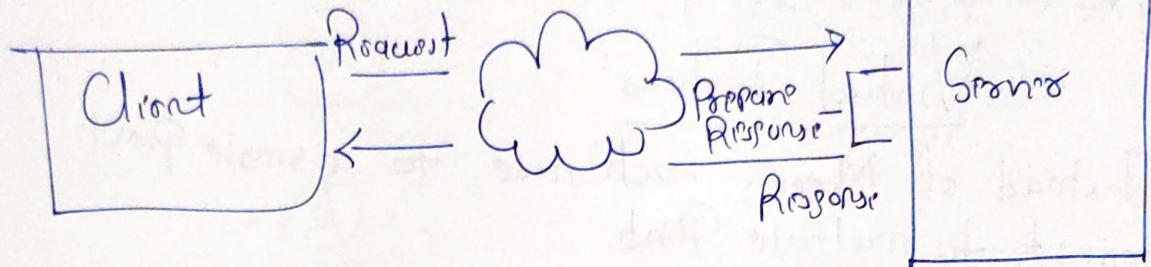
- \rightarrow Instead of Mapping such node to a single Point we map it to multiple Points.
 \hookrightarrow (More number of Replicas)
 \hookrightarrow More equal Distribution
 \hookrightarrow Good load Balancing)



Long Polling v/s Web Sockets w/ Server Events

Client Server Communication Protocols

HTTP Protocol :



#AJAX Polling:

Clients Repeatedly polls servers for data.

Similar to HTTP Protocol.

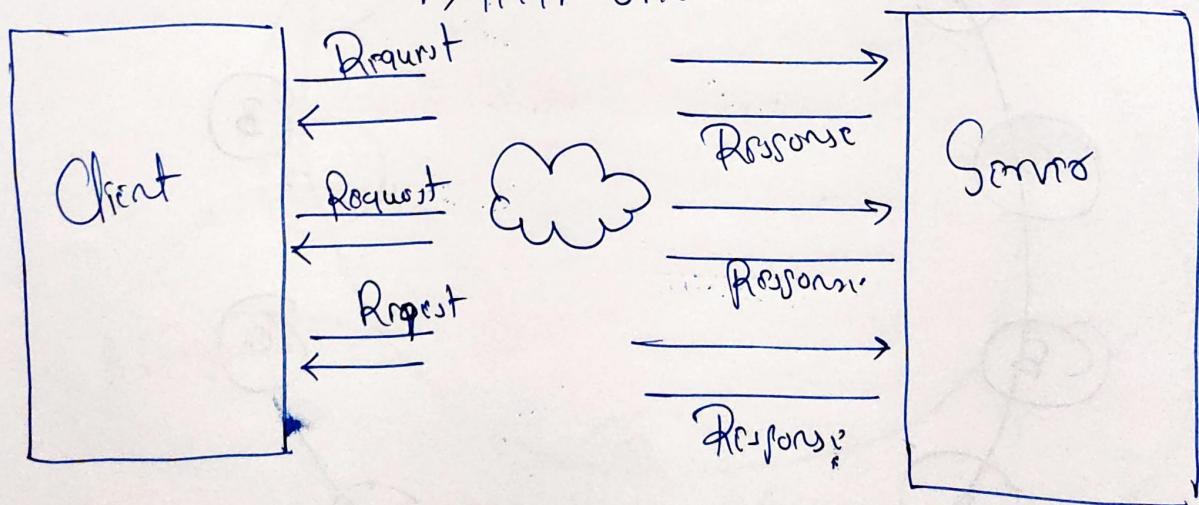
↳ Requests sent to Server at Regular intervals.
(0.5 sec)

Drewbooks:

Client keeps asking the Scrum now what's

↳ List of responses are 'Empty'

↳ HTTP Overhead



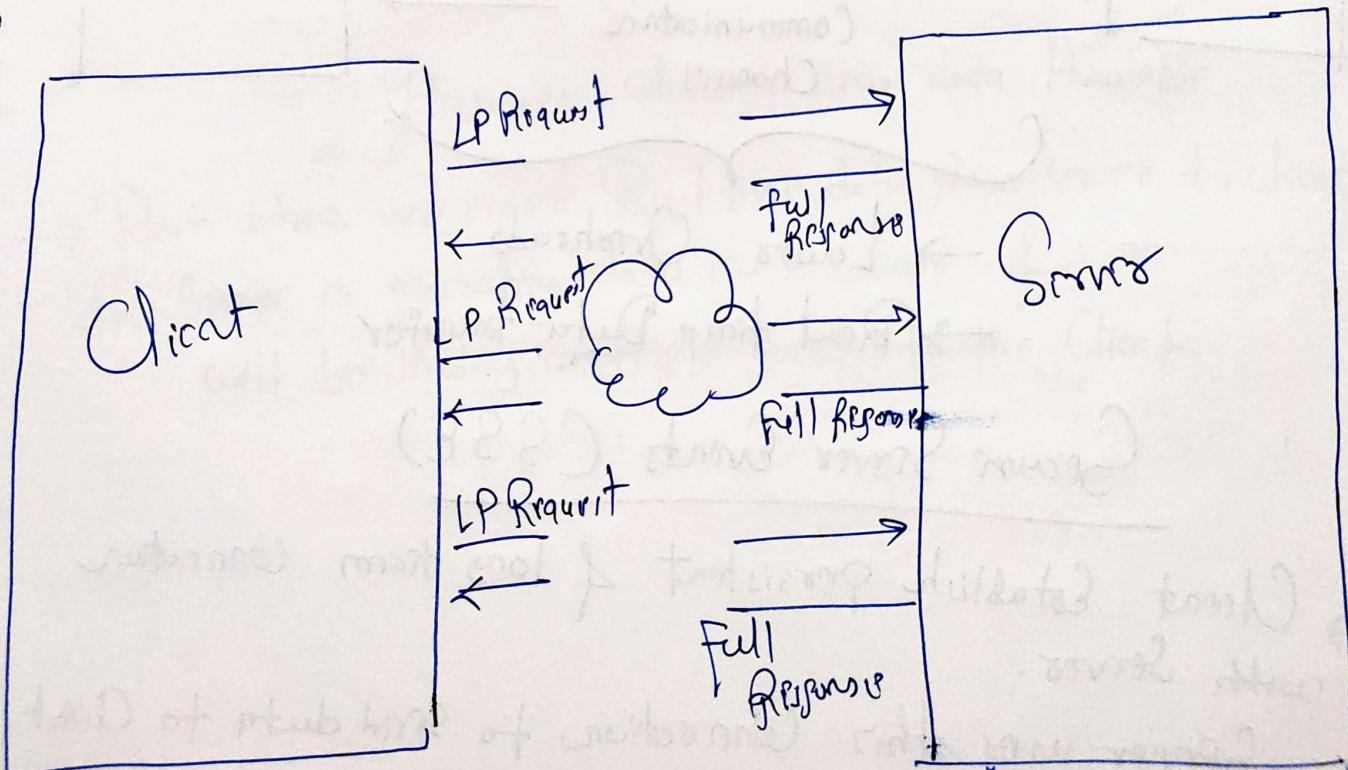
HTTP Long Polling: Hanging GET

Server does not send Empty Response.

Pushes response to clients only when new data is available.

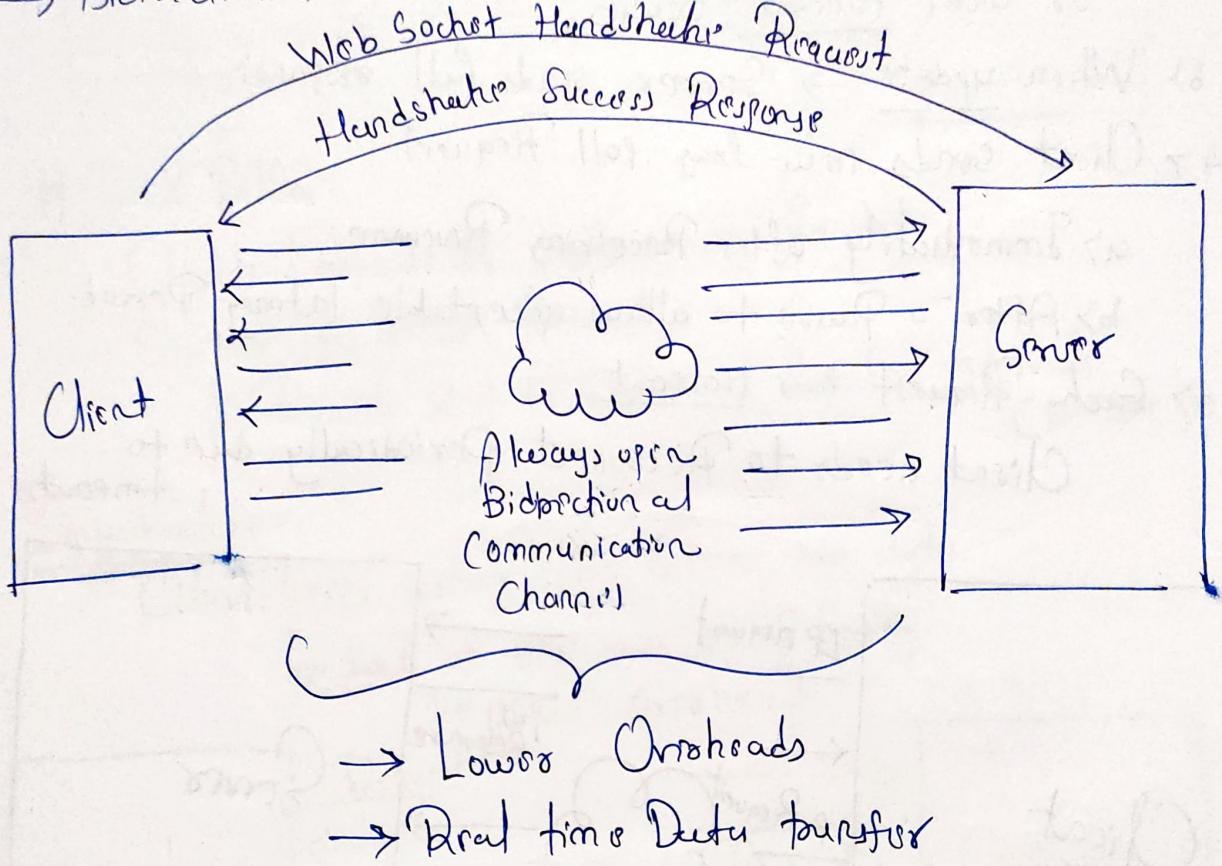
- 1) Client makes HTTP Request & waits for the Response.
- 2) Server delays response until update is available.
or until timeout occurs.
- 3) When update → Server sends full response.
- 4) Client sends new long-poll Request.
 - a) Immediately after Receiving Response.
 - b) After a Pause to allow acceptable latency Period.
- 5) Each Request has timeout.

Client needs to Reconnect Periodically due to timeout.



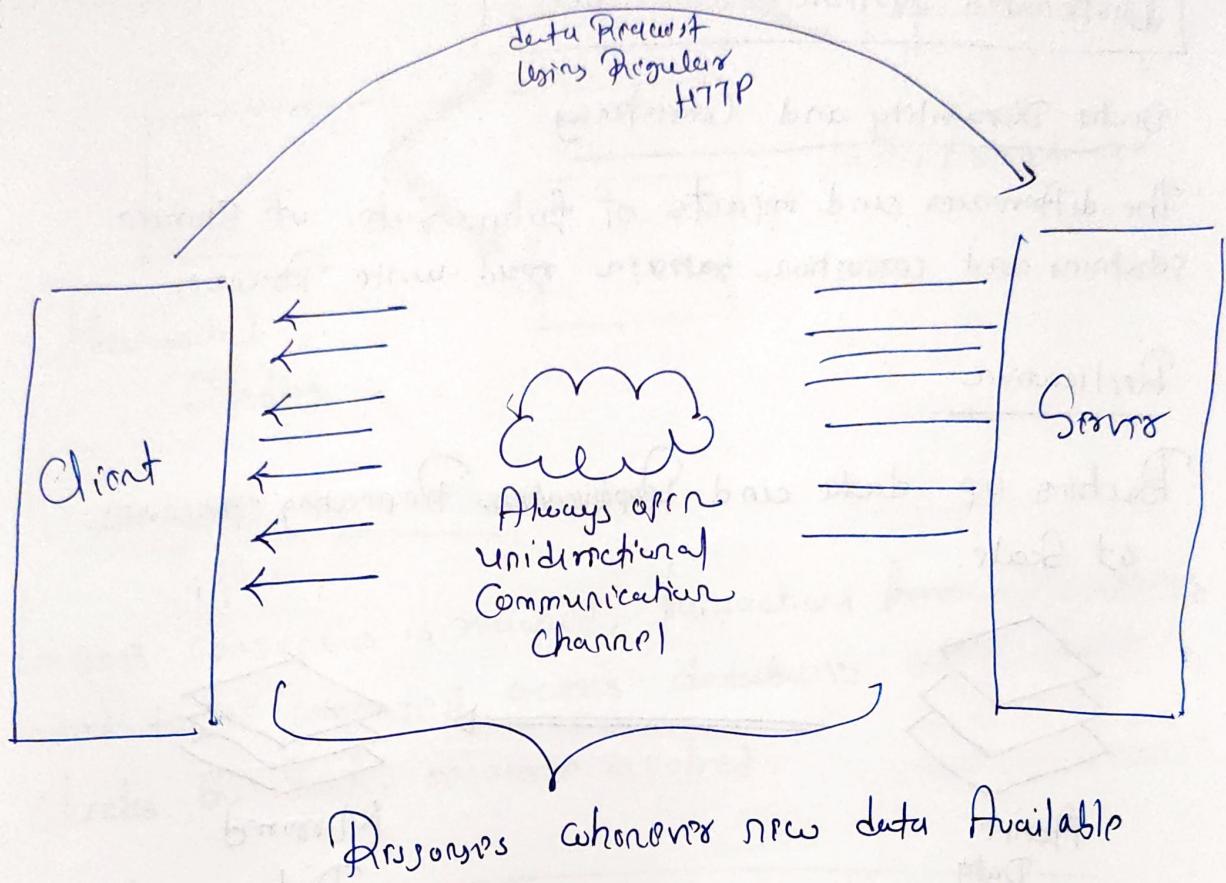
WEB Sockets

- Full duplex communication channel over single TCP connection.
- Provides Persistent Communication (Client & Server can send data at anytime)
- Bidirectional Communication is always open channel.



Server-Sent Events (SSE)

- Client establishes persistent & long term connection with Server.
- Server uses this connection to send data to Client.
- * If Client wants to send data to Server.
 - Requires another technology / Protocol.



- But when we need Real time data from server to client
- Server is generating data in a loop & will be sending multiple events to the Client.