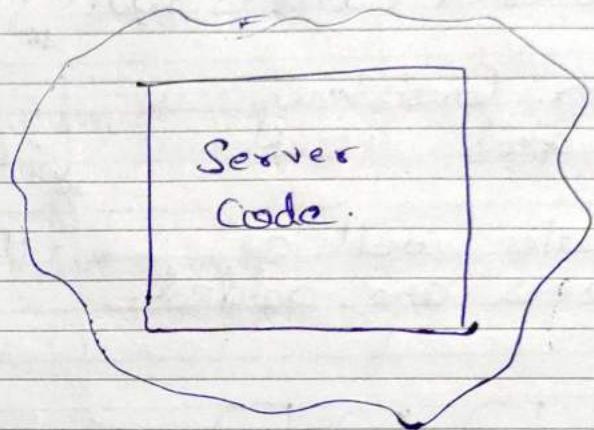
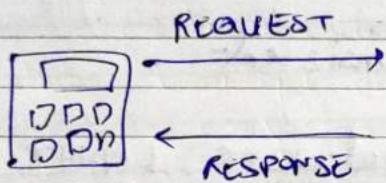


13/10

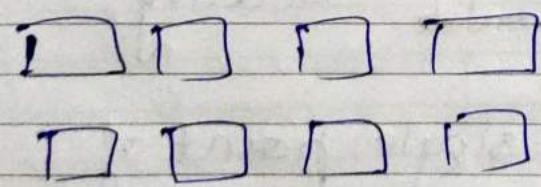
## System Design



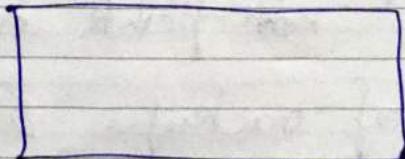
## Scalability

- ① Buy BIGGER Machine.  $\Rightarrow$  VERTICAL SCALING.
- ② Buy MORE Machine.  $\Rightarrow$  HORIZONTAL SCALING.

### HORIZONTAL



### VERTICAL



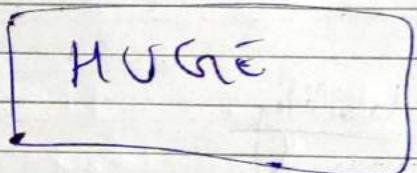
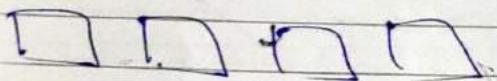
- Load Balancing Required.

- N/A.

★ → Good Qualities.

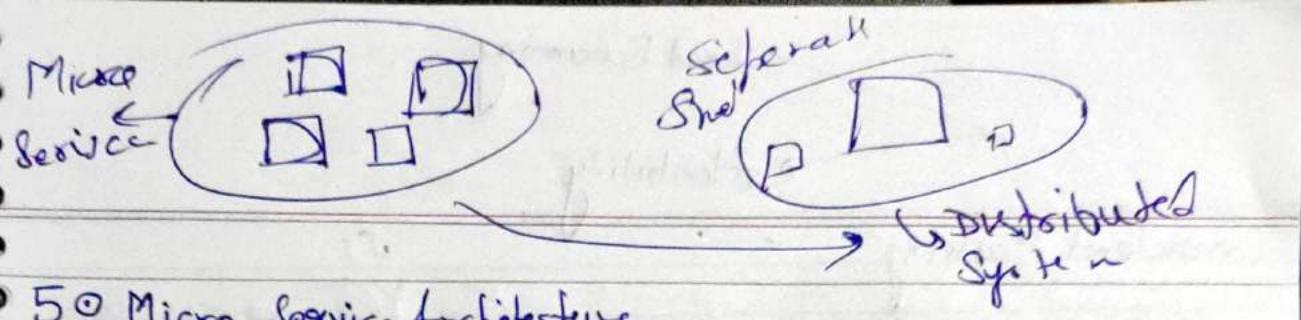
Resilient → preventing failures.

- Resilient
- ★ • Network calls (slow).
- Data consistency is a real issue.
- ★ • Scales well as users are added.
- Single pt. failure.
- ★ • Interprocess communication (fast)
- Consistent
- Hardware LIMIT.



### PIZZA SHOP EXAMPLE

- 1① Optimise processes and increase throughput using the same resource.
  - 1 chef.
  - Vertical Scaling.
- 2② Preparing before hand at non-peak hour. → CronJob
- 3③ Keep backups & avoid single point of failure. → Backups
- 4④ Hire more resources. → Horizontal Scaling.



## 50 Micro-Service Architecture

### 60 Distributed System (Partitions).

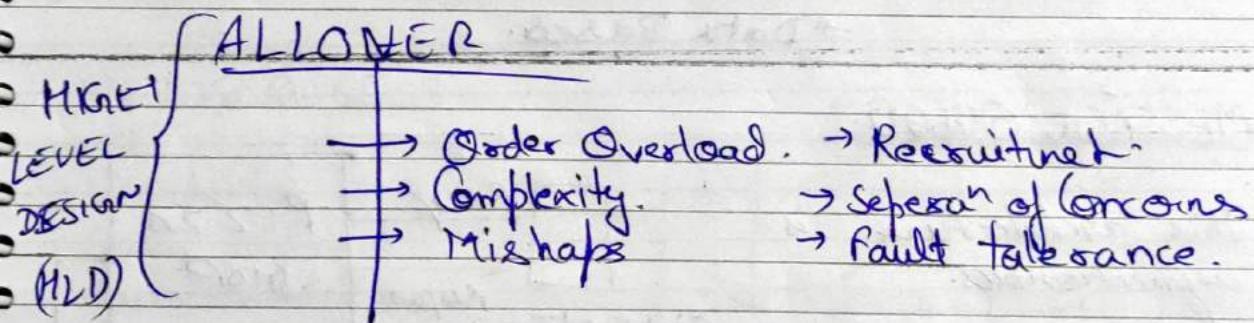
70 Load Balancer (sends req to diff servers). |

- Scalable
- Fault tolerance
- Extensible

80 Decoupling → Separation of responsibility. | Separates system more responsiblity.

90 Logging & Metrics

100 Extensible.



## Load Balancing.

Balancing load among servers.

- This uses Consistent hashing.

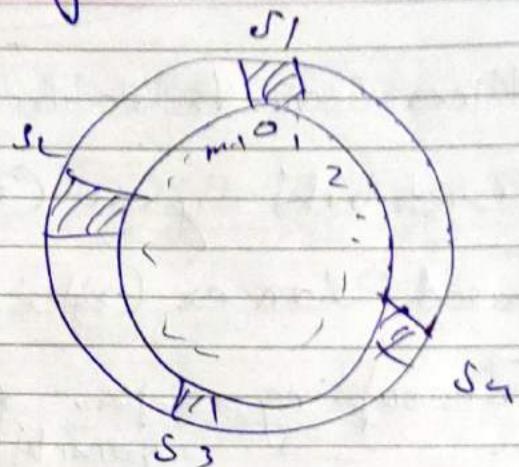
Using hash functions to map request to N-servers.

→ Load Balancing.

→ Flexibility.

## Consistent Hashing.

We use cyclic structure to reduce request change on when servers are added.



We can use k-hash function.

This avoids skewed load on each server.

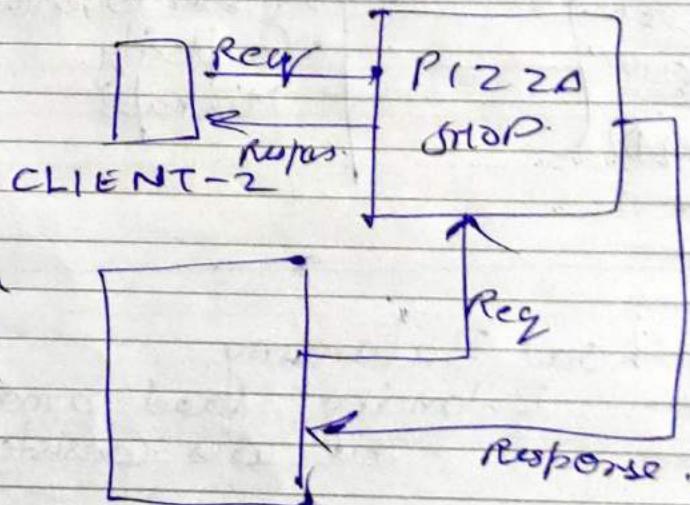
Extensively used in distributed systems.

→ Web caches.

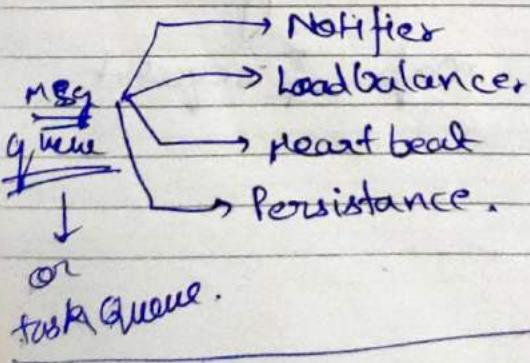
→ Data Bases.

## MESSAGE QUEUE

This architecture is asynchronous.



Principles of load balancer also ensures that there is no duplicates reg. on a server.

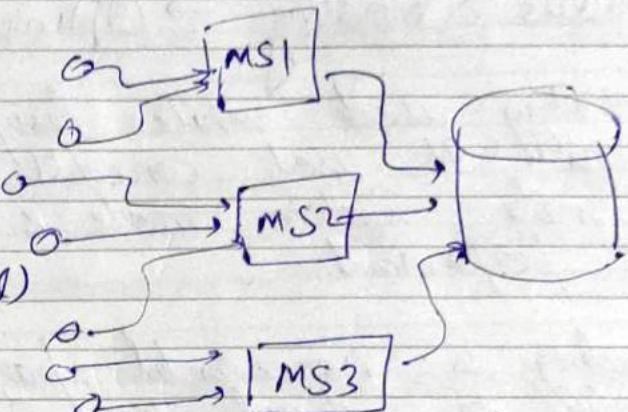


Eg: Rabbit MQ, JMS, ZeroMQ.

## MONOLITH

### Advantages

- Good for small teams.
- Less Complex.
- Less Duplication.
- Faster (due to procedural call...)



(e.g. Stackoverflow)

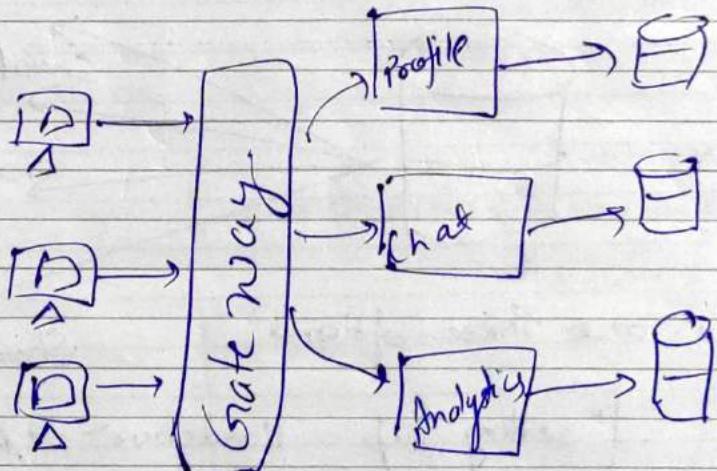
### Disadvantages

- Complicated Deployment.
- Not decoupled.
- Single pt. of failure.

## MICROSERVICES

### Advantages

- Scalability.
- Easier for new team members.
- Easier to reason about (which is used more).



### Disadvantages

- Harder to design.
- Needs skilled Architects.

Good for large scale systems.

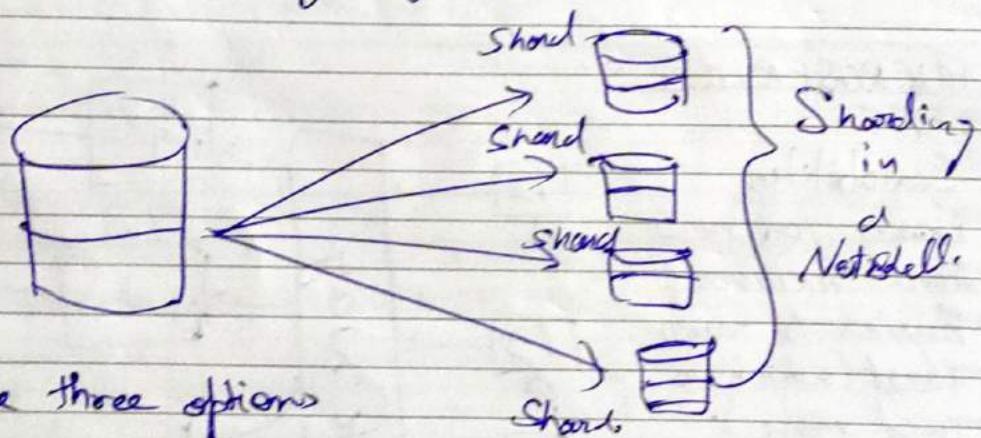
Eg: Google, facebook.

## Database Sharding → Optimisation Technique.

A strategy used while designing server side systems. It uses concepts like sharding to make system more scalable, reliable and performant.

Sharding is horizontal partitioning of data according to shard key. The shard key determines which database which database the entry to be persisted is sent to.

Some standard strategies for this are reverse proxies.



There are three options

- Setting up hardware → foot - Money
- Adding replica → foot → "Eventual Consistency"
  - ↓
  - MASTER  
SLAVE  
ARCHITECTURE
  - because  
of  
W  
R  
R
  - difficult to implement.
- Sharding

Shard → pieces of DB

User ID / Customer ID is a good candidate for shard ID.

We need a hash func" which is crafted in such a way that we can introduce the concepts of versioning.

In shard we have a table which map entries to databases.

What if this goes down?

Increased complexity as data is split b/w different databases.

Need to introduce intermediate layer

Responsible for mapping.

Poss

- Scalability.
- Availability & fault tolerance.

Com.

2 ways

Part in Mappn.

• Complexity ↗  
NON UNIFORMITY OF DATA.

we need to build.

Routing Layer.

Rersharding  
is req in this case.

Leads to "Scatter Gathers"

{ Occurs when one client has lot of data }

Codec → A way in which we compress video.

### Best practices:-

- Creating Index on Shards

↳ Meaning: People in New York > 50 years  
will lie in one Shard

### ④ HOW NETFLIX UPLOADS NEW VIDEO?

There are several things involved.

1. Different Format (F no. of formats)

→ high  
→ low  
→ med } → codec.

2. Different resolutions (R no. of resolu<sup>n</sup>)

→ 1080p  
→ 720p  
→ 480p

$$\text{No. of videos processed} = F * R.$$

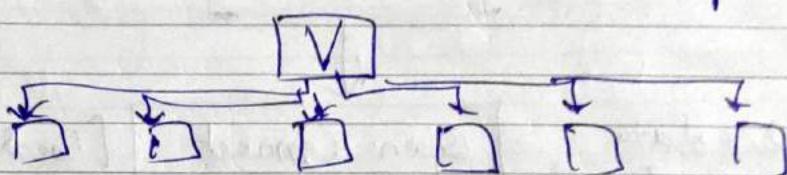
Automatically adjust video quality

↳ Adaptive Bitrate Streaming

Risk: What if System shuts down?

Solution: What Netflix does?.

→ Break a video into multiple chunks.



→ Then we process each chunk.  
(3min each).

Video file:   
0 3 6 9 12 15

What if it breaks on some action scene

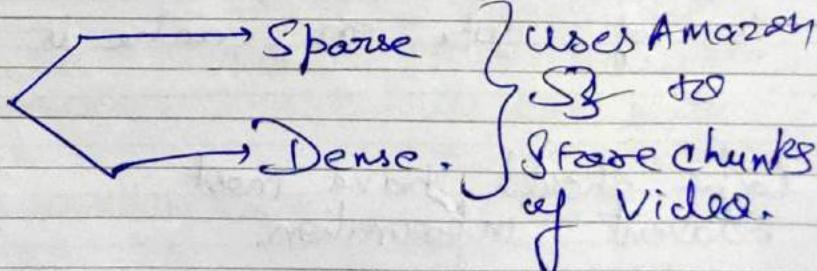
So what it does is ↴

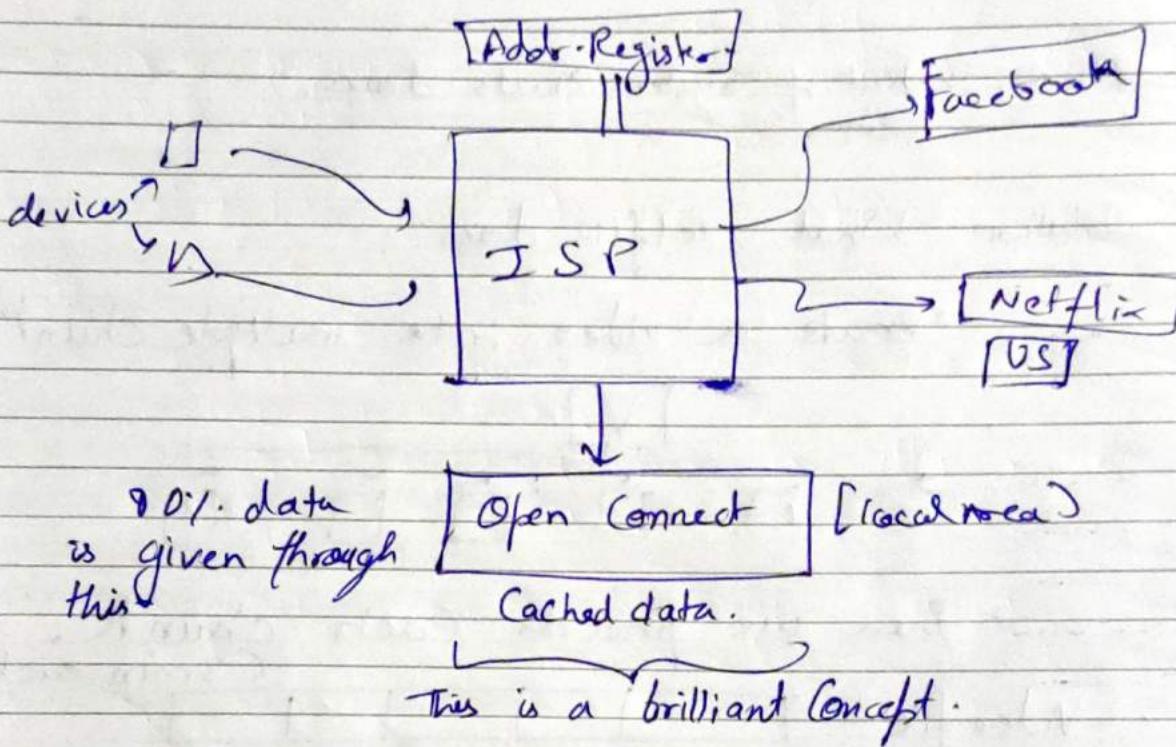
NOT to break on TIMESTAMP  
↓

INSTEAD break on scenes.

New file:   
0 12 14 20

Uses "recommendation" algorithms with 2 scenarios.





## ~~THUNDER~~ as a MICROSERVICE Architecture

### Distributed Caching

Why :-

- To save network calls.
  - Avoids recomputations.
  - reduces database load.
- } speeds up responses.

Cache stored in SSD, which is extensive.   
 ton of data on cache is counterproductive

Cache should have most relevant information.

↓  
increases search time.

## Cache

Policy → The way you decide to load and delete data from cache.

Eg: LRU - Least recently used.  
- Prog. LIFO.

LRU

Sliding Window based policy. → developed by Google.

Thrashing: When you save and delete data from cache without using it.

## Problems :-

- Extra calls.
- Thrashing.
- Consistency.

Where to place

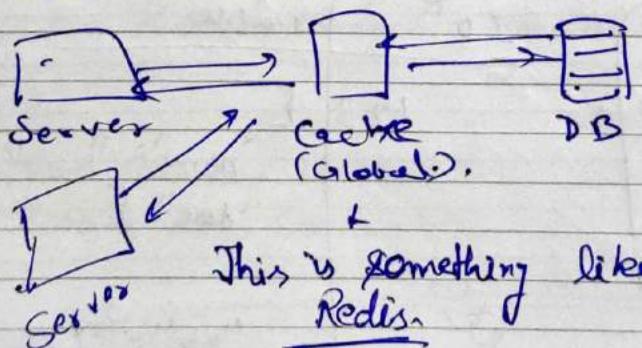
Near to database

Near to servers

## Near Server

→ If small cache → Keep In memory cache  
 ↳ If server fail, cache fails.  
 ↳ but its faster, and simple.

## Global cache



- Slightly slower.
- More Accurate.
- We can use redis cache.

- ① Write-through → while updating info, update on cache as well,
  - ② Write-back.      then DB
    - ↳ Update DB and then cache.

## Problem

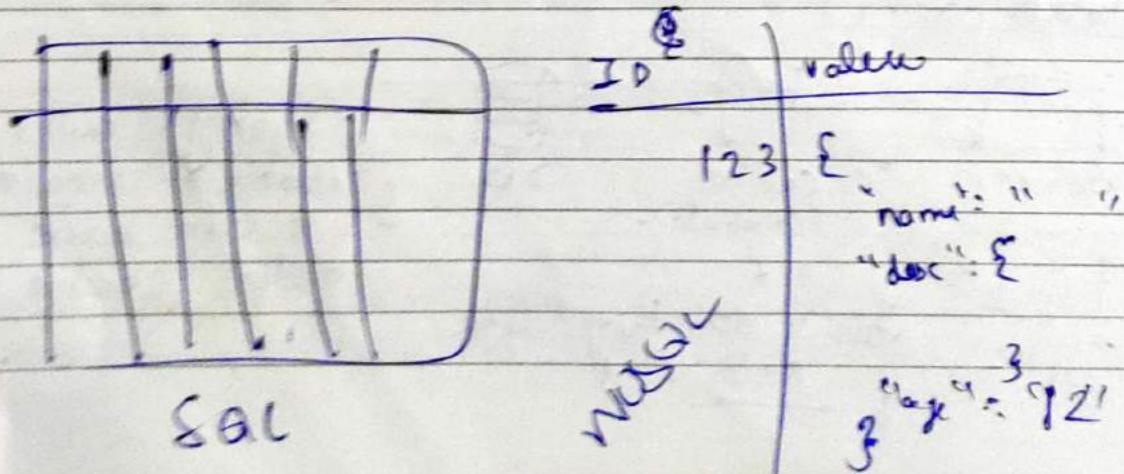
- What if there is other servers in with in-memory cache.  
~~With all~~ So updating only one is silly, we DONT use this.
  - In write-back, the only issue is speed.

How can we do better?

↳ Use both, depending on the info.

Critical data → Use Write-back  
Normal data → Write-through.

## SAL and NO SAL



(\*)

## Benefits of NoSQL.

- ① Insertions and retrievals are cheaper.
- ② Schema is easily changed.
- ③ Build for scale.
- ④ Good for Metrics (finding aggregate).

## Drawbacks of NoSQL

- ① Not build for Updates.
- ② Consistency Issue. (ACID not guaranteed)
  - ③ Main reason why financial system don't use these NoSQL databases.
- ④ Read time is slower.
- ⑤ Relations are not implicit.
- ⑥ Joins are HARD.

Eg: Same ID with diff data.  
SQL give ACID property

Consistency



<u>S A P</u>	<u>A C I D</u>
Consistency	Atomicity
Availability	Consistency
Partition tolerance	Isolation
	Durability

- Two Generali Problem.
- Two phase Commit.

## Database Design Tips

choice of DB → It doesn't impact functional requirement

### Factors :

- Type of Data ↗ Structured ↘ Non - " .
- Query Pattern
- Amt. of Scale.

# Always use some caching solution.

↓  
How it works → we have a key & value.  
Eg: Redis  
Memcached.  
etc.

# File System (for static content, eg, image, video, etc)

↓  
for this case → Eg: Amazon S3 +  
we use BLOB Storage      CDN. (along with S3)

Primarily stored in S3 ↗  
and cached in Data centers. ↗  
Distribu" content is in various locations. Worldwide

# Text Searching Capability. (Product like Netflix)  
Search movie.

↓  
Text Search Engine.

Eg: ElasticSearch or Solr



Also support FUZZY search  
↓

Both are build on top of Apache Lucene.

Provide searching capability

• AIRPORT } Identity  
    .            | this  
    .            | AIRPORT

The Above are NOT databases, these are search engines.

Database gives guarantee that data will not be lost.

# For Metric Data.



Time-Series Database. → Extent of RDBMS.



Sequential update in append-only mode

We query for last few min/hour etc.



Eg: Influx DB  
open TSDB

## Soln Pattern

- Sharding Data
- Replication type

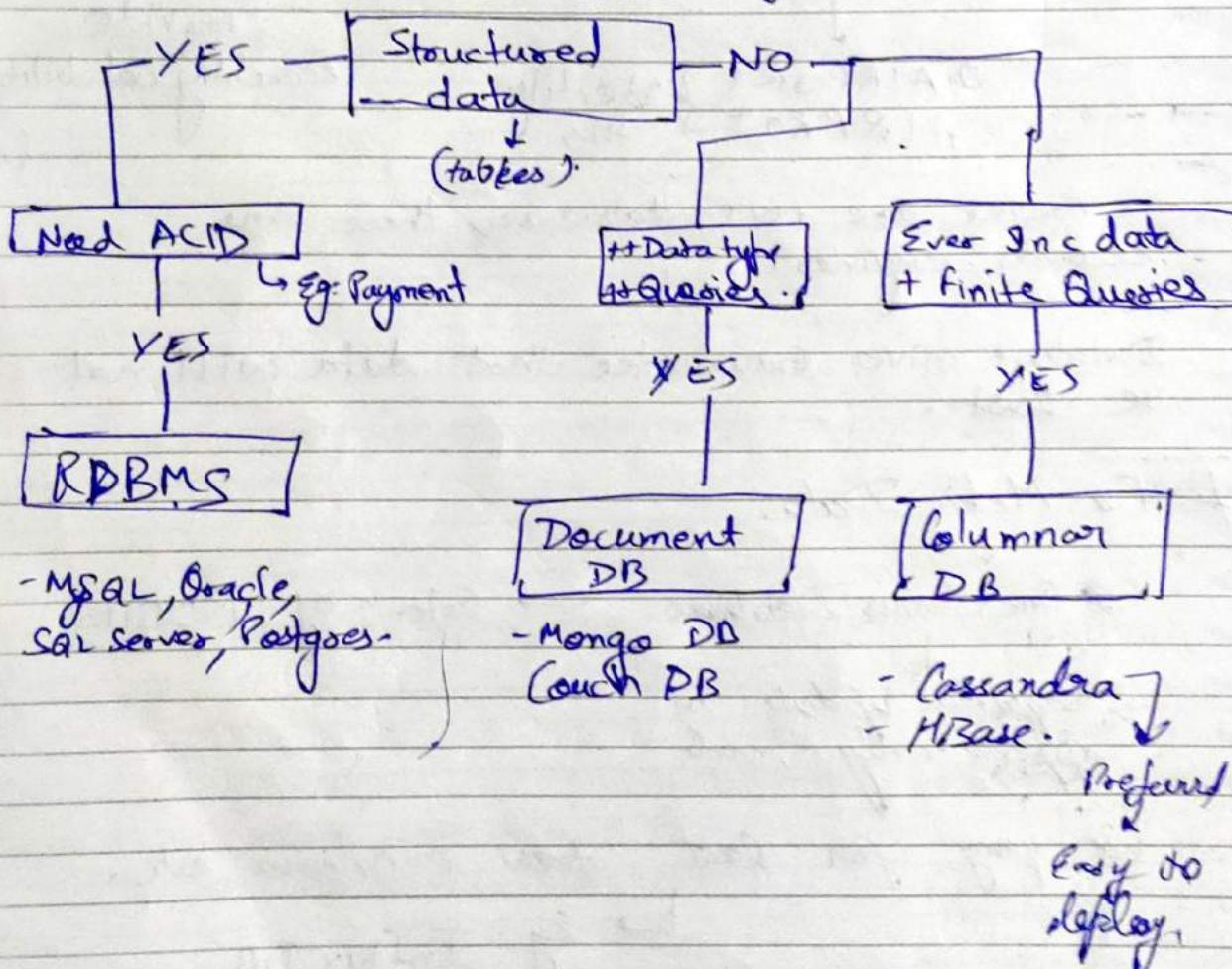
- Write ahead Logging.
- Separating data & Metadata storage.
- Load distributions.

# Large DB for Analytics (Amazon, Uber, etc)

We need DATA WAREHOUSE

→ Used for offline reporting.  
 → eg : MapReduce  
     ↳ for batch (Slow) processing.

In RDBMS → last option is Postgres.



Low Cost → Go for openSource.  
High Cost → Production Service.

Helps us choose tech  
in requirements.

To ensure  
no lag, low  
latency & high  
availability.

### Questions to ask.

- Location.
- Type of Data.
- Type of Queries.
- Scalability.
- Users / Customers.
- Scale (read & write)
- Performance.
- Cost

### Extra things

- Shards
- Monitoring
- Backup.
- Log.
- Caching (Netflix; OpenConnect)
- Sharding & Master-Slave.
- Load balancer / message que.
- Images (FileSystem) + CDN.

CPU Utiliz.  
Disk Utiliz.  
Throughput  
Log in Roll.  
add more nodes

### Non-functional (Write down).

- Scalable (1,000 req/sec)
- High Performant ()
- High Available
- Recall CAP (talk about())
- Cost (of everything)

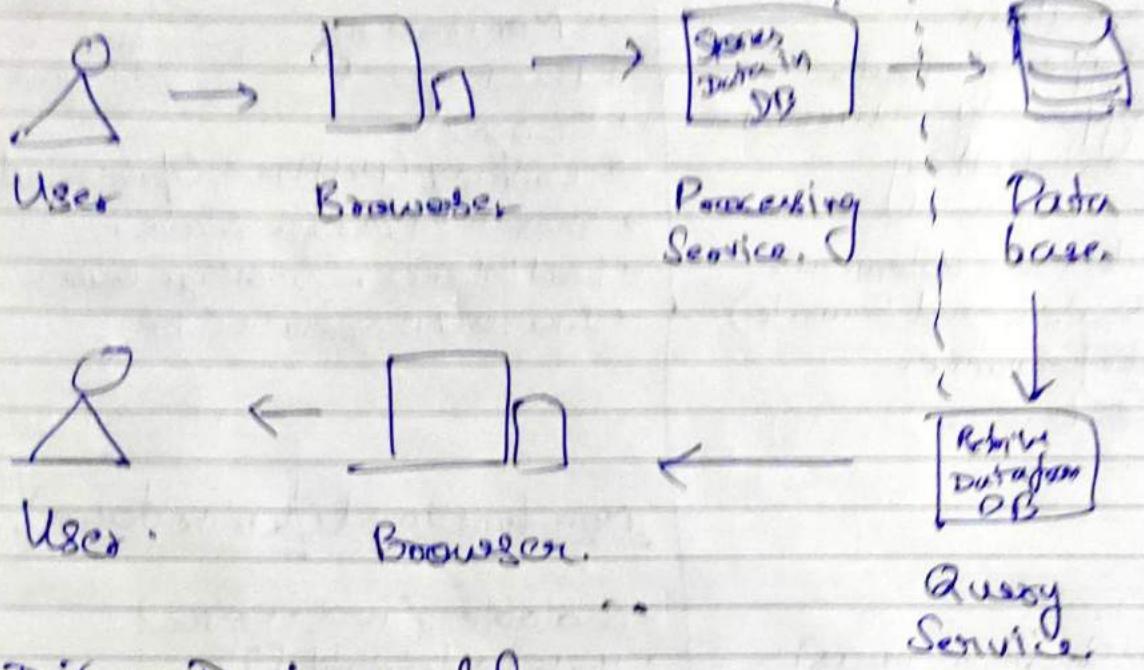
### High level

Next Page

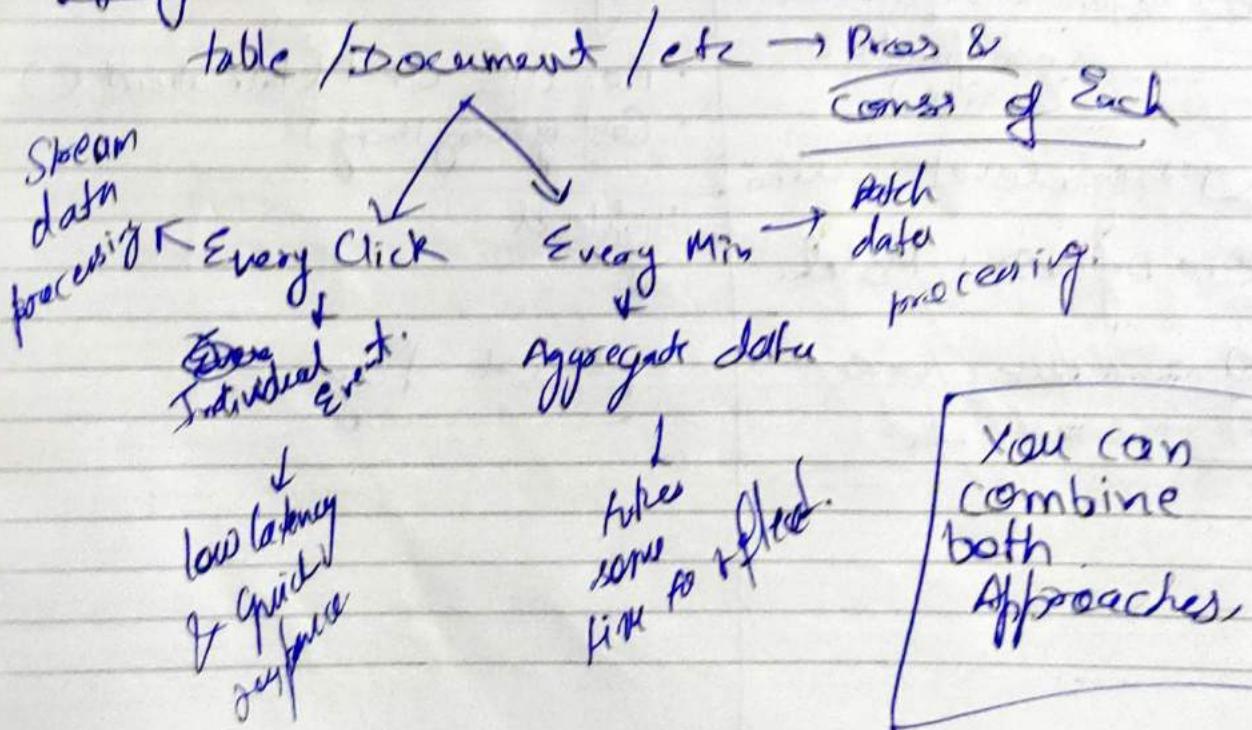
When you do off, they do get best of both worlds

Point 11

High level



## Defining Data model

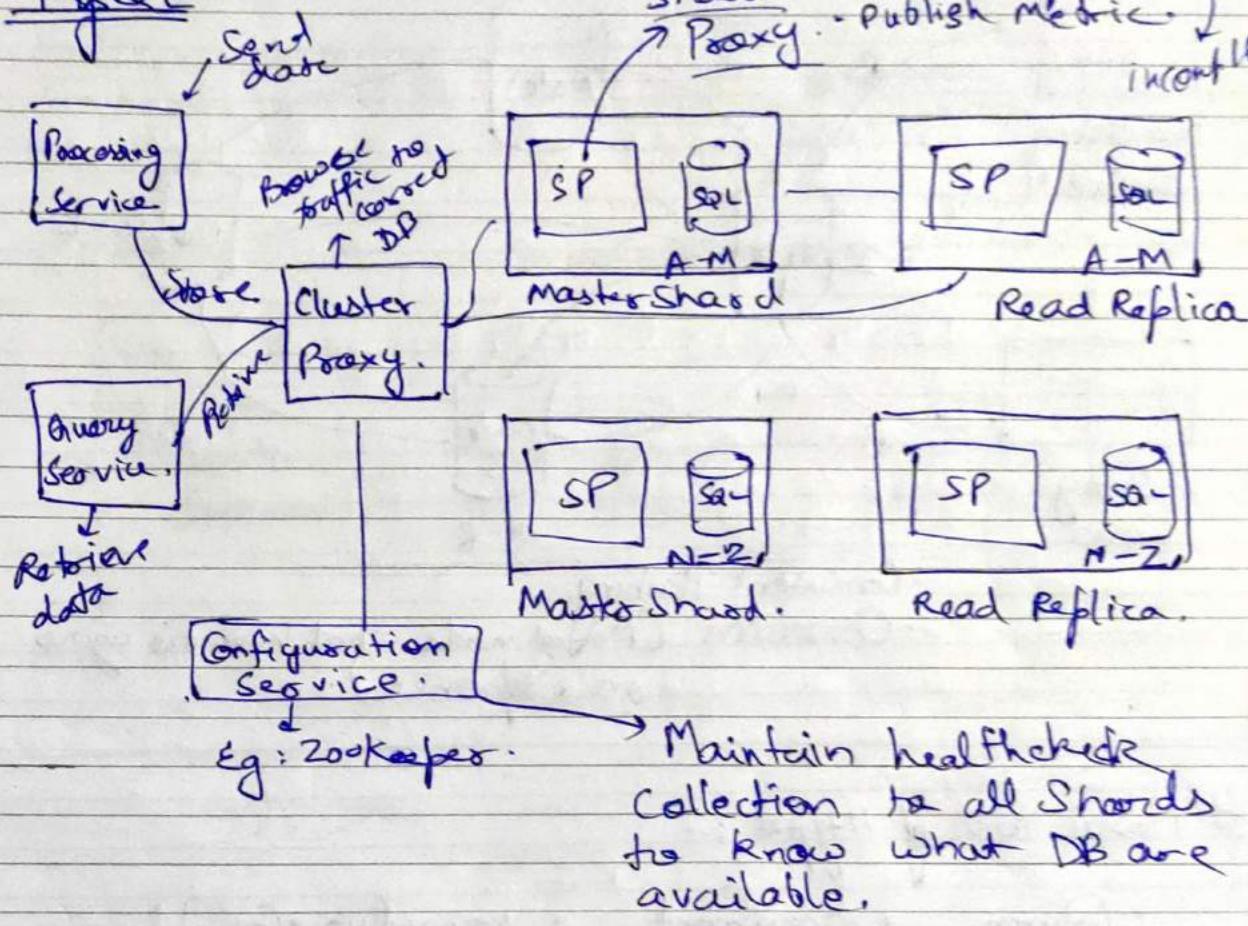


- Scalability
- Availability
- Performance

{ } //

## Database

~~MySQL~~

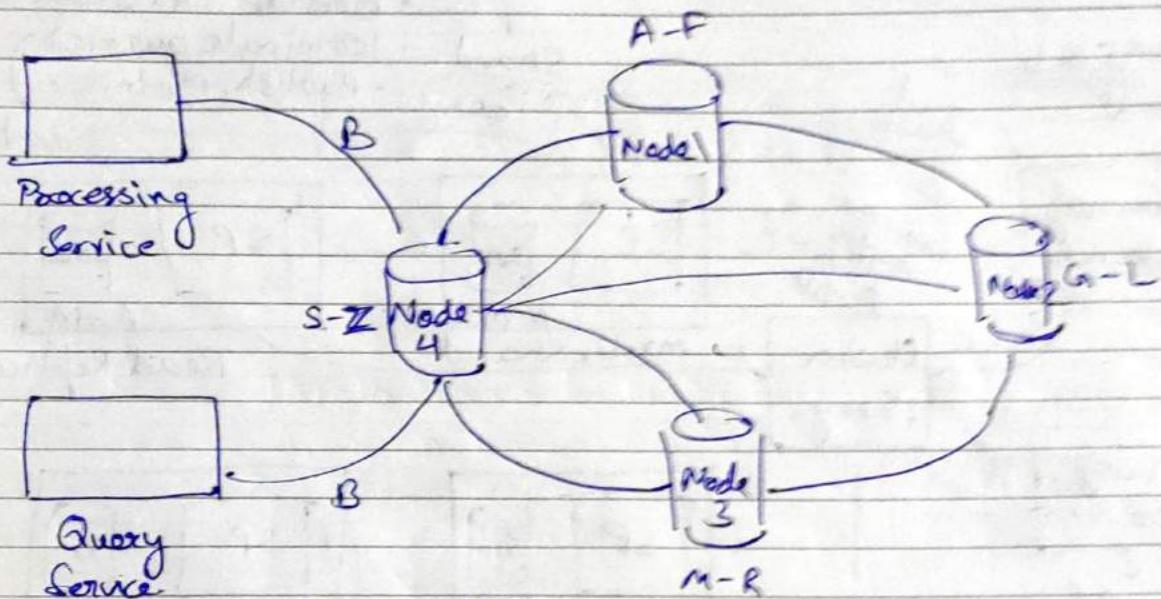


Above Design addresses all things

Scalability  
Availability  
Performance

Problem of Master-Slave is known for "eventual consistency".

## NoSQL (Apache Cassandra DB)



- Consistent Hashing
- Quorum (Majority of nodes that have to agree on response).

NoSQL are of 4 types :-

- Column
- Document
- key-value
- Graph.

BUILD SOLUTION IN

FIRST 20 MINS.

(Latency of throughput) → Numbers

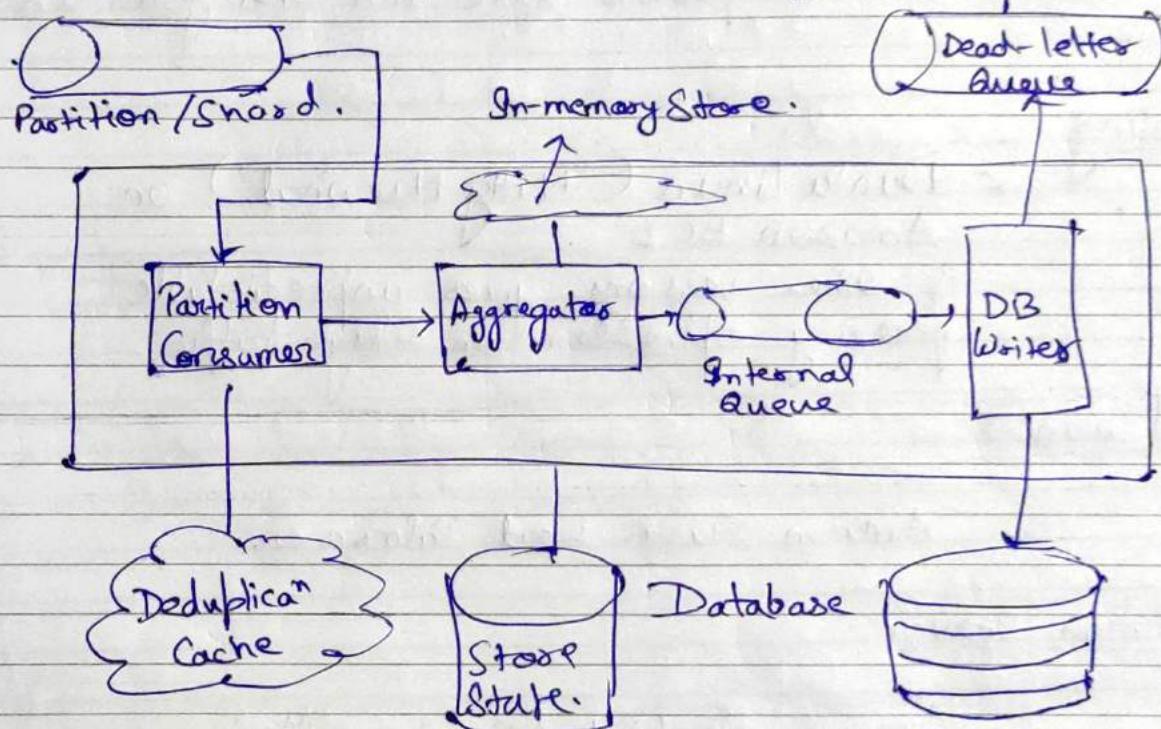
• Learn 'adv/disadv. of all database designs'

## # Processing Service

Scalable = Partitioning  
Reliable = Replica & checkpointing  
Fast = In-memory.

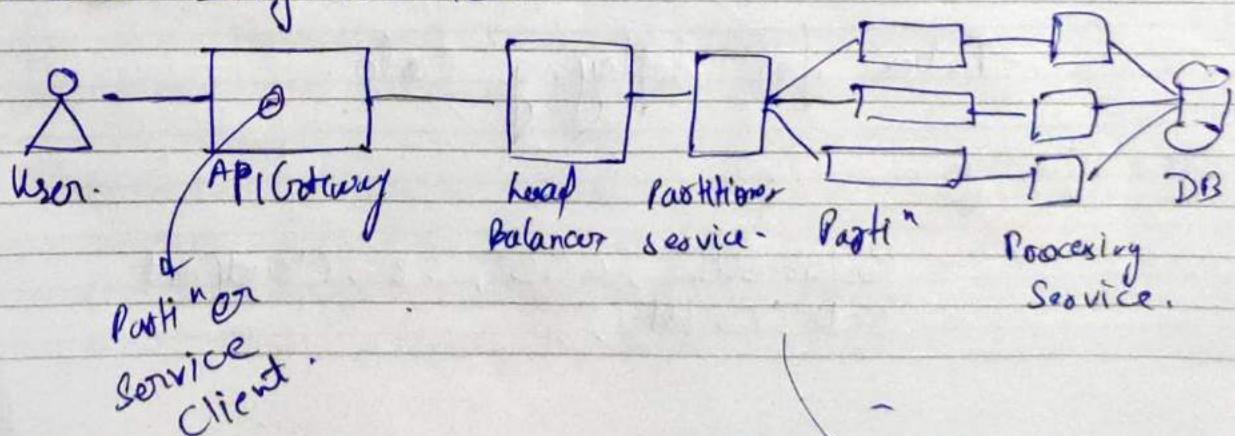
- Checkpointing
- Partitioning

to protect from availability issues.



## # Data Ingestion

### # Data Ingestion Path



## Instagram

### Principles

- Keep it very simple.
- Don't reinvent the wheel.
- Go with proven solid tech when you can

## Hosting

- Ubuntu Linux ("Natty Narwhal") on Amazon EC2

P.S.: Other versions gives unpredictable freezing episodes on high traffic.

## Load Balancer

- Amazon Elastic Load Balancer.

## Application Server

- Amazon high-CPUs Extra Large Machines
- WSGI → Gunicorn (easy to configure) which is less CPU intensive
- Fabric for deploying code.

## Data Storage

- PostgreSQL, EBS, S3, CloudFront
- Redshift (Analytics)

90% customers  
use IoT Workflow

## Caching

- Redis, Memcached
- Amazon, Elastic cache.

## Task / Push Notification Queue

- Gearman for Task Queue
- pyapns for push notification.

## Monitoring

- For Graph Metric → Munin.
- For External Monitor → Pingdom.
- For hard "notifi" → PagerDuty
- For Py Error → Sentry

④ Run a script whenever we send a new msg.

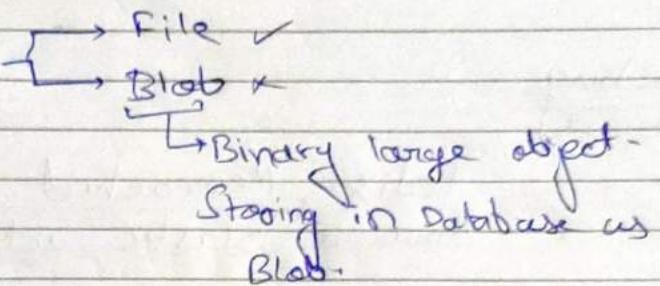
```
CREATE TRIGGER name
WHEN EVENT
EXECUTE
ON table-name TRIGGER_TYPE
EXECUTE stored-procedure
```

} Skeleton

```
CREATE TRIGGER sendMsg
WHEN INSERT
ON Message FOR EACH ROW
EXECUTE stored-procedure.
```

System design interview is not about building a working system, it's about measuring my aptitude

## ① For Storing Images



What DB provides :-

- 1) Mutability
- 2) Transaction → ACID } → NOT Required
- 3) Indexes (Search)      } for Images.
- 4) Access Control.

Adv of File System:- → BEST

- 1) Cheaper
- 2) Faster
- 3) CDN.

In this we need to store URL in DB.

### Protocols to Use

- Internal → GRCP
- External only → HTTP
- Video ( ) → RTMP → live streaming,

Kassandra  
Elastic Search

KNOW  
ABOUT  
ITF

Learn #  
about  
trade off.

FOCUS  
ON  
CORE

- Ensure high availability
- Improves cost efficiencies
- Helps define & meet SLAs
- Provide insights

$$\begin{aligned}
 2 \text{ hr} &\rightarrow 4 \text{ GB} \\
 &= 400 \text{ MB} \\
 200 \text{ MB} &= 1 \text{ MB}
 \end{aligned}$$

## Capacity planning: YouTube

$$1 \text{ B} \times \frac{1}{1500} \times 10 \text{ mins} = 10^7 \text{ min}$$

$$= 10^7 \times 3 \text{ MB}$$

$$200 \text{ MB} \rightarrow 1 \text{ hr.}$$

$$= 30 \text{ TB}$$

$$\frac{200}{60} \text{ MB} \rightarrow 1 \text{ min}$$

Art. of data  
to keep raw

We need  
in multiple  
resolution.

$$3 \text{ MB/min}$$

footage of video  
on storage disk.

We need copy

for fault tolerance & redundancy

For 3 copies of data.

$$\rightarrow 90 \text{ TB} \rightarrow 180 \text{ TB} \rightarrow 0.2 \text{ PB}$$

Alternate approach → We can store video in form of image frames.

- Estimate daily data usage
- Estimate what the cache & how many system of 16 GB's are required.
- Estimate no of processor to process video.