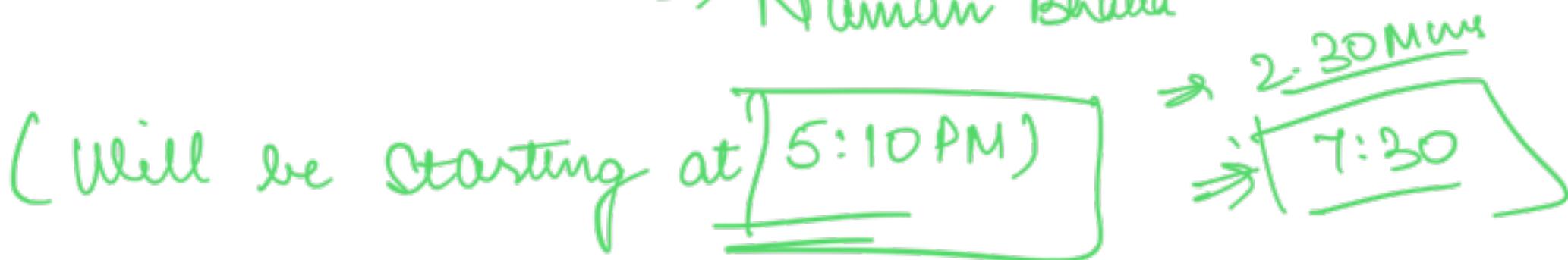


HLD of Live Streaming Services

HLD

→ Naman Bhalla



Break will be around
6:30 PM

6:30 PM

[Scaler.com] events

Naman Bhalla

YouTube live } Hotstar

1cr

How do OTT Platforms

[SOL]

Scale for such a
large userbase

HLD Master class

(High level Design) System Design

→ No code

→ Not a Frontend architecture

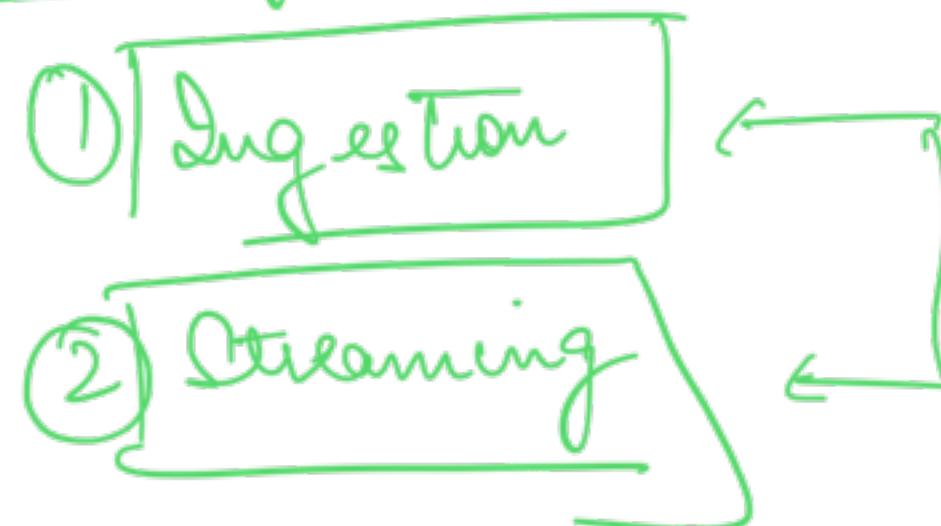
(Backend & UI)

Flowchart Theory -

Agenda

- What is HLD
- How to approach HLD

- Live Streaming Platform



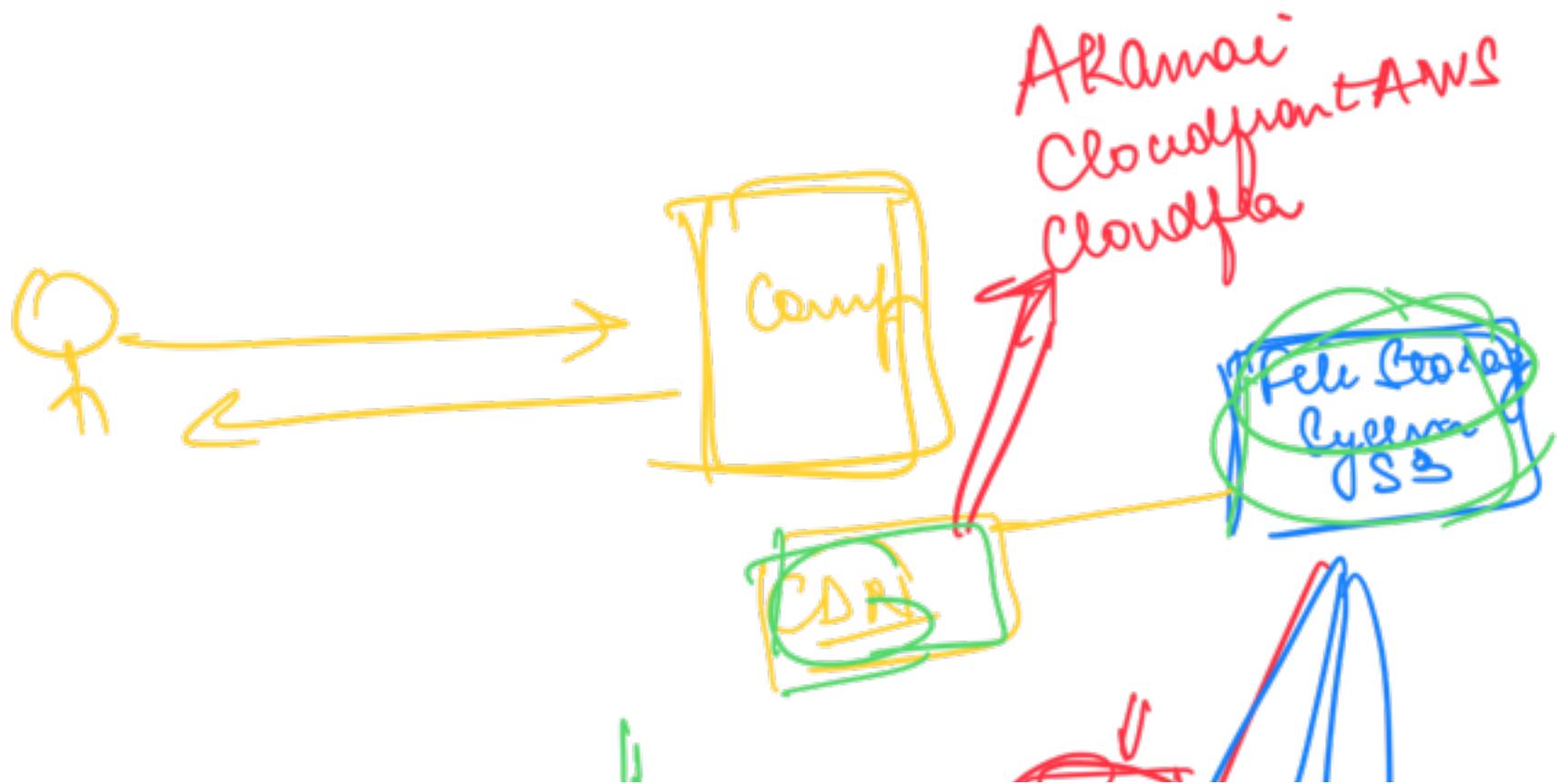
What is HLD

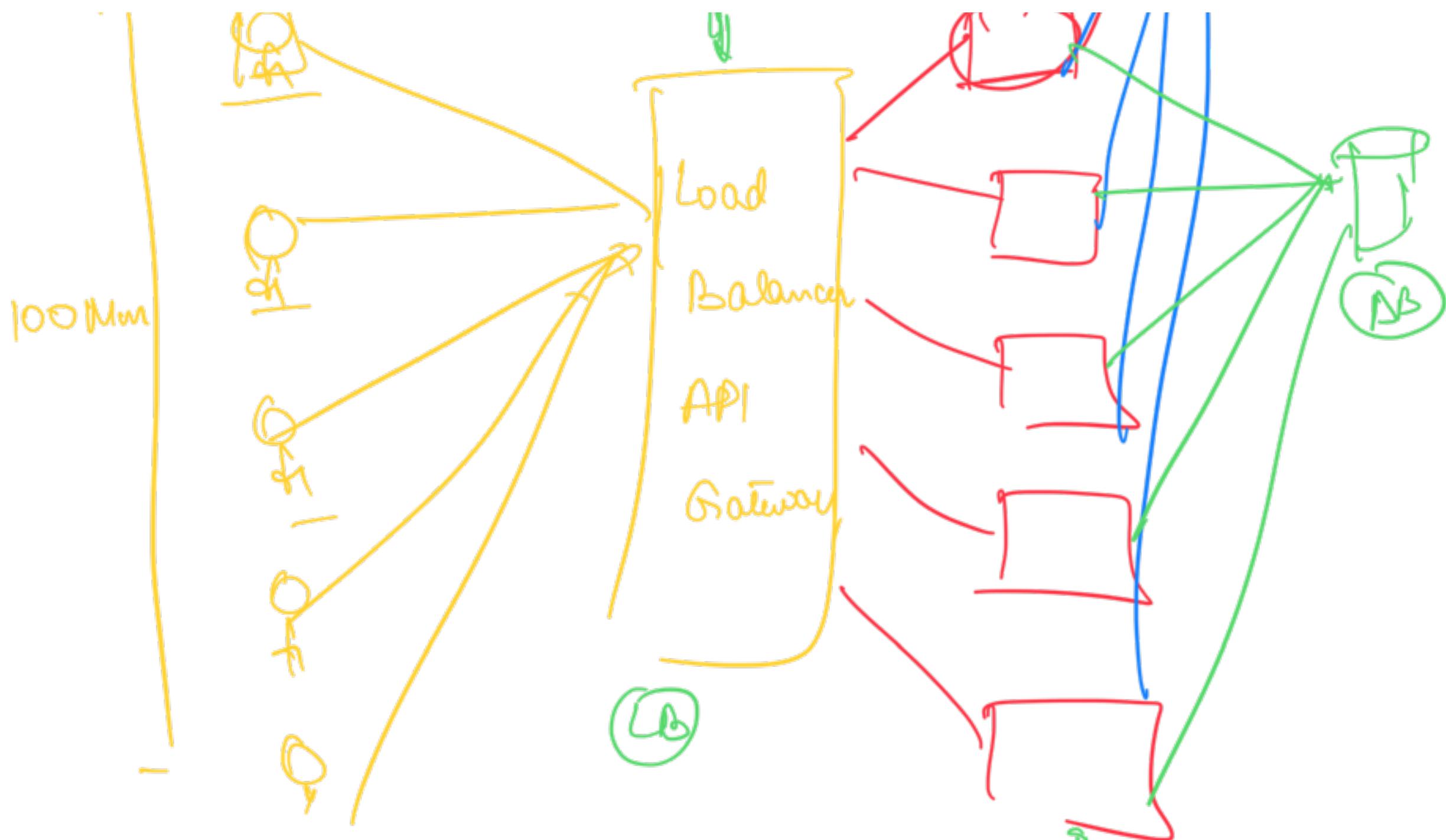
→ Bird's Eye View
Overview

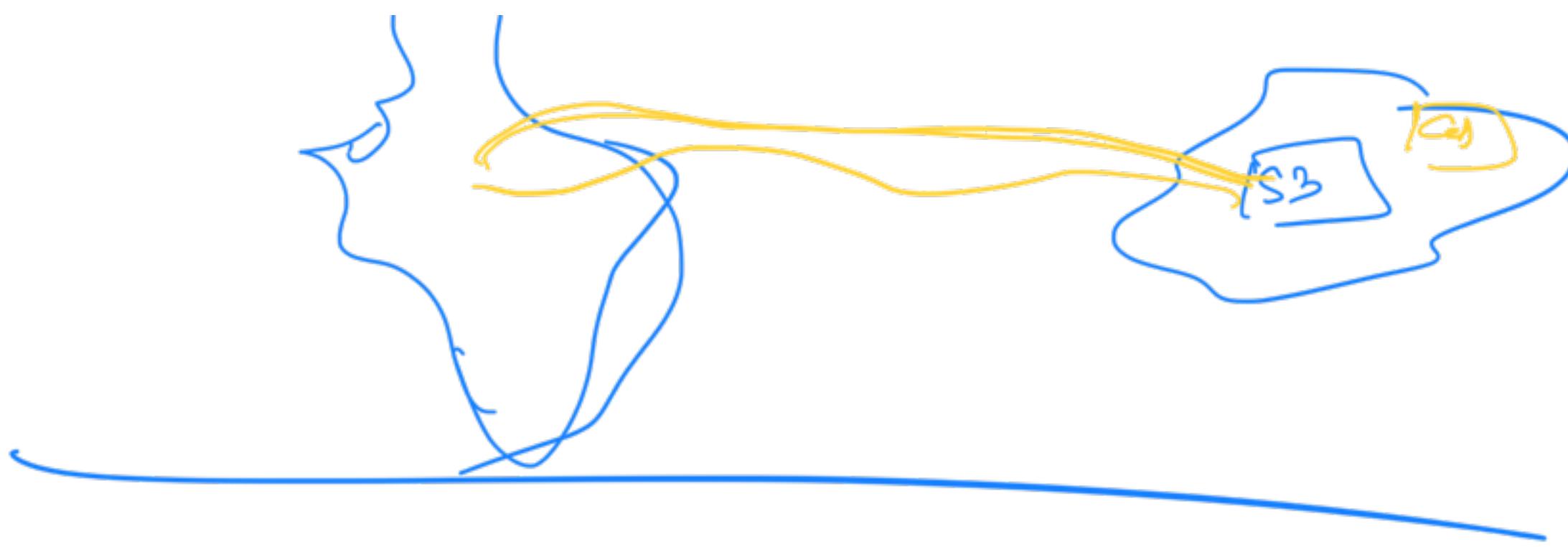
HLD: High Level Design

System Design

Bird's Eye View | Overview of a Software System







Infrastructure layers

HIS | System Design: deals with the design and interaction of different infrastructural layers that work together to serve you ... deal with expected

the app[™] at a very ---
efficiency

<10Sec

How to approach HLD

① Understand the Problem Statement

a) Functional Requirements \approx
(diff features of the system)

b) Non Functional Requirements \approx
(Qualities (Expectations))

Netflix

- ① Auth
- ② Streaming
- ③ Recommender

WhatsApp

- ① Call
- ② Multimedia Mess
- ③ Groups

CAP Theorem

Latency (delay) → Consistency & Availability

② Back of Envelope Calculations

→ Real SW Job X

→ Based on the requirements,
quantify the scale
→ guessimates

- a) Size of data → Total
- b) N/W Bandwidth → Total
- c) Concurrent Users → ≈
- d) Usage of a feature →

Worst Case

1 MB data

If you want to build an application

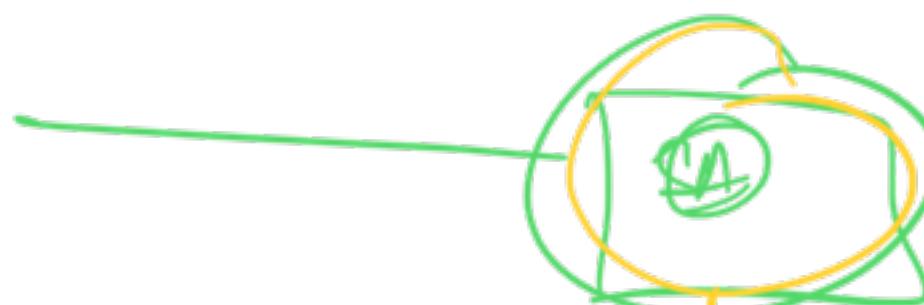
that supports Partition Tolerance

you can either have Consistency

or Availability. Not Both.

20 GB data

Name

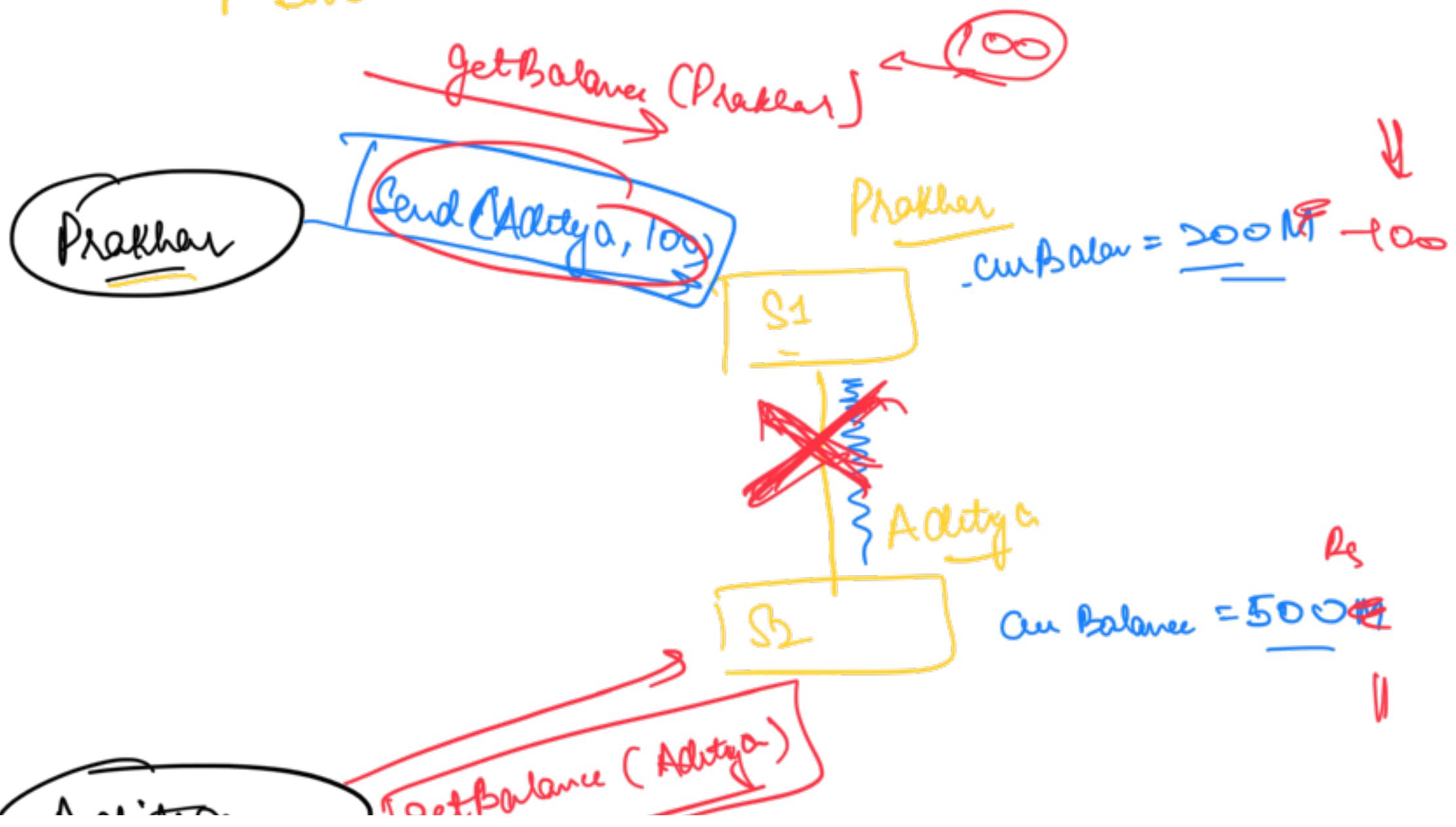


⇒ Consistency
⇒ Availability





P lends 100 Rs to Aditya^a



Audrey 170° 500

PT ✓

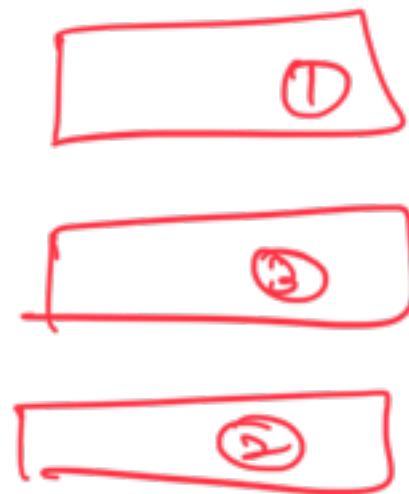
Consistent ✓
Available ✗

CAP Theorem

(A)P → Availability → Hotstar, WhatsApp

OR

(C)P → Consistency → Bank



WhatsApp

Max # of users at same time

$2BN$ people who use

$1BN$ people will have day
+ some time

For \Rightarrow 10x open at same time
 \Rightarrow 100 Mn concurrent user
 \Rightarrow 500 Mn

2 Bn people use daily
5% of people
20 messages
1% of people \Rightarrow 500 mes

③ APIs to support
/send /~~receive_id~~ , message)

/watch (video-id)

④ Schema

~~What are going to be diff entities that I will store~~ and attributes for those

entities

Data Models

||

users

col	name	phoneNo	pin

⑤ Design

a.) Start with a basic design

b.) Iterate (finding a bottleneck and resolving it)



Functional Req

→ Streaming at different qualities

(
360p
480p,
1080p)

→ Upload video

→ Replay

Non functional Req

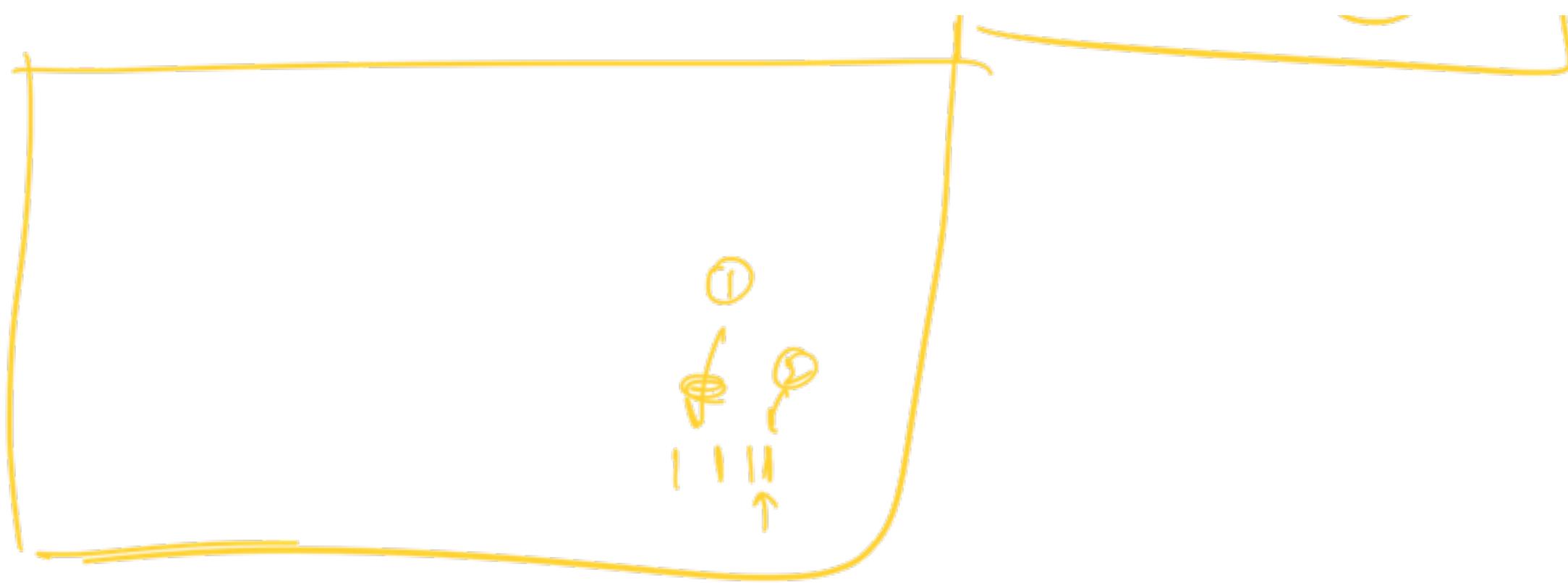
① Streaming

- Very low latency
- High Availability

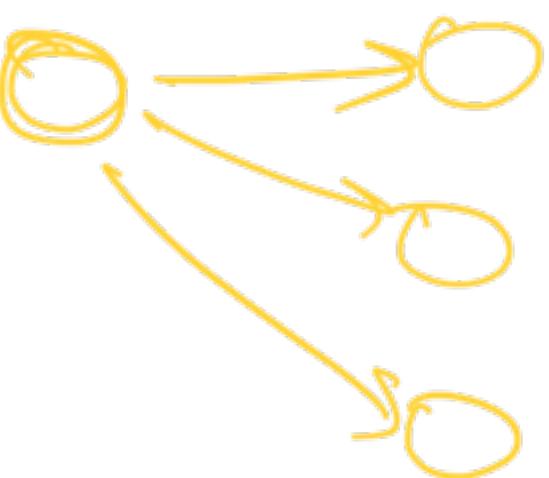
② Upload Video

- High Availability
- Very low latency

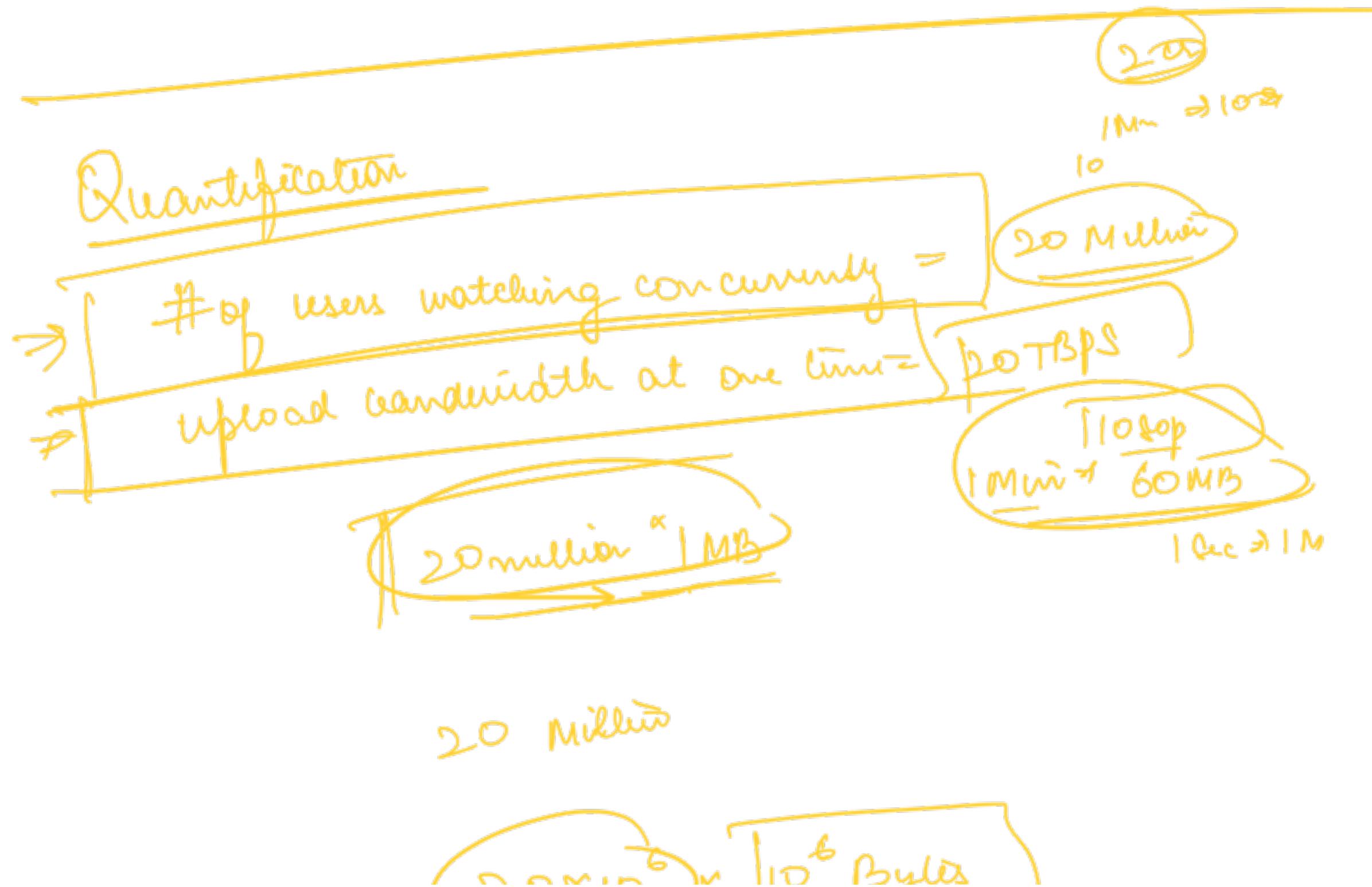




Eventual Consistency



①
6



20 TB

20×10^{12} bytes

20×10^3 GB

20 TB

1 Ma claim can have 1 GB/s





| Machine \Rightarrow | Gbps

| Gbps \Rightarrow | Ma

$$20 \times 10^3 \text{ Gb} \Rightarrow 20 \times 10^3 \text{ Machn}$$



every video \Rightarrow 4GB \rightarrow 4100h

Avg size of a video \approx 30MB

Avg size / min \approx 60MB

Storage

$$400 \times 10^3 \times 30 \times 60$$

MB

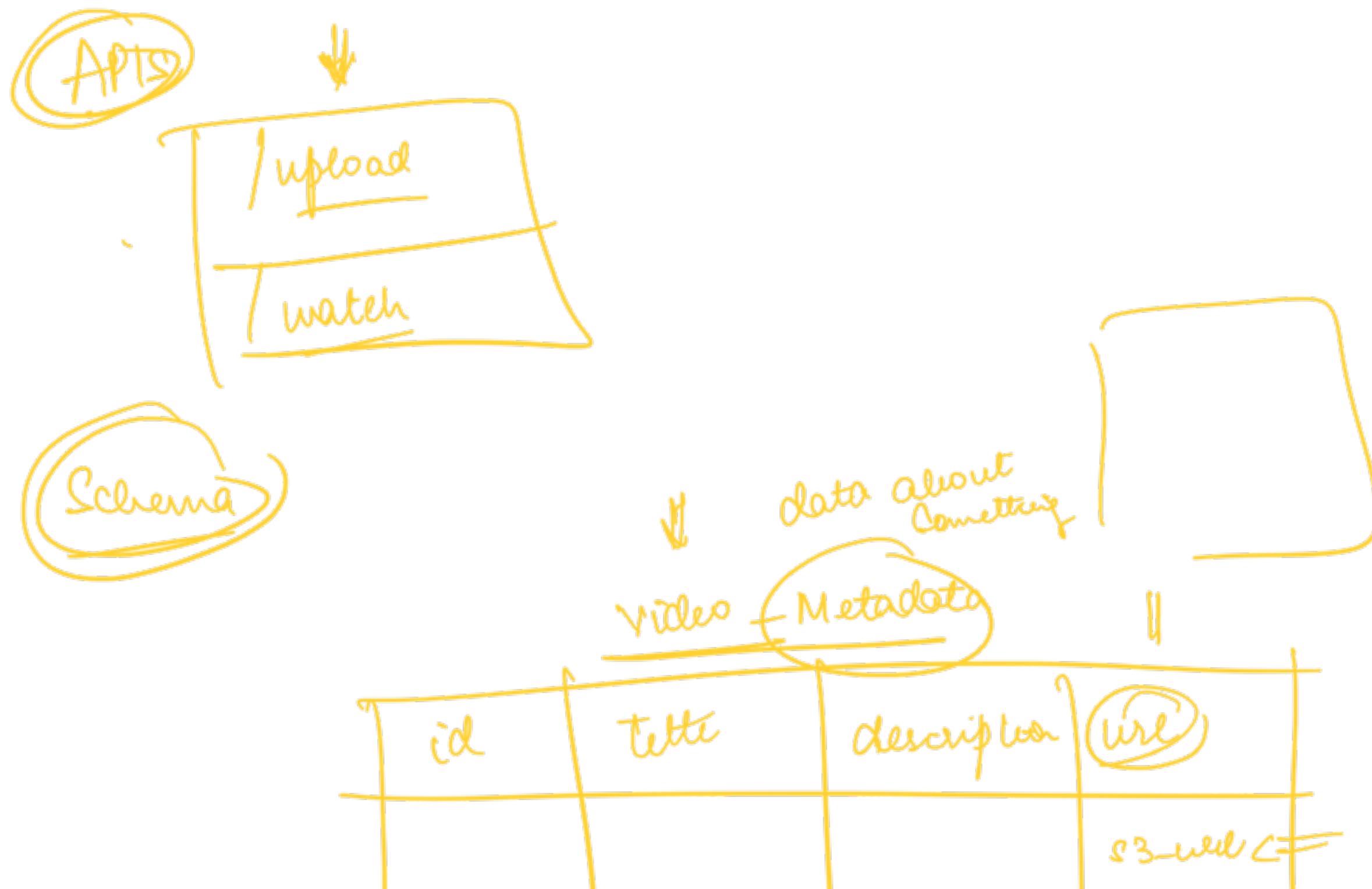
$$400 \times 30 \times 60 \text{ GB}$$

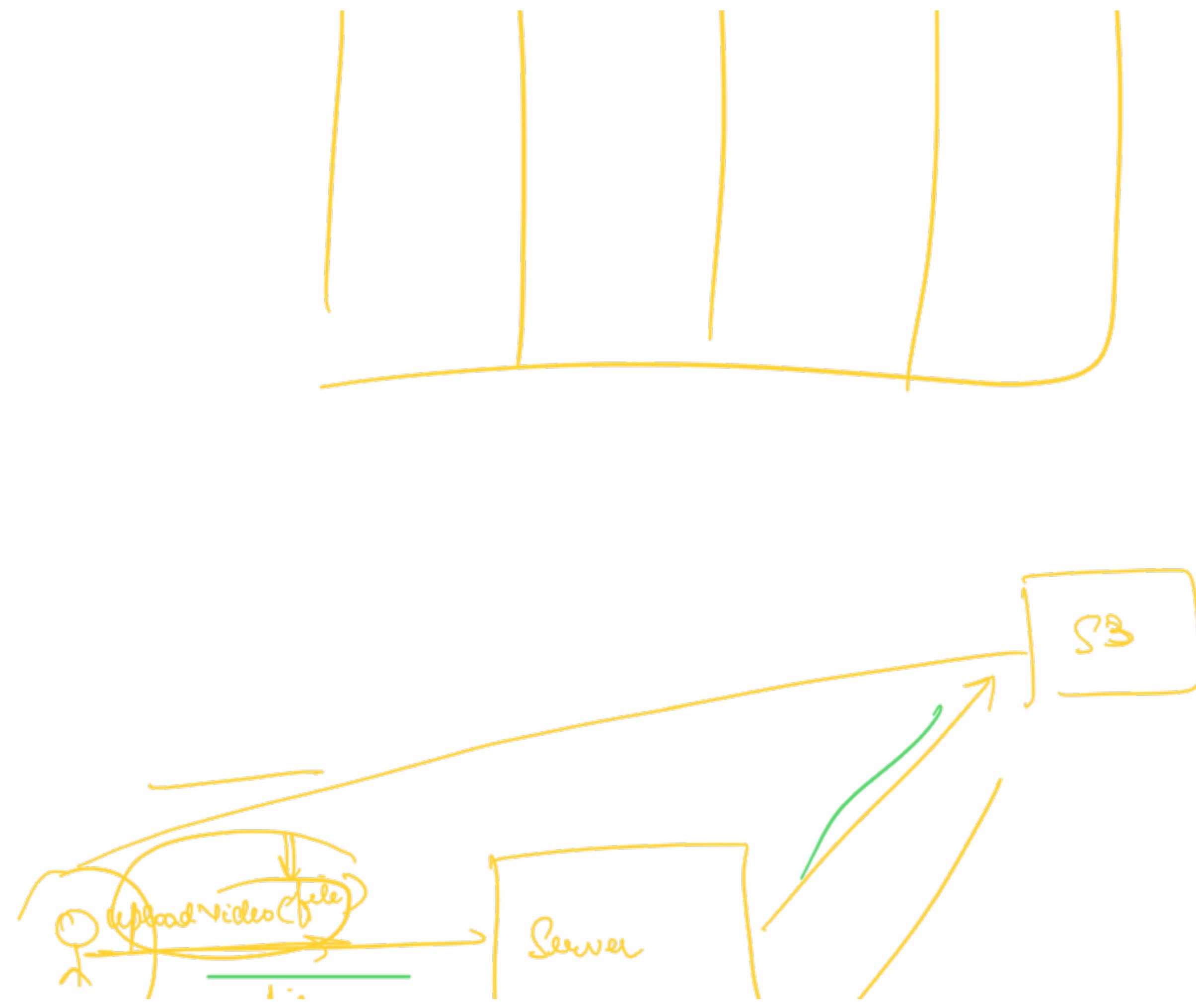
40KB x 6TB

720 TB

~ 1PB

AWS
S3

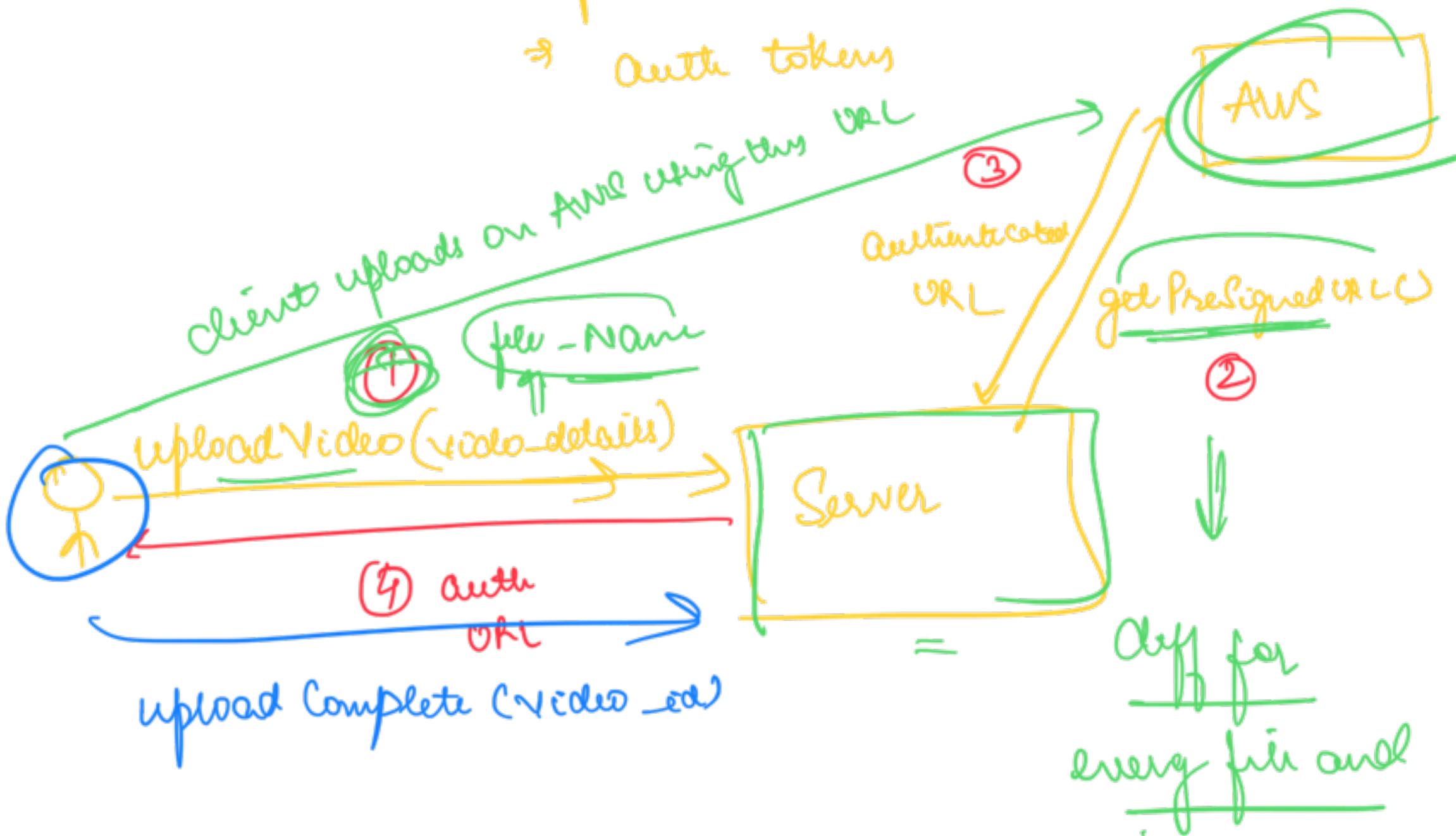




file

l

Pre Signed URL \Rightarrow feature of AWS
 \Rightarrow provides authentication
 \Rightarrow auth tokens



it is valid

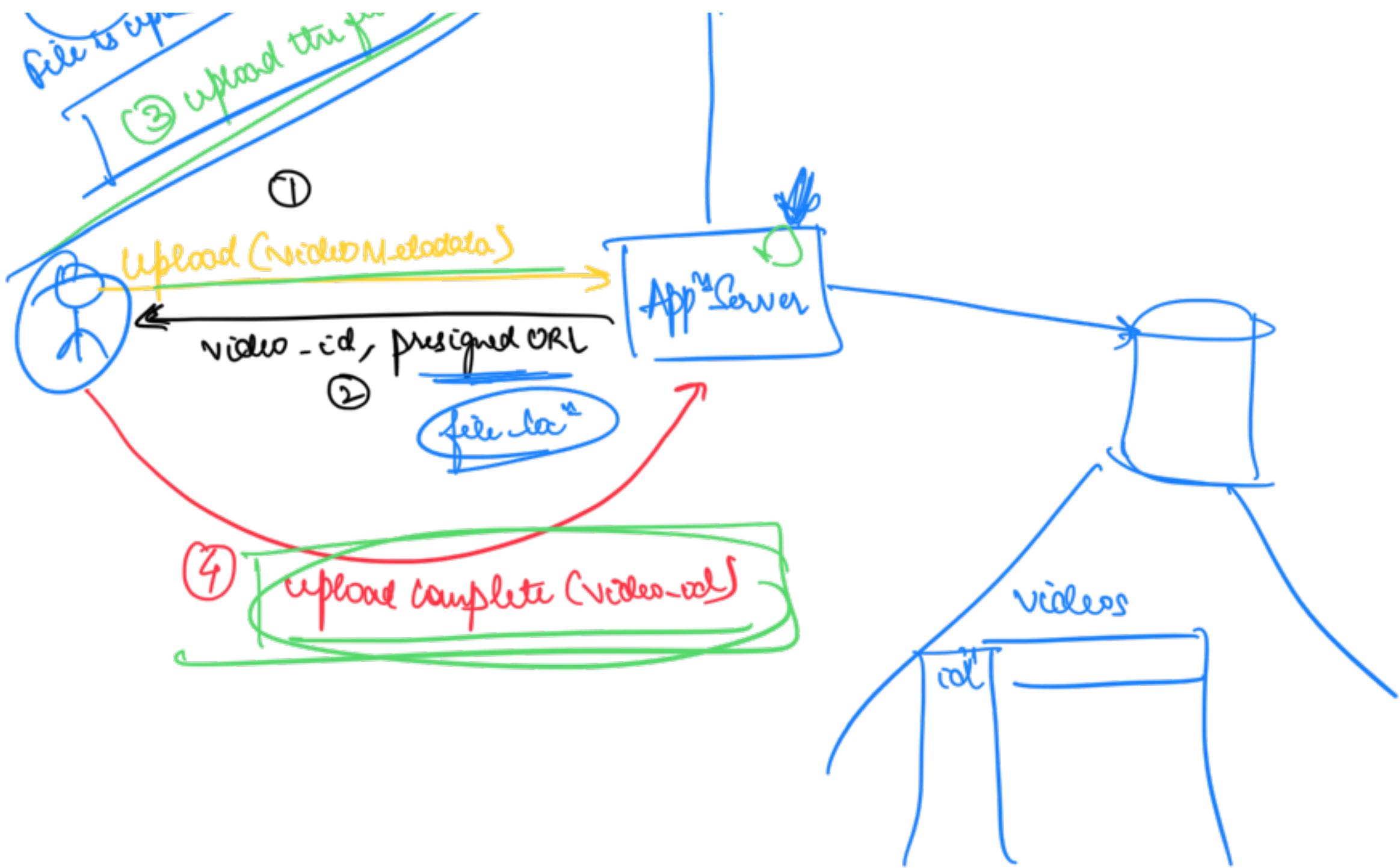
only for Lambda

Flow of video upload

- ① Send request to Server with video metadata.
 - a) Video MD gets stored in DB
 - b) Server creates a pre-signed URL and
Sends to client.
- ② Client uploads directly to AWS.

I have saved dual bandwidth usage.





When a video is uploaded,

- ① I might want to transcode it

~~compressing~~
~~changing format~~
Converting to diff results

Event Driven Systems

(SNS by AWS, Kafka, Redis)
RabbitMQ

Messaging Queue

uploadComp()

App's Server

Publisher

publishes
an event

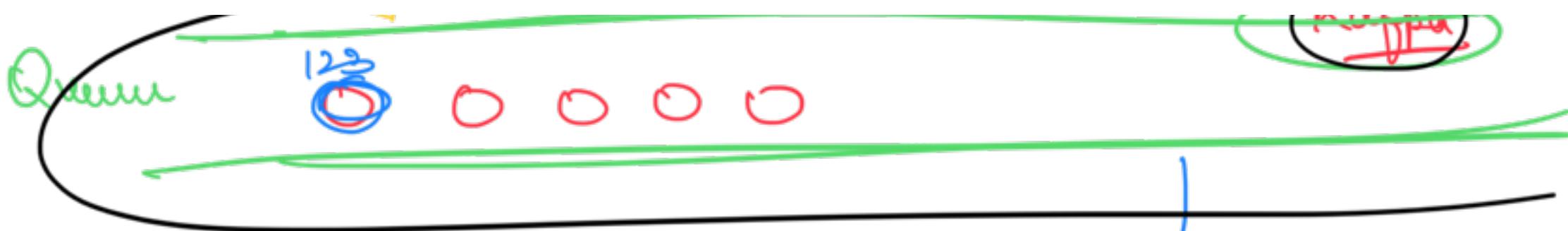
upload Video
Completed
(video.mp4)



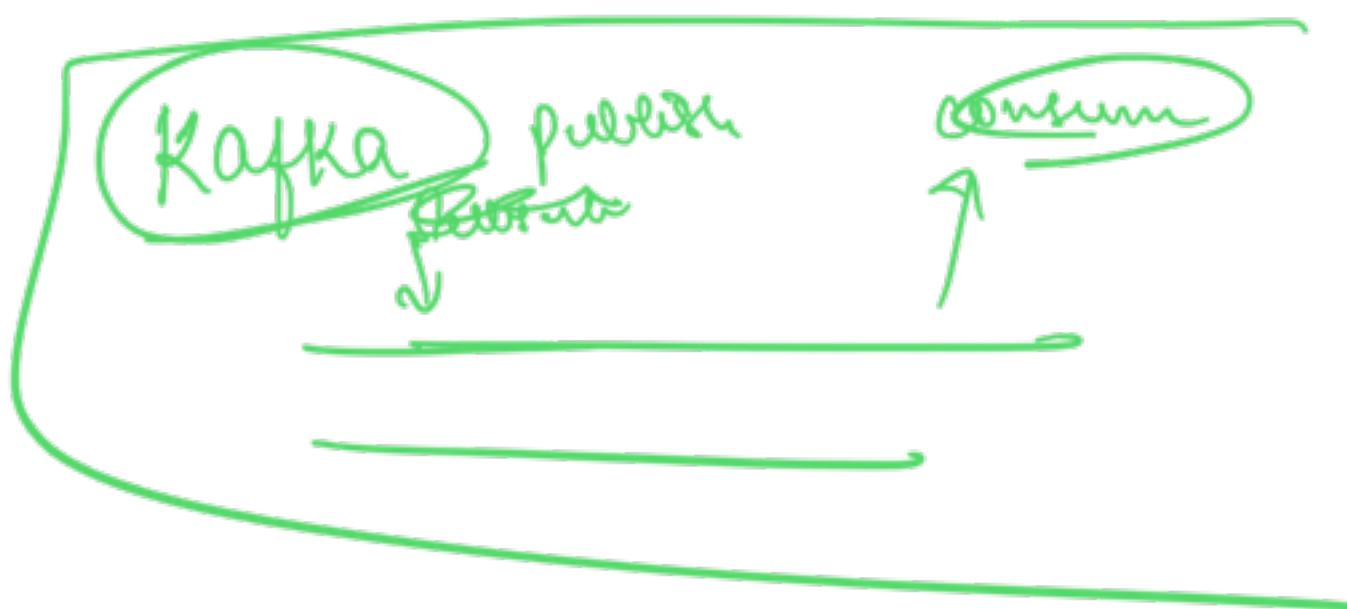
Subscriber

stop compressions





Any event will
be delivered
to ONLY
1 server
per group



What does each worker do?

① Downloads the raw video from AWS

② Compresses, converts to resolution

③ Uploads to AWS (as chunk)

④ Stores the details in DB

(along with chunk details)

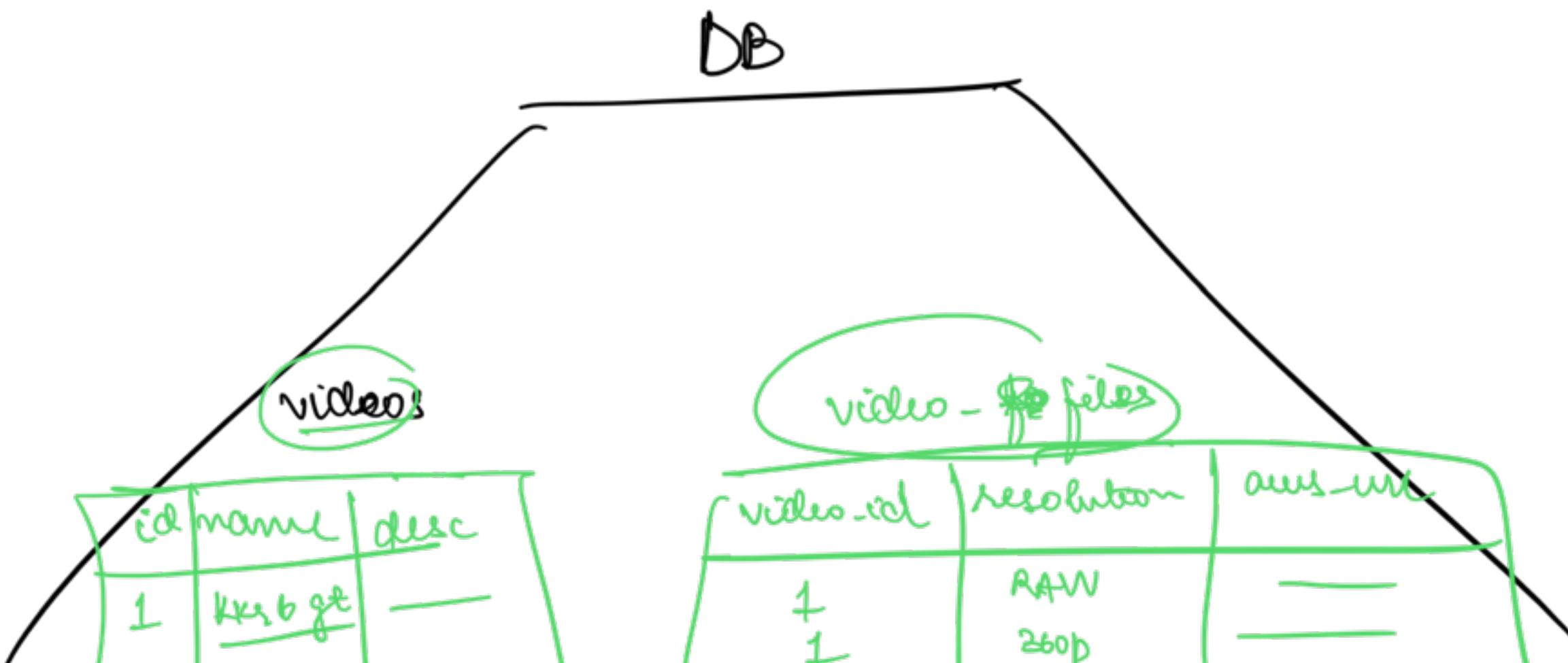
DB

Break the video
into multiple chunks

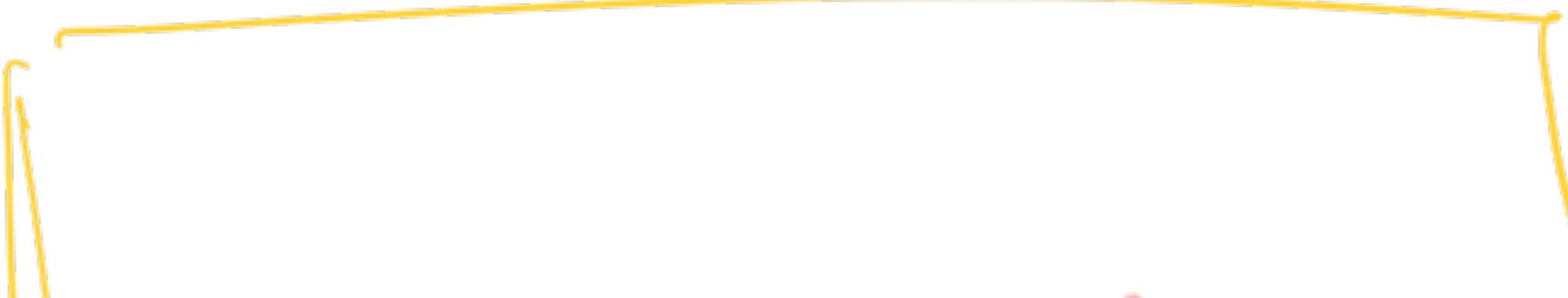
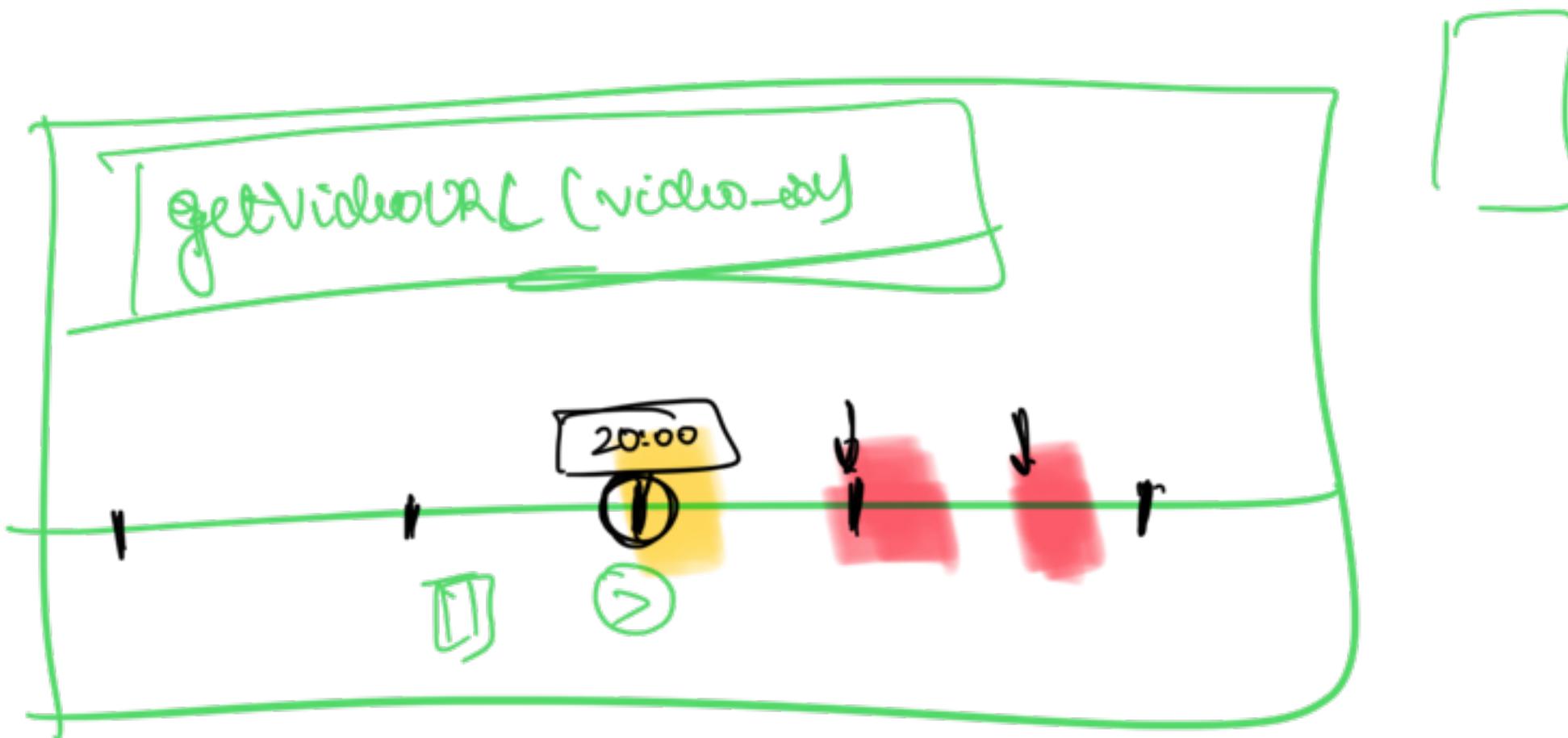
(every chunk of 10s)

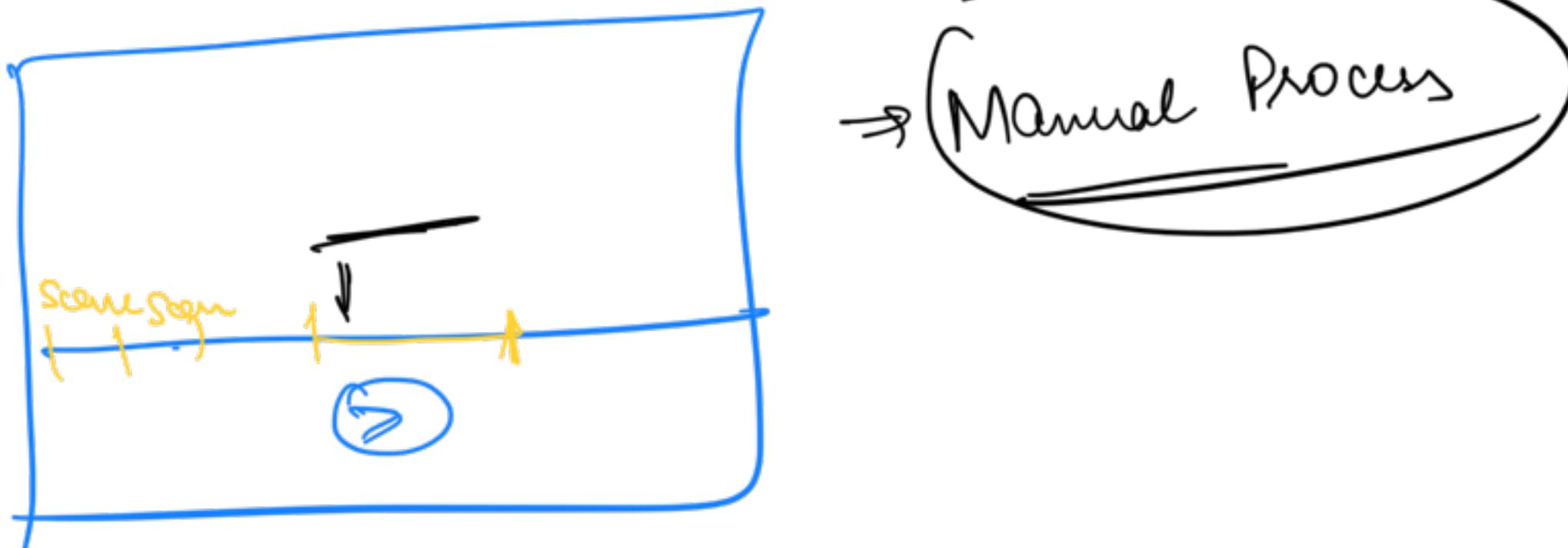
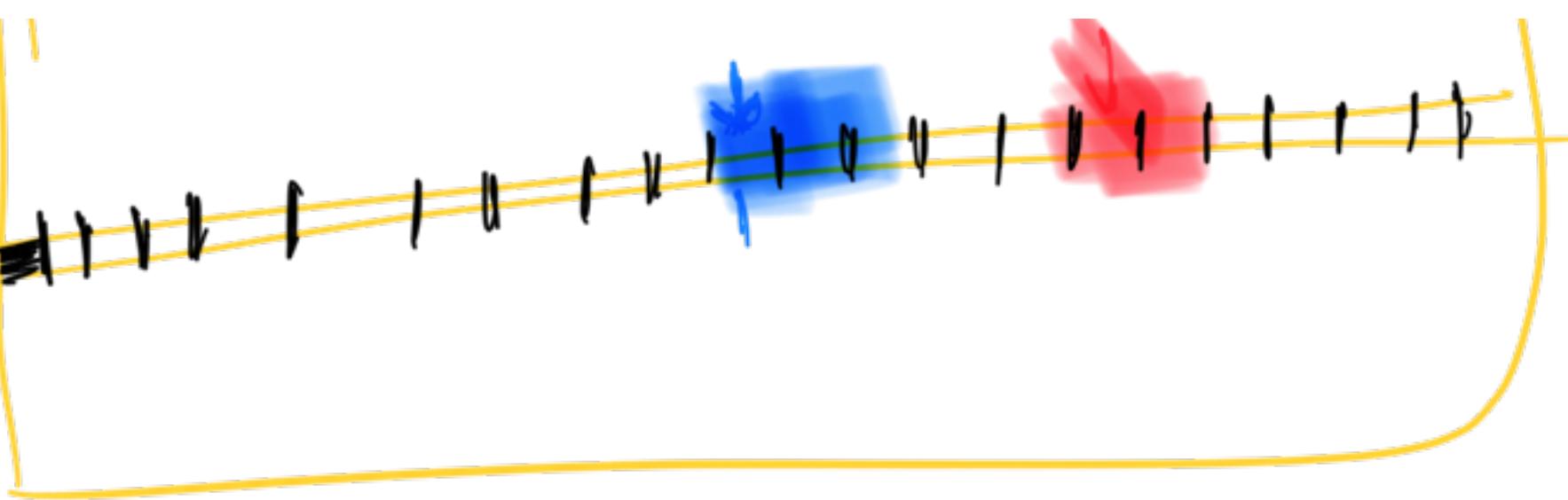
videos

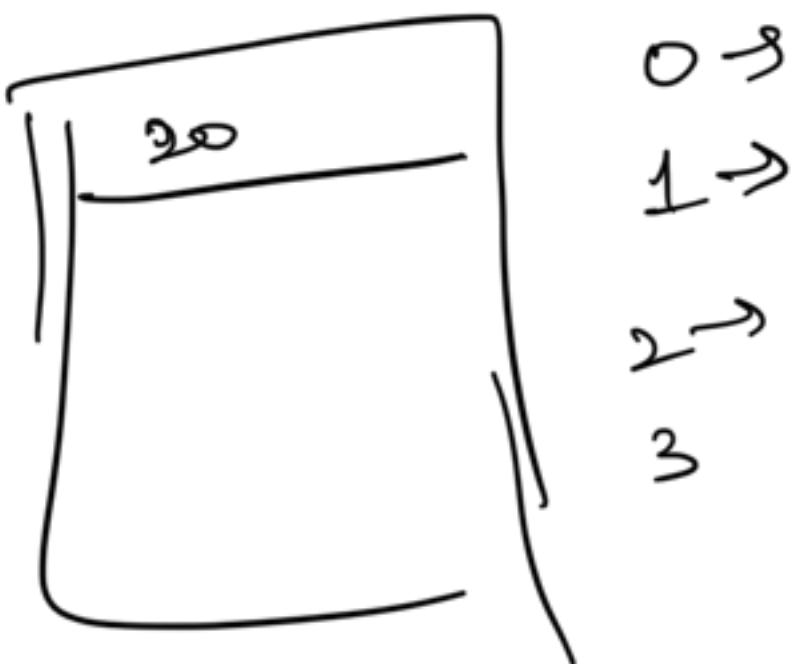
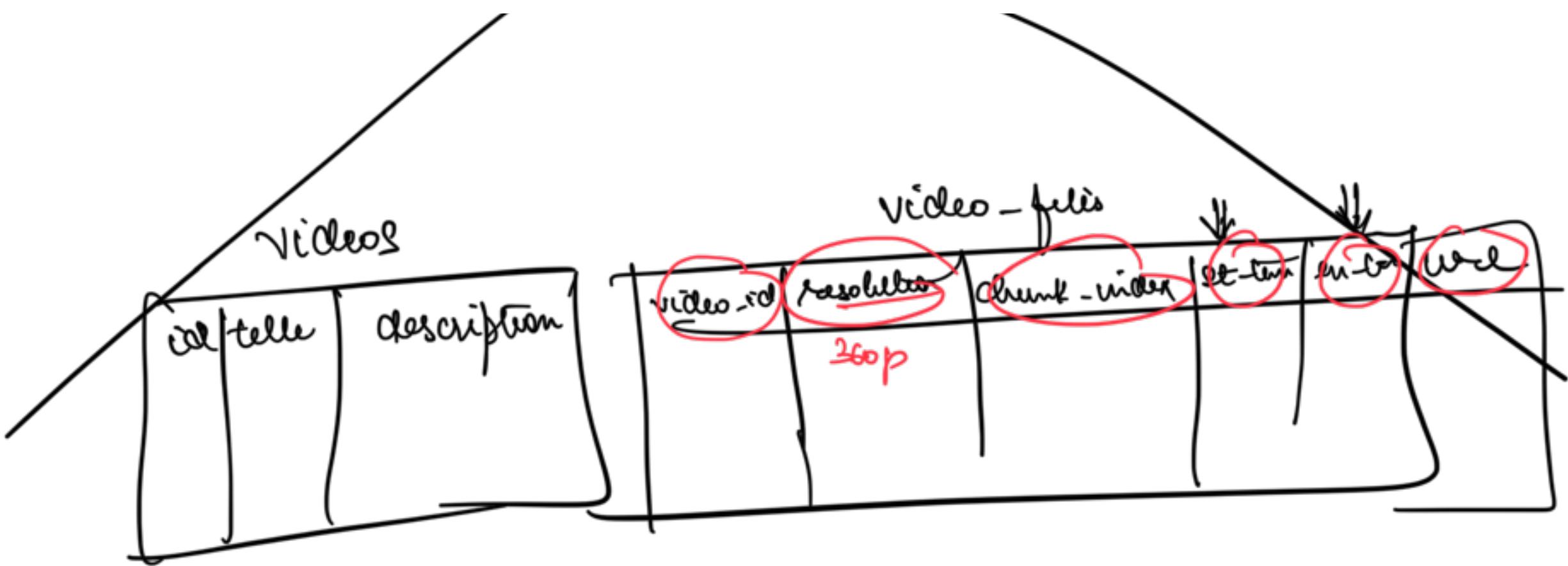
id	Name	Description	url

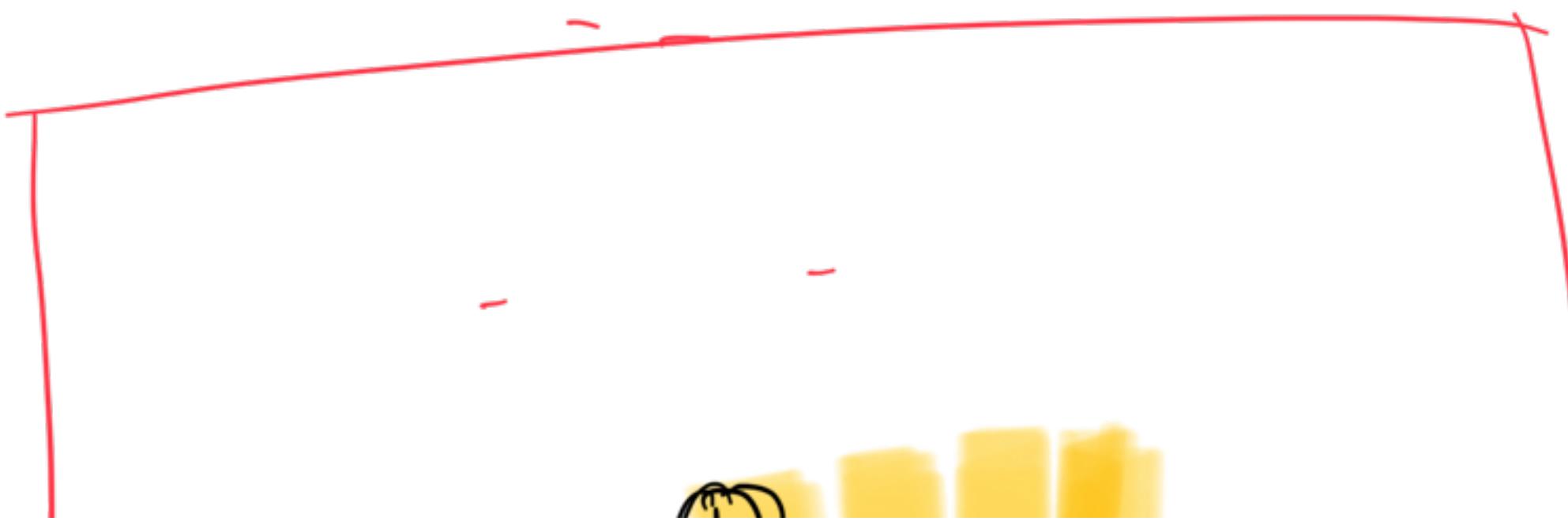
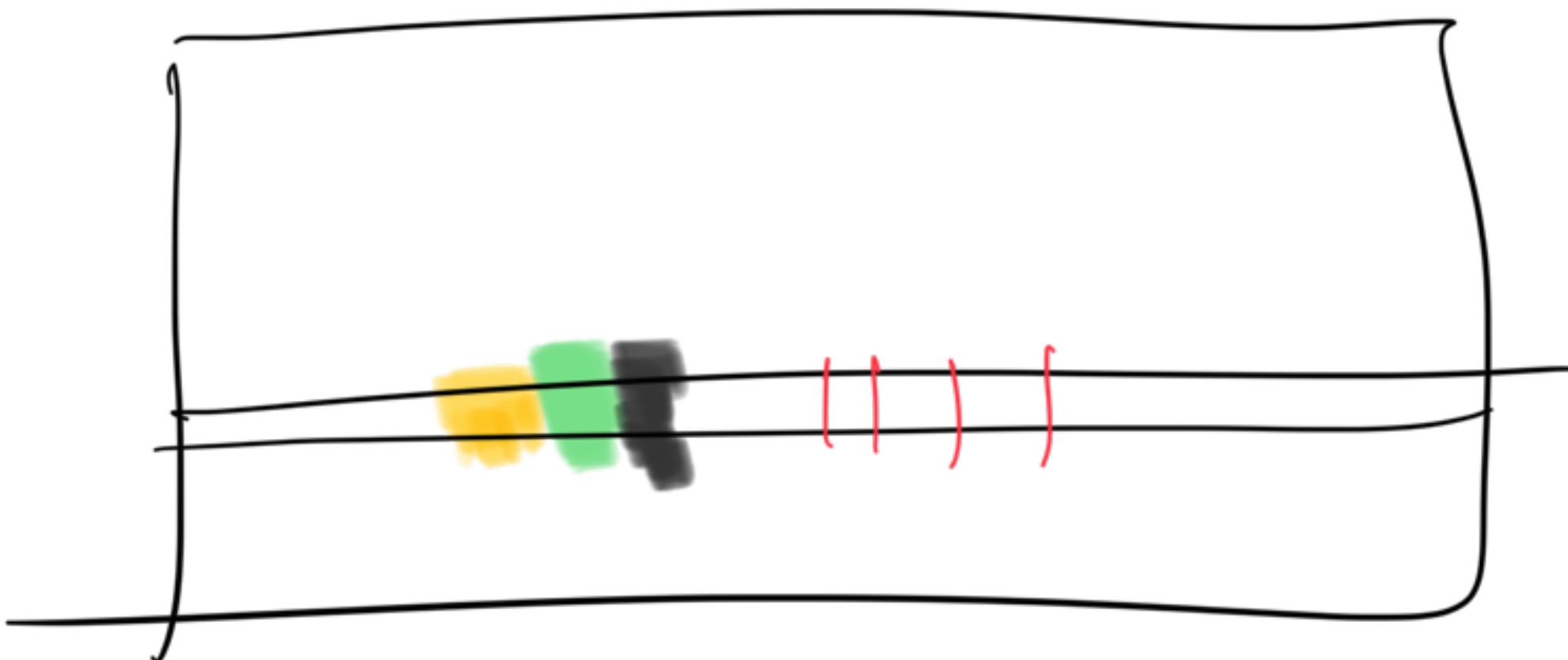


' U i | trip =







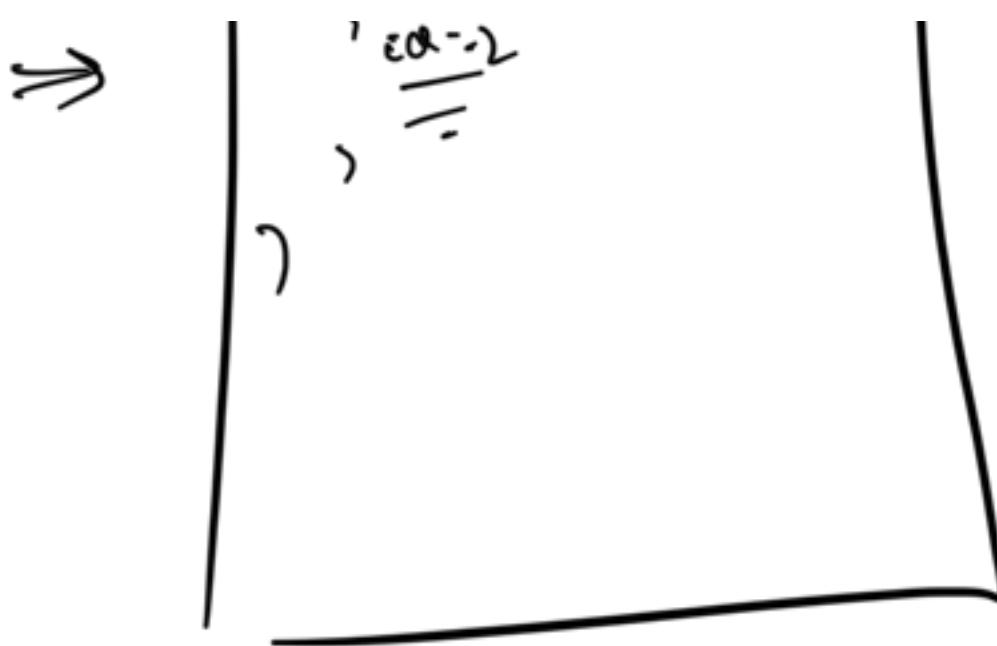




How Client Watches a Video:

- ① Client opens video Player
- ② Video Player sends Request:
get Video Details (resolution, video id)
- ③ Server sends whole lists of chunks for that video
for that resolution

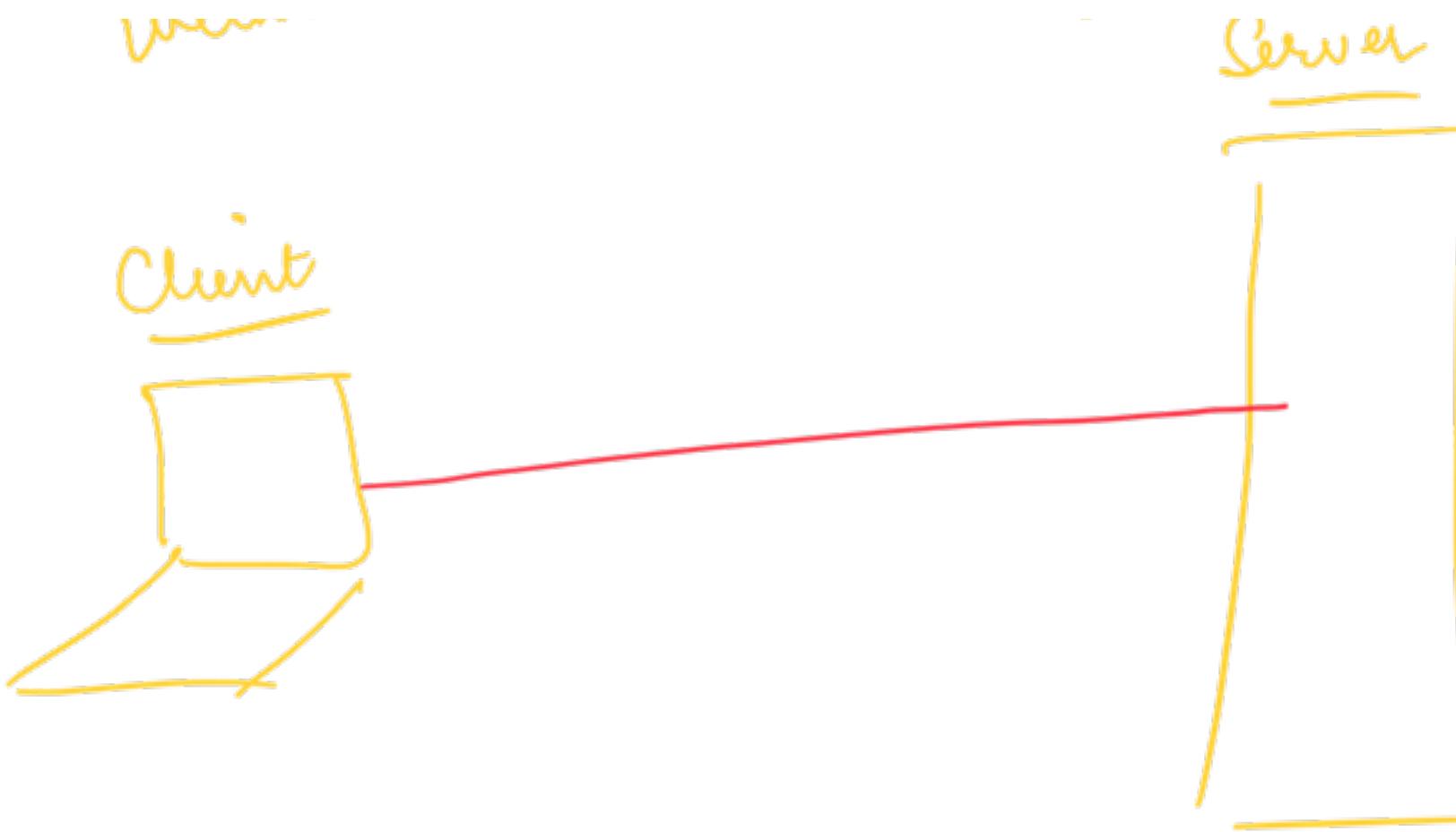
Chunks: [
? Ed = 1
St - tm = 1
end - tm =
? ure = -
]



- ④ for whatever timeframe we are at in the player
 sequentially player starts downloading ~~total~~
 chunk

Web Sockets → Picture

WebSockets allow servers to send info to client
 ... without client requesting

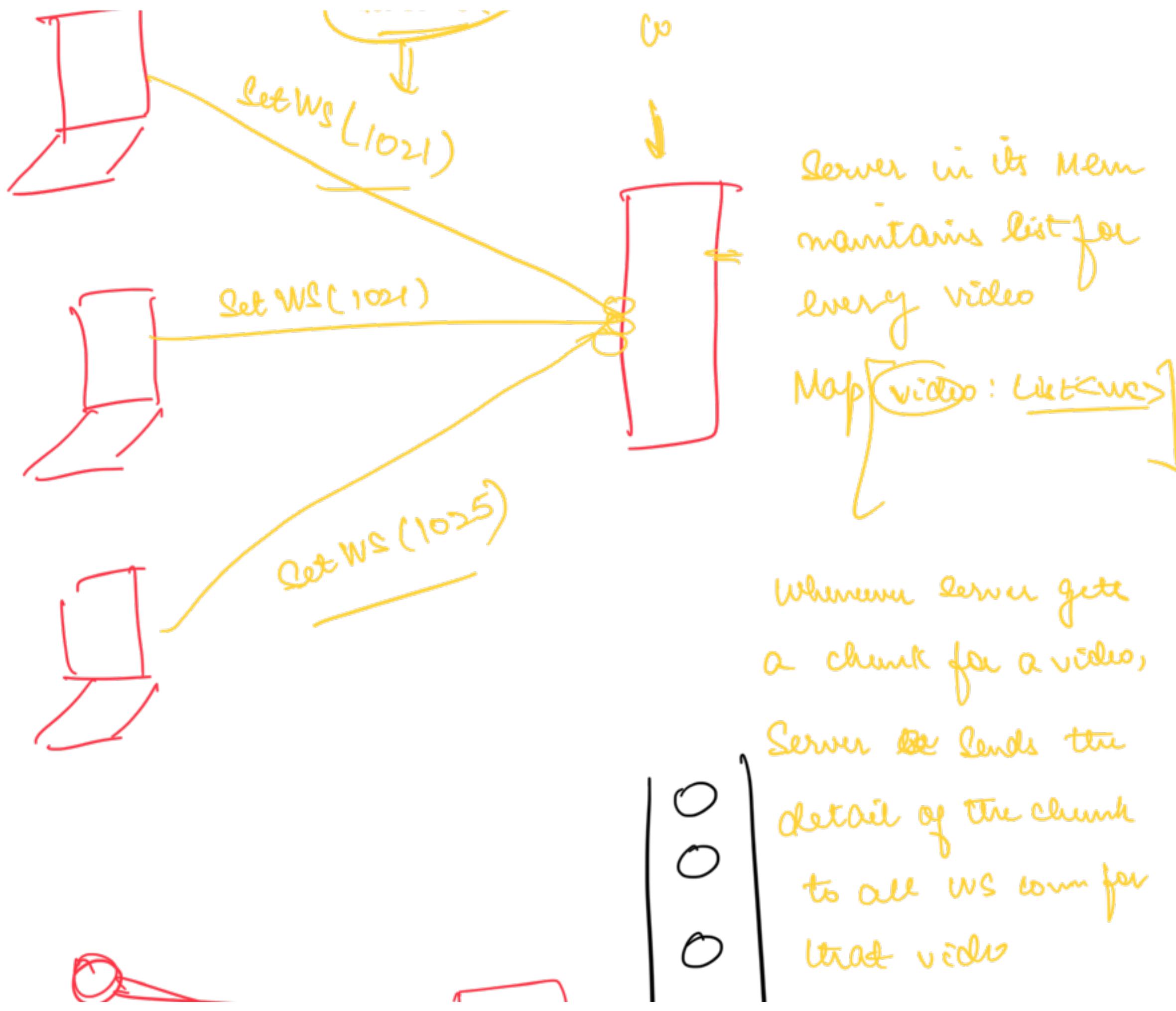


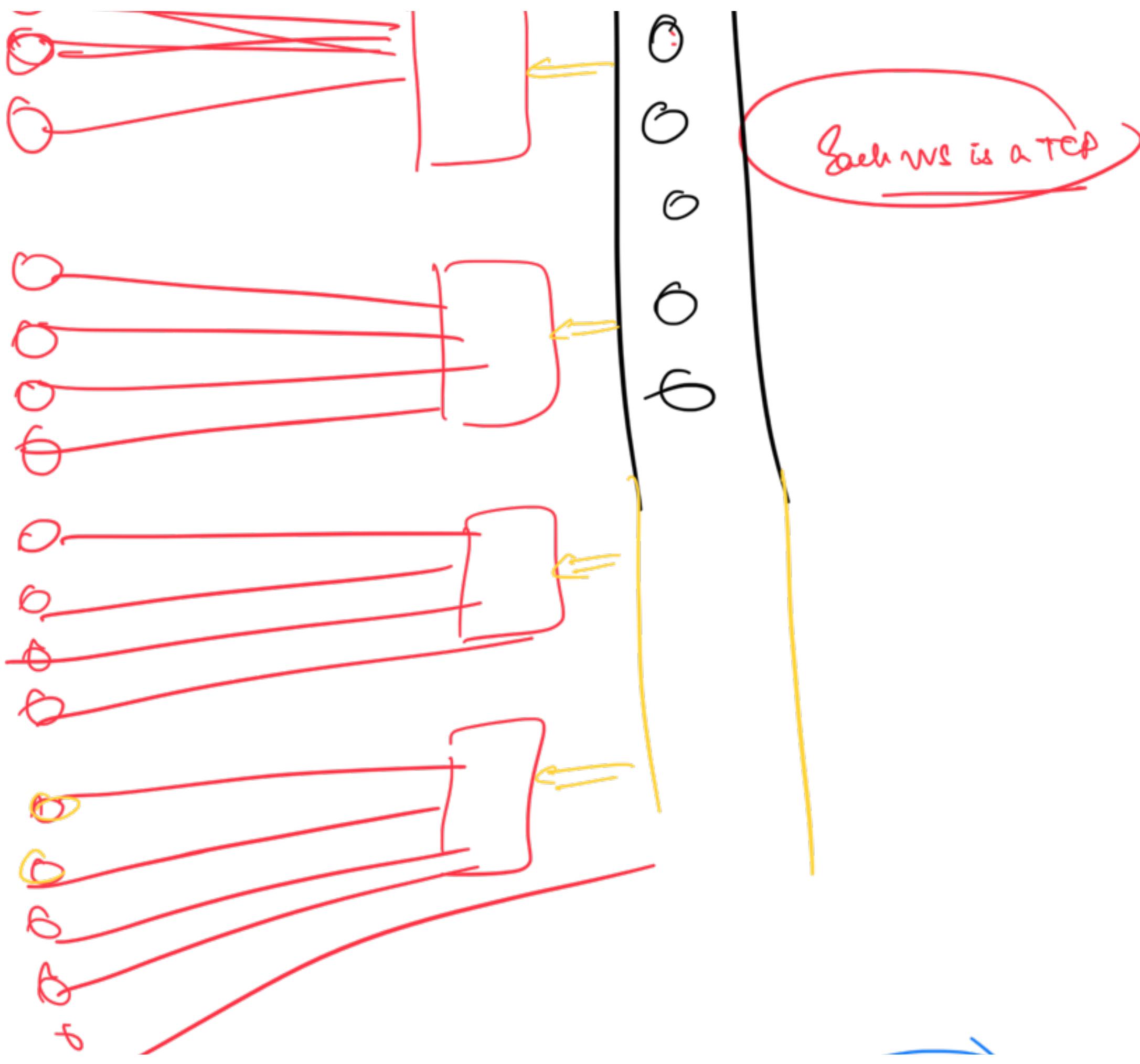
When a video player loads:

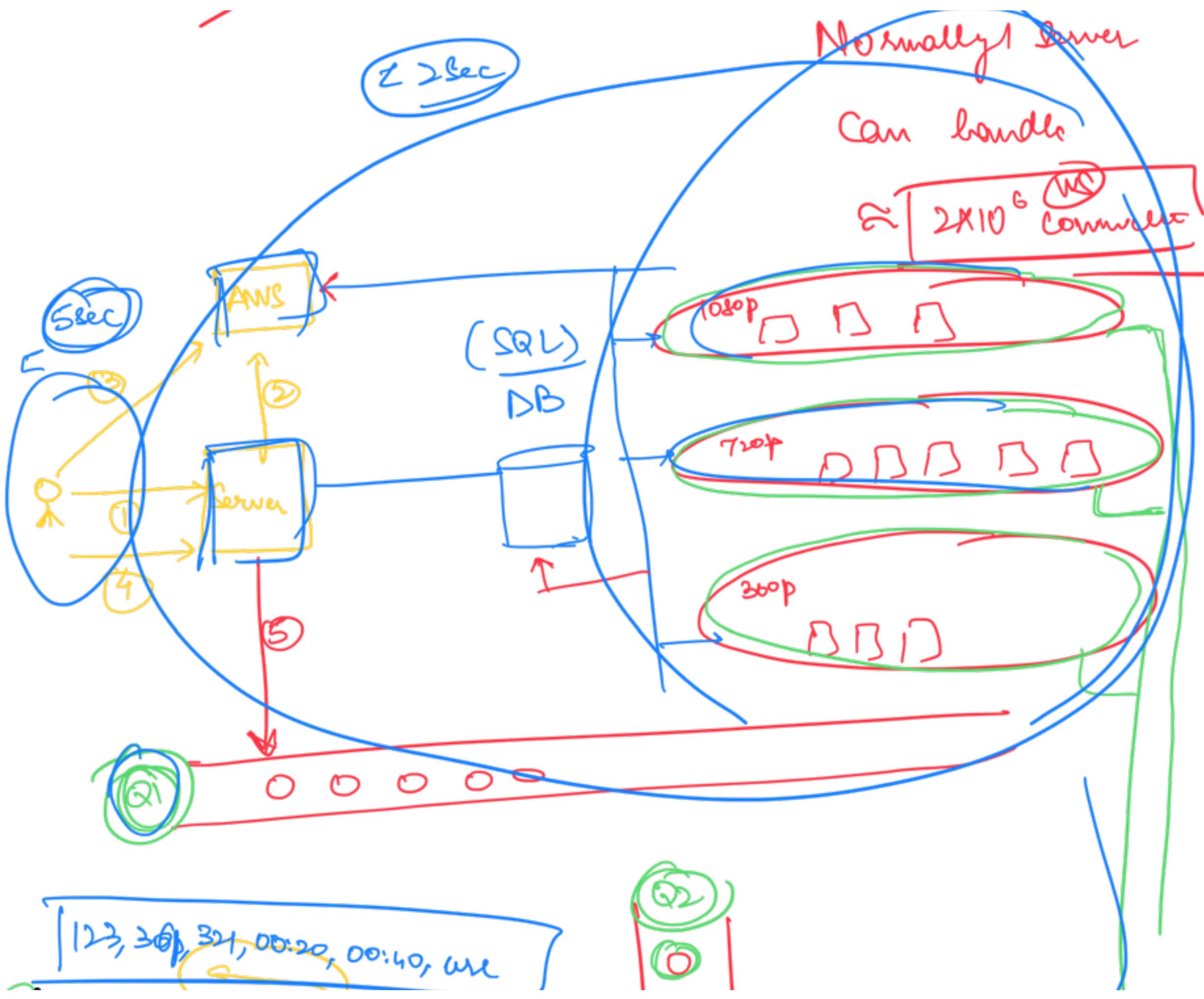
- ① Client lets a WS:|| with the server
- ② whenever server has a new chunk for the video
server sends that via a WS:|| conn

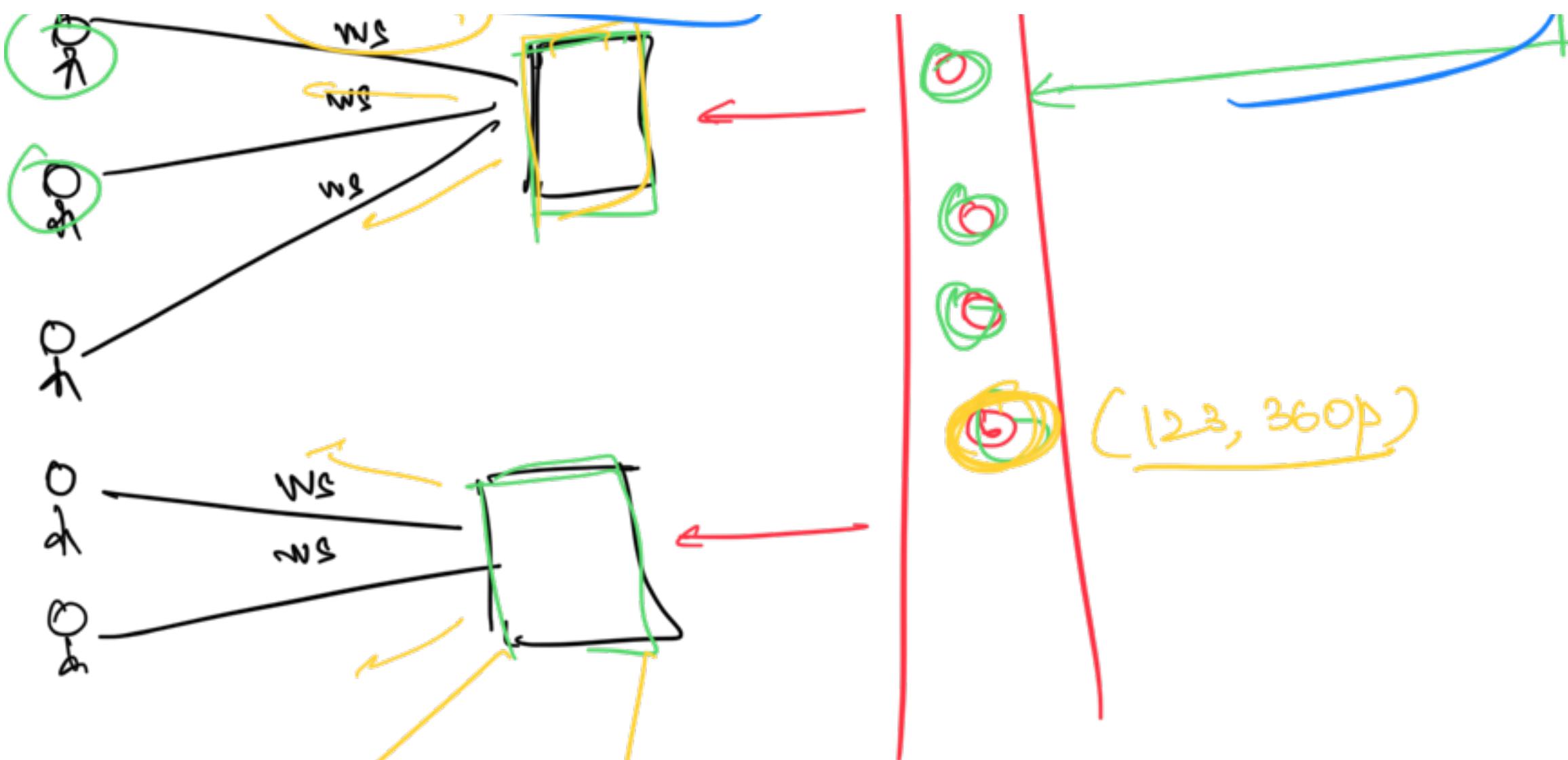


Video-id >

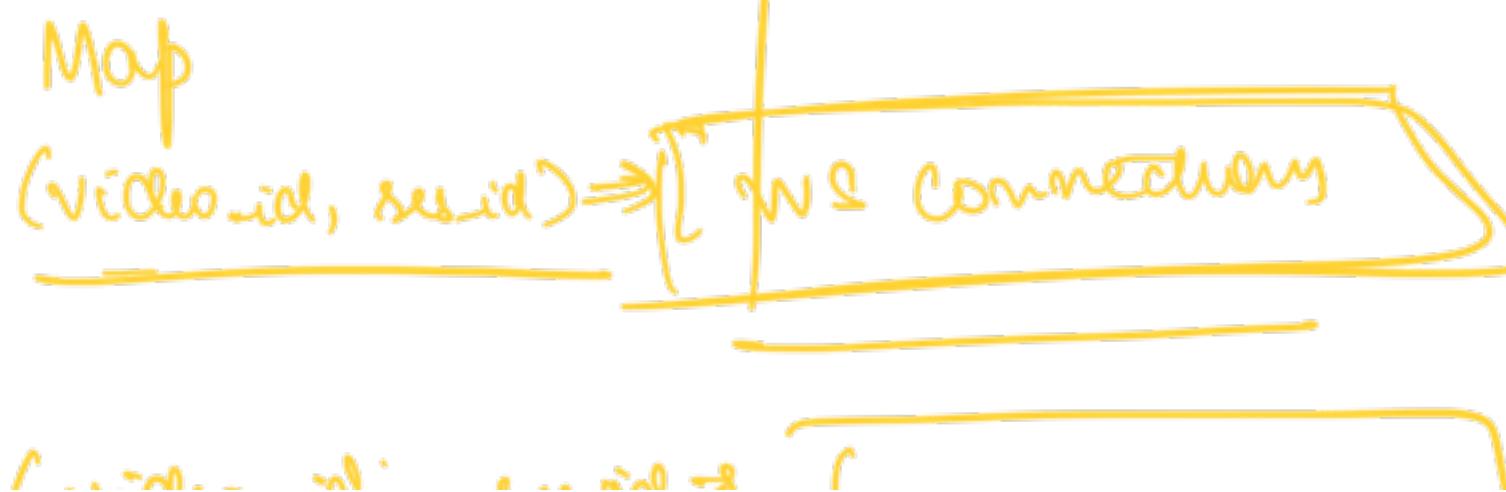








① All servers need every event



[View, Map]

