


ML System Design - 2

[Start at 9:05]

→ End-to-End Perspective

⇒ Design A video streaming platform like Netflix

⇒ 1. Ask requirements

- Provide your assumption

→ Functional Requirement

- User / Team / Brand uploading
- Stream / Interact with video
- Personalised Recommendation

Non-Functional

- High availability
- Low Latency →
- Scalability & Efficiency
- Ease of Use UX

→ HLD → High Level Design

↳ Blueprint of services / Tech / models

LCD → Low Level design

- Specificity of services in HLD

System Design →

- 1) Break down the problem → identify what services needs to be created
 - Upload Service
 - Streaming Service
 - User can search
 - View
 - Like / dislike
 - Recommendation Service

→ Tip → ask for validation.

2) Estimation/ Assumptions →

- How many videos would be uploaded
 - Total users → 10M . Total videos → 5000
 - Beginning → Most users would be in US

3) Data →

- Where the data will come from
- What all data you will have access
- Where will you store the data



Type of data →

① User data → Name, email, location, gender, subscription

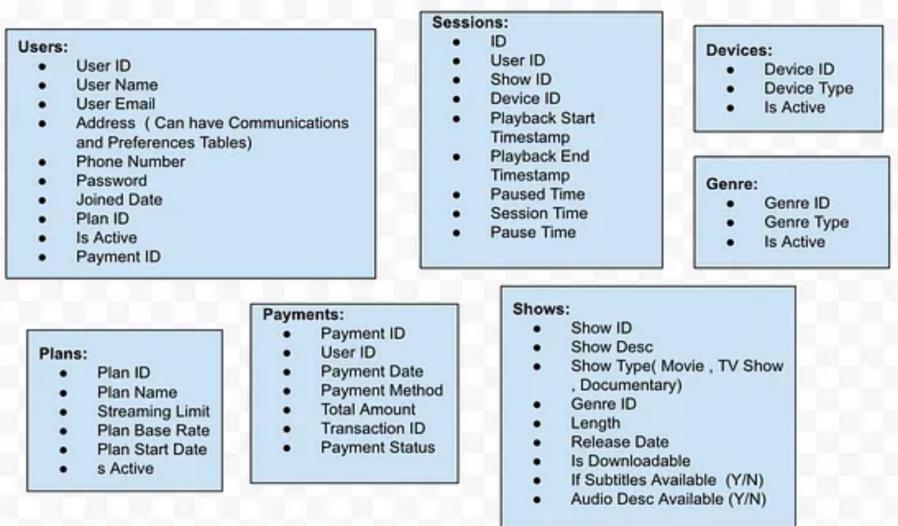
② User interaction → views, likes, time spent
• clicks, convert, time of day

③ Video data → ^{meta} genre, cast, release year
• avg rating, url, pics

④ Video →

⑤ Contextual → location, Device used,

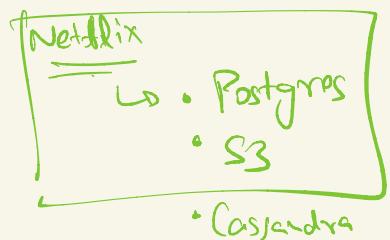
⑥ External → social media information
• reviews



① How to store data

- One SQL database → PostgreSQL
- one nosql db → MongoDB + S3

② relationship b/w data



→ Distributed databases & Warehouses

→ Optimise the Schema design

→ Caching mechanism → to reduce latency

→ AWS RDS → SQL database

• AWS DynamoDB / Apache Cassandra → No SQL database

• Apache Kafka / AWS Kinesis → Streaming data

→ Data will grow very rapidly → Scale | Scaling

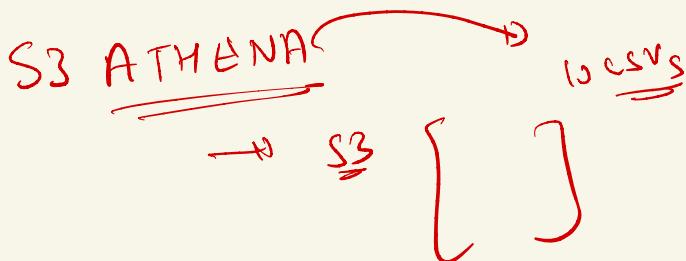
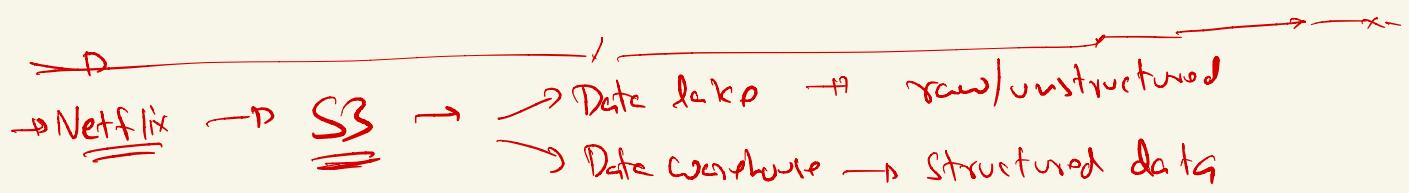
Scaling

Horizontal

- Increase the number of machine
- easier to upgrade
- Harder to implement
- Costlier

Vertical Scaling

- Increasing the machine Specs
- Harder to upgrade
↳ may have downtime
- easier to implement
- cheaper



ATHENA can query
underlying csv using sql

→ S3 is damn cheap →



→ Most (99%) services → Want

⇒ Compute & Orchestration →

→ Cloud

- Scalability
- Security
- Ease to use

vs

On-prem

↳ sensitive data

→ Most Companies → have multiple Services

→ Loan Recom →

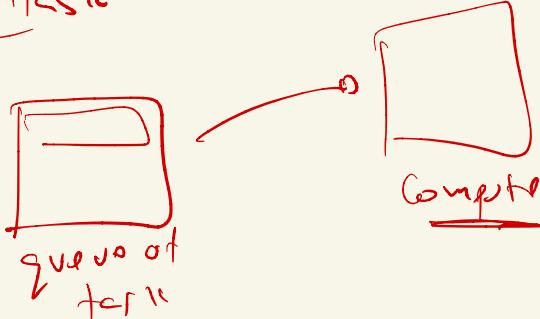
(Input user info) → Loc-Rer

→ Car price prediction

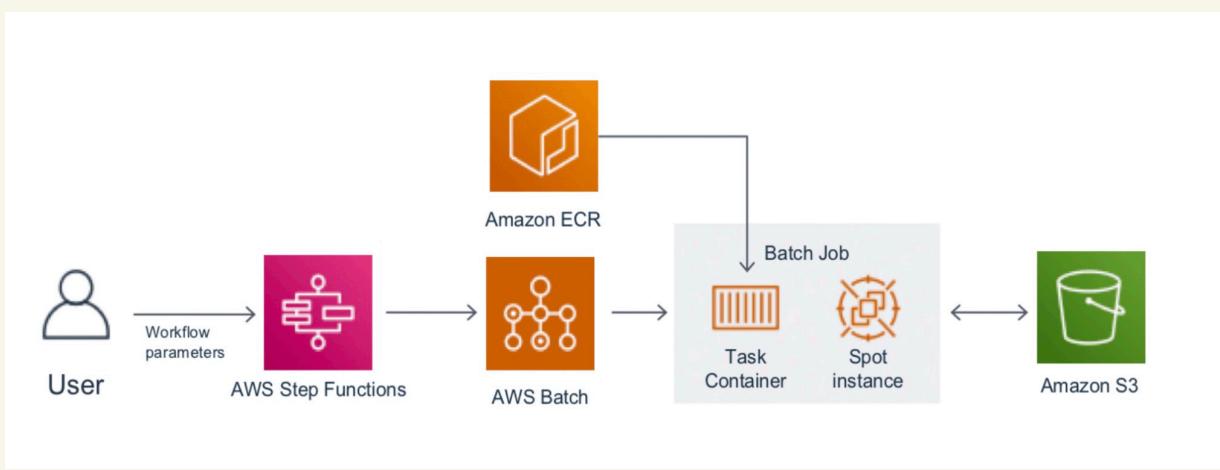
⇒ - few services → parallel
→ sequentially -

→ A lot of deep learning tasks →

→ Batch process →



for optimum usage.



Stepfunctions → workflow automation service

↳ How many times to run a service

• what resources will it take.

Aws Batch → provide necessary compute to task.

- High availability → Ensure that your service is up 100% time.
- • keep proper scaling -
 - • Load balancing



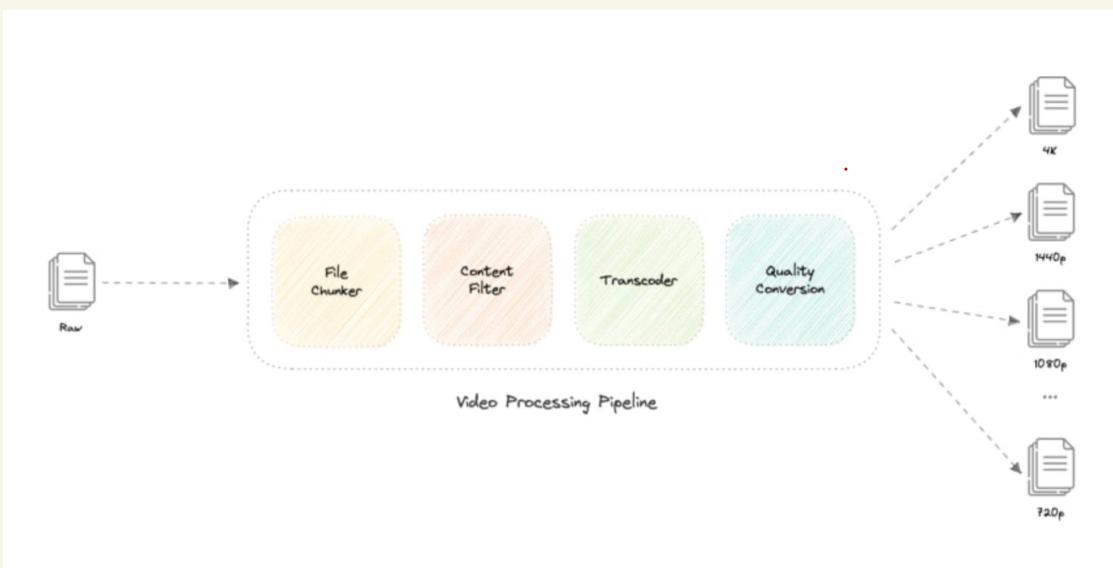
- Lower the pixel → lower the file size
- Store the data in multiple resolution -
- 

Good
Big
File.

→ Compress → serve to user 

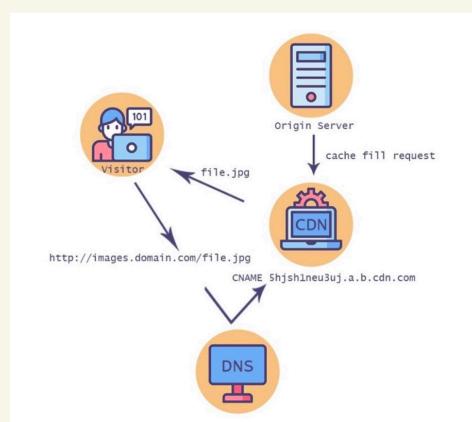
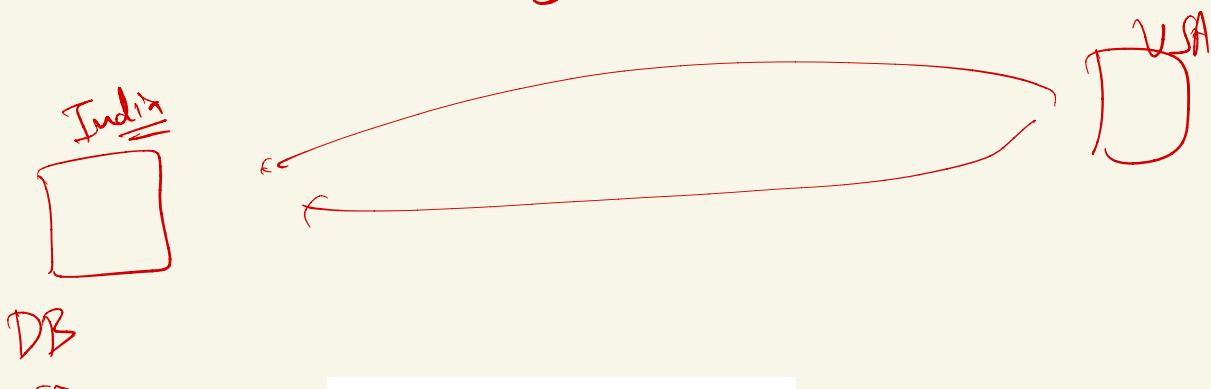
Big
file

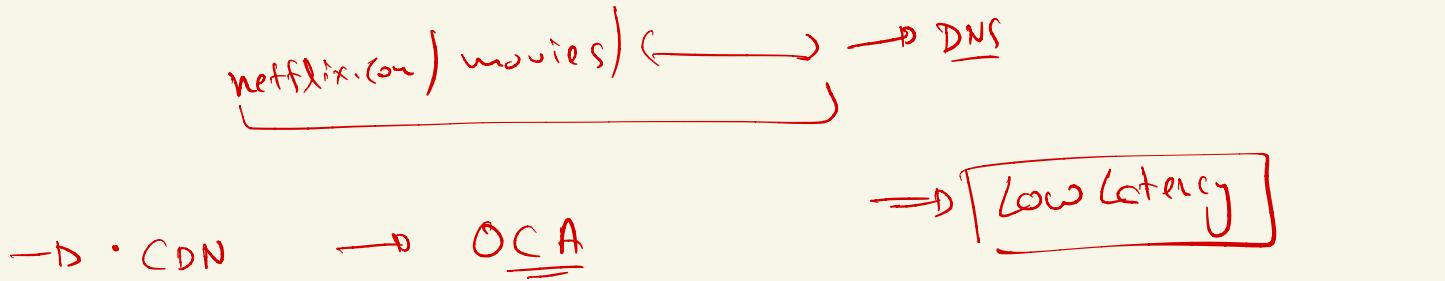
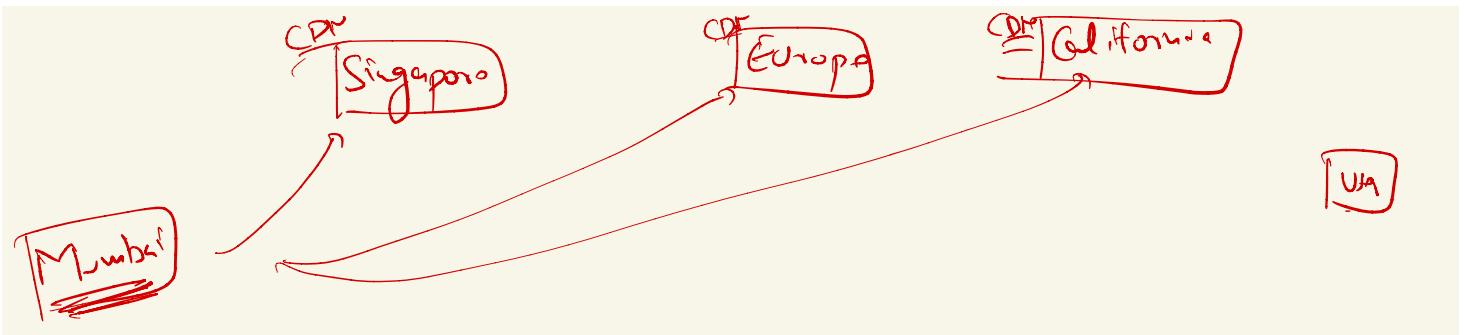
→ Compression → Store it → serve to user 



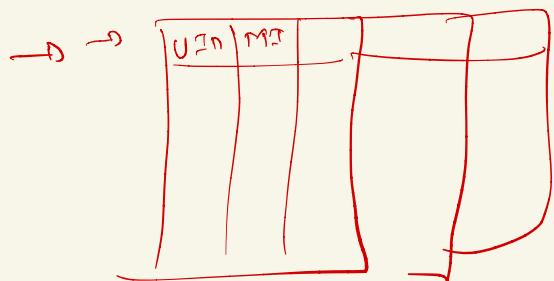
→ Speed things up → Low Latency

→ CDN → Content Delivery Network





	MOVIE - 1	MOVIE - 2	MOVIE - 3	MOVIE - 4
U ID - 1	LIKE - Yes TIMESTAMP 2:00:00 END No RECOMMENDED Yes	LIKE - No TIMESTAMP 1:20:00 END Yes RECOMMENDED No		
U ID - 2	LIKE Yes TIMESTAMP 1:58:21 END No RECOMMENDED No	LIKE Yes TIMESTAMP 2:03:00 END Yes RECOMMENDED Yes	LIKE No TIMESTAMP 1:20:00 END No RECOMMENDED No	LIKE No TIMESTAMP 1:40:00 END Yes RECOMMENDED Yes
U ID - 3	LIKE Yes TIMESTAMP 3:00:33 END Yes RECOMMENDED Yes			



→ Recommendation System →

→ ~~Algo~~

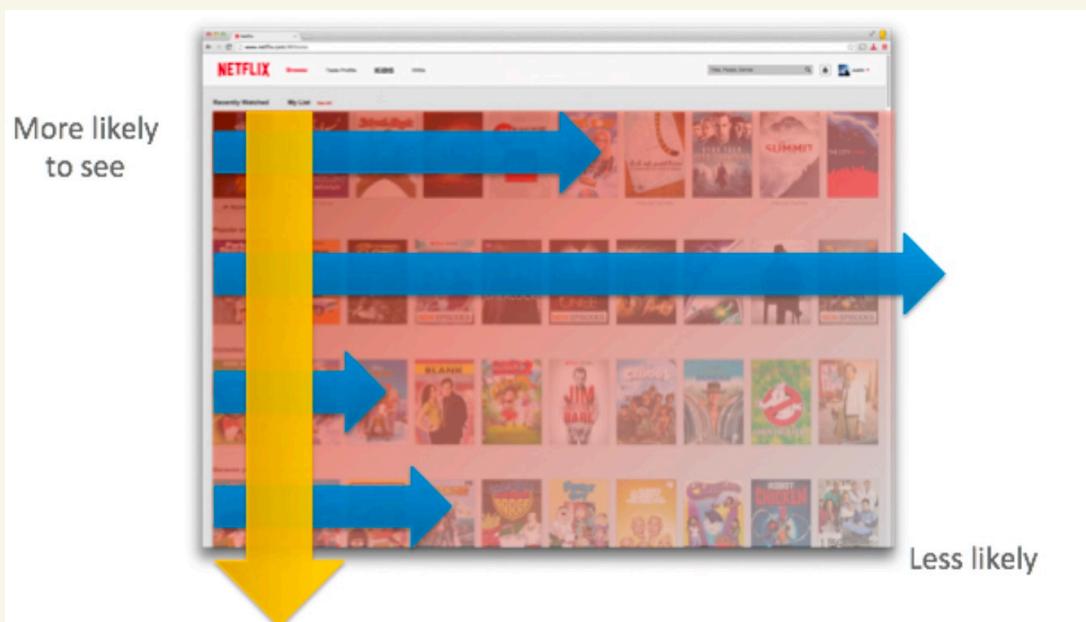
→ System Design →
↳ How many RS

→ D Row level RS

→ 1 RS to define the
Segment of Prod

Content

→ • For every row you have
are allowed
→ • Similarity → based recommendation



→ 1) Identify less popular data pts to use in model

- User interaction
- User - User Similarity
- Time of day ←
- device ←
- avg time spent ←

→ 2) Identify the algo

→ We will do AB Testing →

↓
on good performance → we will deploy it