

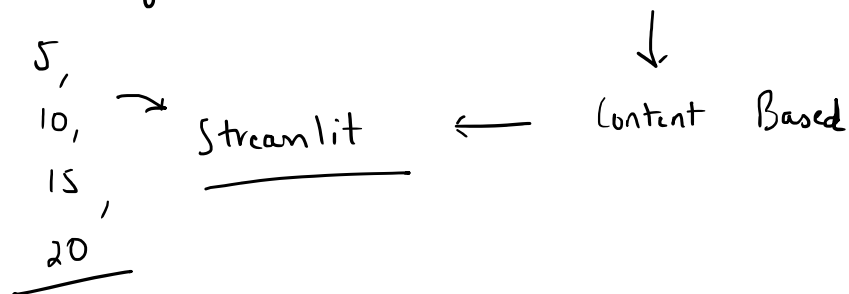
Hybrid Recommender system.

Phase I → Content Based Filtering

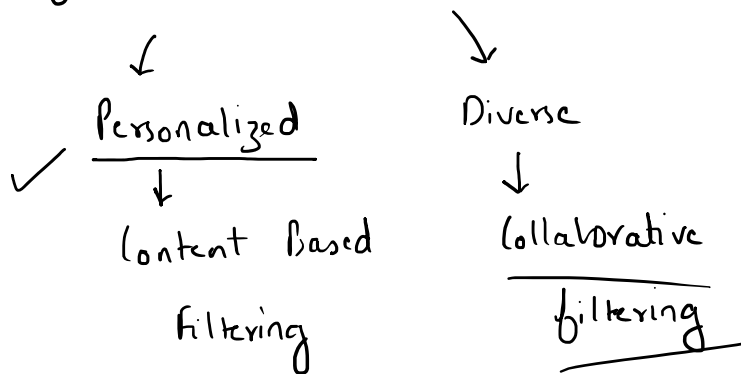
Phase II → Collaborative filtering

Phase III → Hybrid Recommender system

EDA, Data cleaning, Songs → Transformations



User engagement, User retention ↑



Collaborative filtering → User history & Content based

Collaboration → Between users → User Based CF

Between Items → Item Based CF

_____ ← Songs artist, attributes, metadata

Song A

--	--	--	--	--	--	--

 ← Songs artist, attributes, metadata

Similarity scores → cosine similarity.

Users liking, Ratings, Playcount → Similarity

User Based CF → Similarity b/w users

Playcount → User history ←

	Song A	Song B	Song C	Song D	
User A	5	4	2	0	4 dim
User B	—	—	—	—	
User C	5	3	2	1	
User D	—	—	—	—	

cosine similarity = Sort Similarity
↑ ↓

User Based Approach / CF → Items >> Users. ✓

Item Based Approach / CF → Users >> Items. ✓

	User X	User Y	User Z	User M	User N
Song A	—	—	—	—	—
Song B	—	—	—	—	—
Song C	—	—	—	—	—

Song 1 —————

Song 1) —————

Item - User Interaction Matrix \rightarrow Item Based CF

Rows \rightarrow Unique songs

cols \rightarrow Unique users.

Song \rightarrow Vector Matrix

Similarity Scores.

$$(5, 10)$$
$$\underline{(1, 10)} \quad \times \quad \underline{(10, 5)}$$

Song A Song A

 $(1, 5)$

Recommend

song A song B

Preview

History Database

30,000 songs

9.7 lakhs

track Id

User Id

Playcount

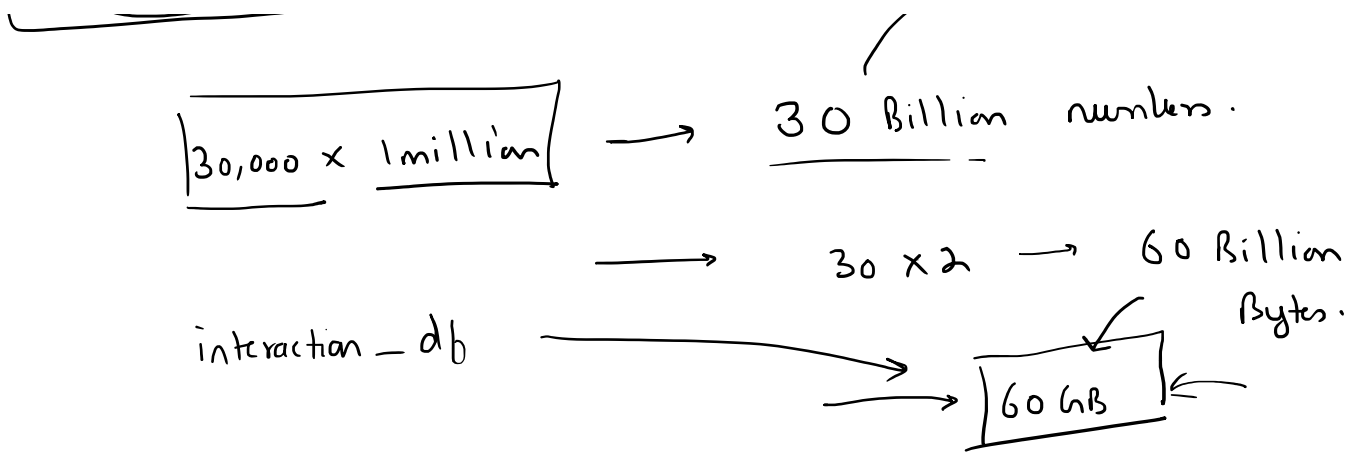
1 core

User Item interaction matrix

(30k, 1million)

2 byte

$$n_{p \cdot \text{int}}(16) \cdot (8)$$



Dask → Big Datasets → TBs

$V_{x \times r}$ (30k)

Songs → 1 million. 10000

Values → Zeros → sparse data. ≪ dense

→ CSR matrix ← cosine similarity

Non zero

Why Dask?

Similar to Pandas and numpy.
dataframes → Pandas
arrays → numpy

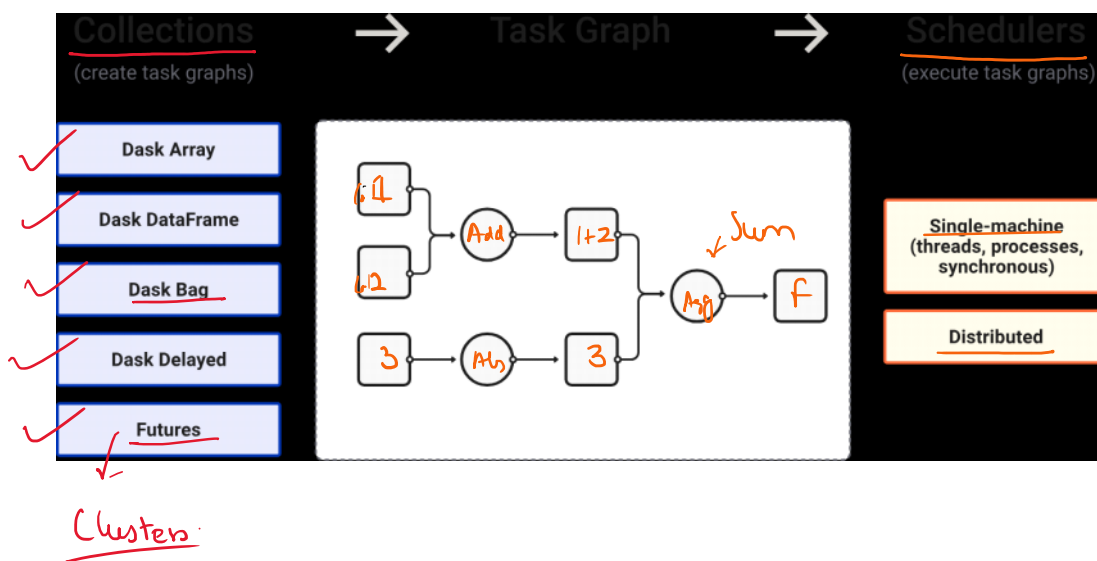
Easy to learn
Easy to implement

1) Familiarity

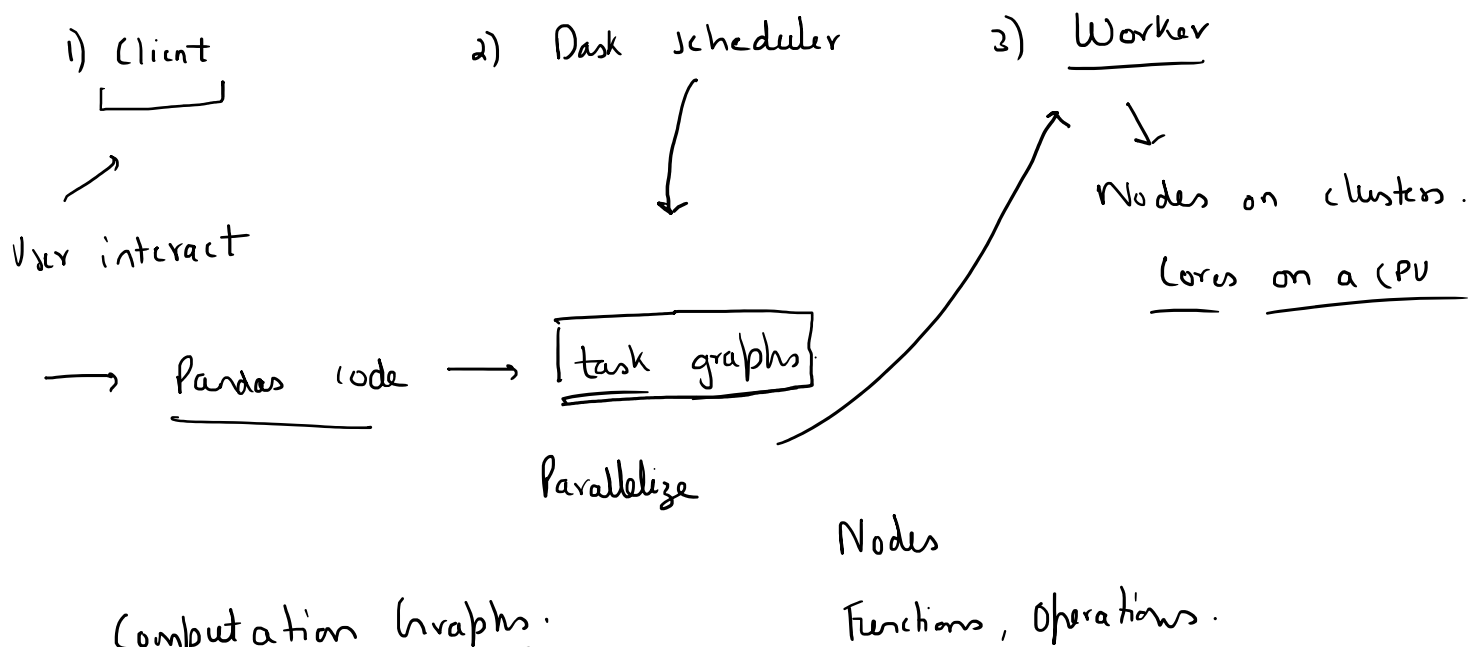
2) Parallel Computing → Multicore
→ clusters.

3) Chunks

4) ML, DL, Boosting Algo
XgBoost, light GBM

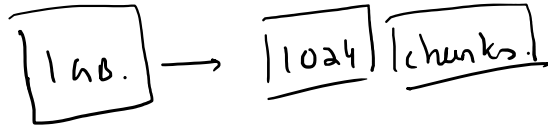


Dask dataframe, Dask Arrays → Collections.



Computation Graphs.

Partitioning of data → chunking



Pandas, numpy → Graph

Chunks. ↙ → sparse matrix

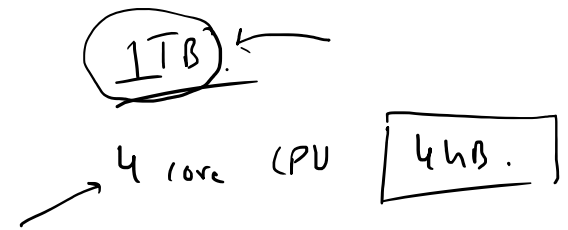
↓

Similarity Scores.

↓

Recommendations

Functions, Operations.



Understand the flow

User history → 10 million rows.

track Id User Id Play count

30k , 1 million → Item Based CF

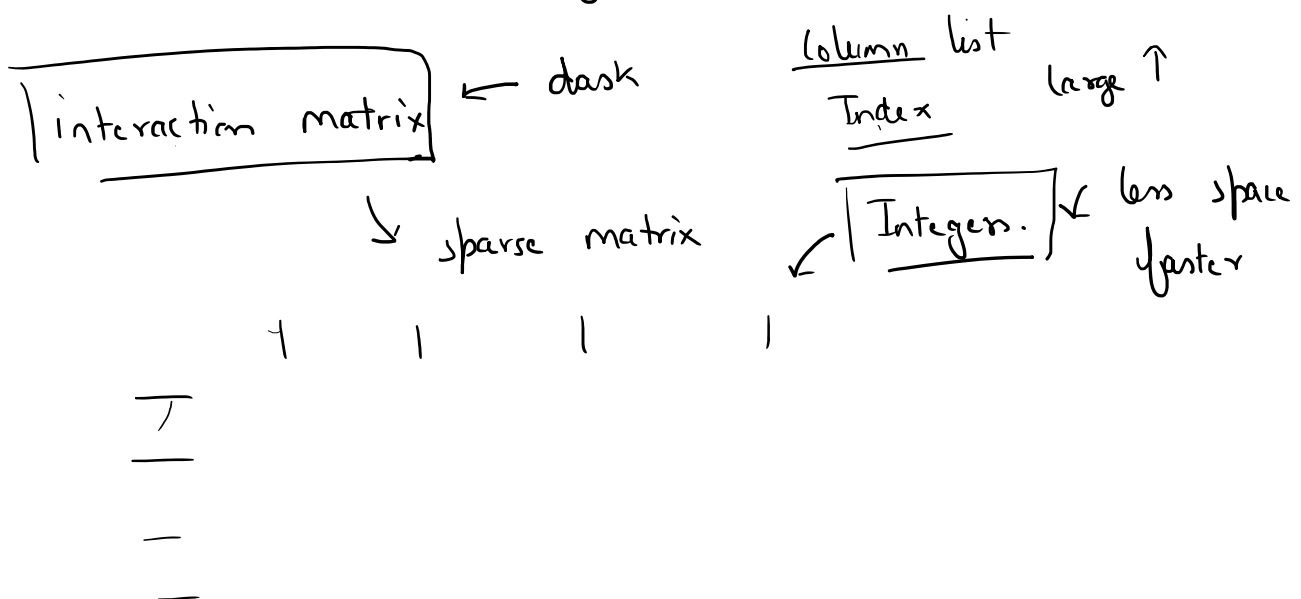
Matrix → Interaction matrix

	User Id 1	User Id 2	3	4
Song A	—	—	—	—	—	—
Song B						
Song C						

60 GB RAM → chunking Disk

matrix → sparse matrix size ↓ Zero

cosine similarity.



Songs datasets → 50k ← Pandas

User history → Dask dataframe

↓

No. of unique songs → 30k

No. of unique users → 1million

unique track ids

song dataset 30k → track id, name, artist, url
↓
streamlit

Interaction matrix → Indexes → trackid
Columns → Userids
Values → Playcount

<u>track id</u>	<u>Userid</u>	<u>Playcount</u>
-----------------	---------------	------------------

Integers encode dask specific
chunk → 9 divide 10 million → User, track

track id 1 ----- 1million

dask → unique songs > Categories
unique users
categorize → [30k] → [0, 1, 2] encode index
[1million]

user history → add new cols.

97 lakh
10 million
5 cols

track id	User id	Playcount	track index	User index
ABC	x	0	0	0
ABC	y	0	0	1
ABC	2	-	0	2

(ABC, x) → add

(ABC, y) → add

(ABC, 2) → add.

10 million track index User ID index Playcount

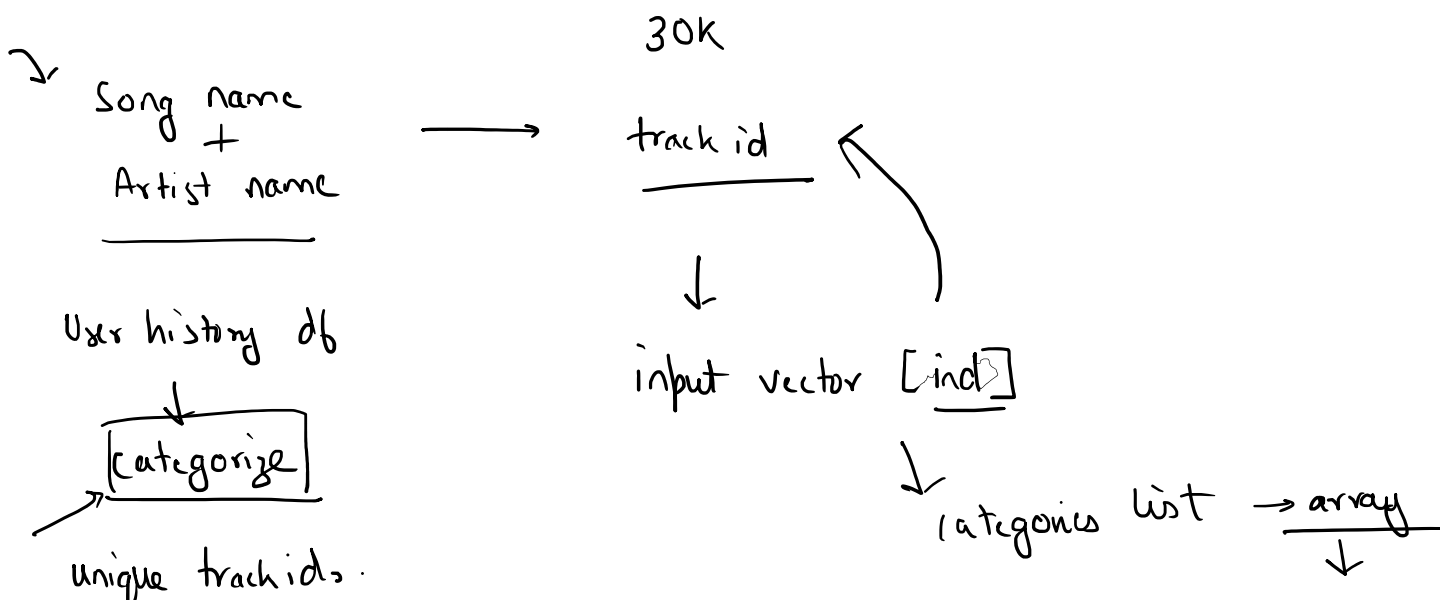
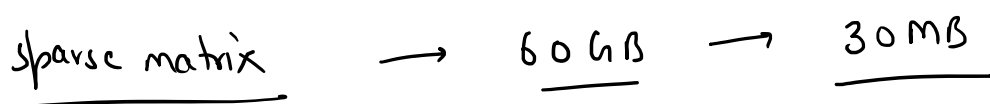
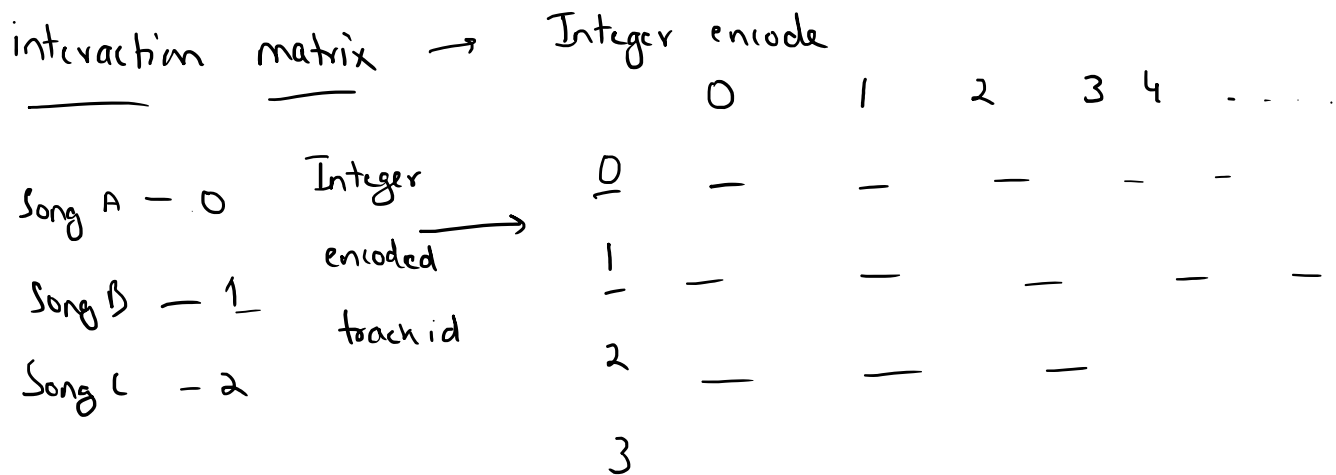
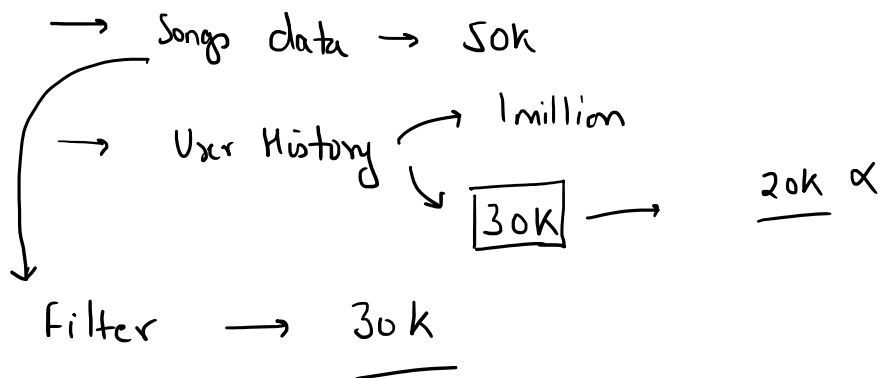
sparse matrix

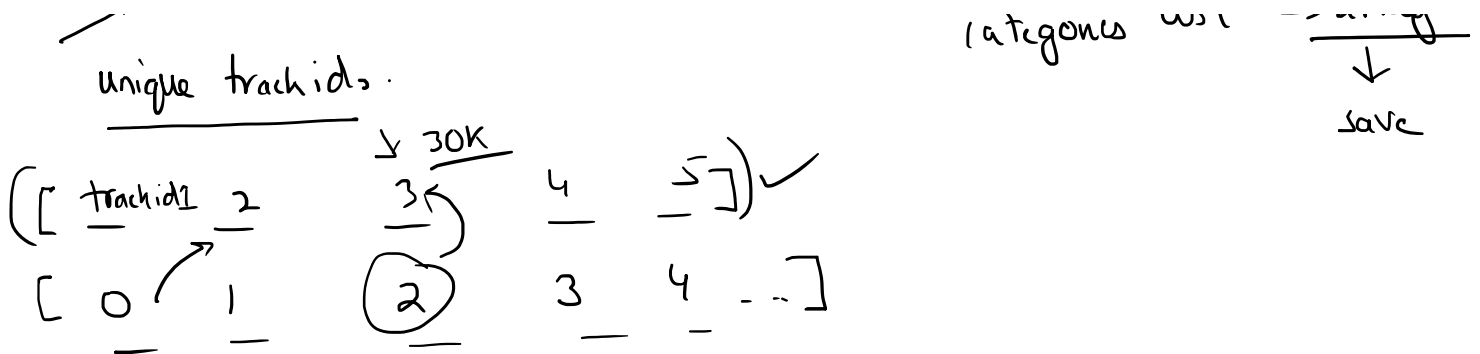
	User 0	User 1	User 2
track 0	—	—	—
track 1	—	—	—
2	—	—	—
3	—	—	—

sparse matrix

Building Pipeline

09 January 2025 15:29





cosine similarity (input vector , interaction matrix)

→ sort Descending → Topk similarities

↓

indexes.

trackid → index → input vector → topk indexes.

↓

top k track ids.

array[top k indices]

↓

top k track ids → Songs dataset → 30k

↓

name, artist, url ← get the song details

Collaborative filtering → (Song name , artist name , 30k songs data , categories array , interaction matrix)

→ 8 MB

↓ small

↓ 30 MB

Song name + → Row from data → track id

