# System Design for Recommendations & Search

⊛ **Model**

    — Deep Learning

                — Neural Collaborative Filtering.
                — Transformers for Recommendation

⊛

⊛ **Model Design** — How it fits in wholistic design?

           → Move from Batch to Real-Time

                           ↓

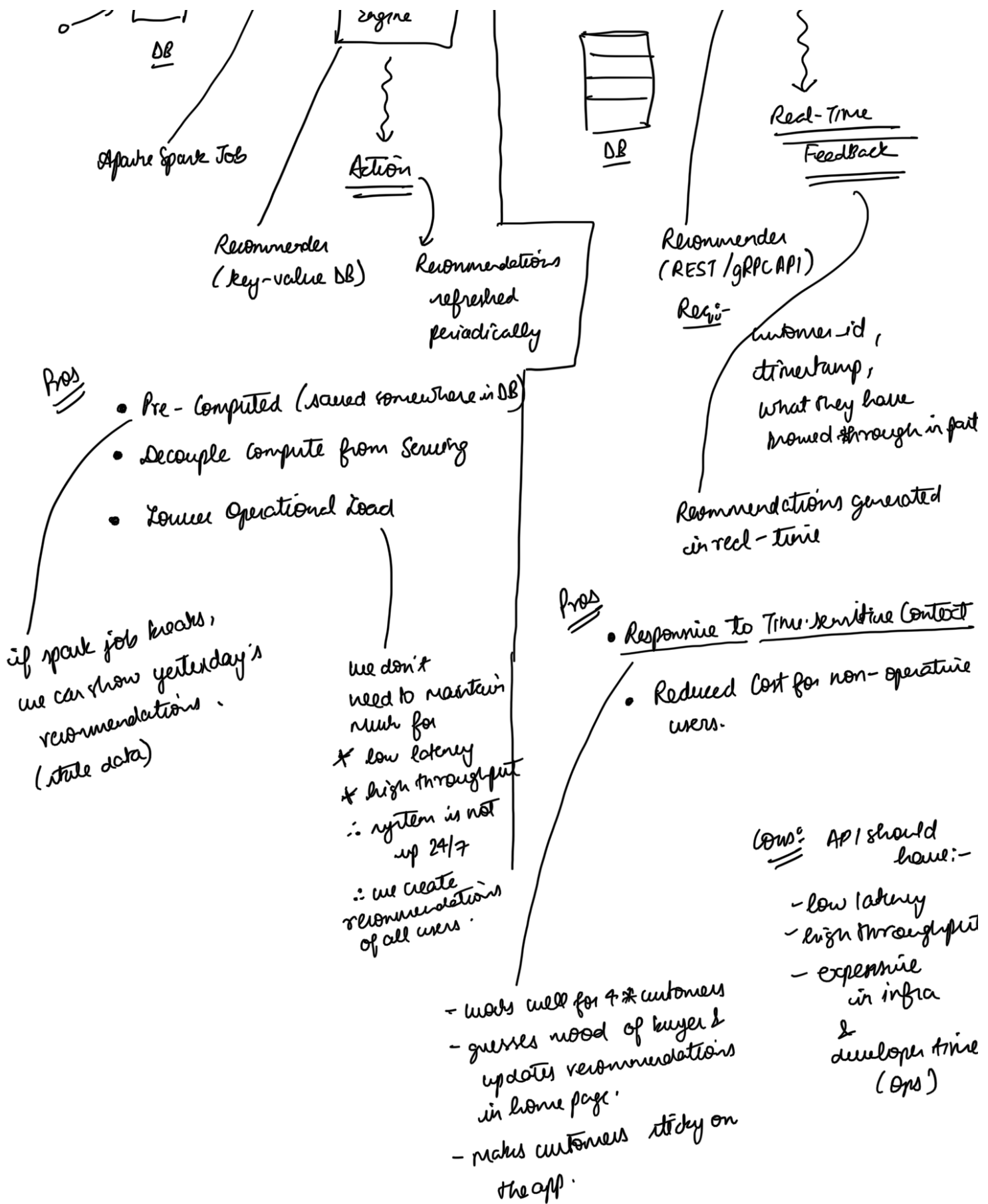                         has constraints

⊛ Basic Fundamental Design that applies to 90% of all
Recommendation Based search.

How we do this in Industry?

                — Candidate Retreival

                — Ranking

## BATCH v/s REAL-TIME



BATCH

REAL-TIME

Data Points

every 2days

Analytis

data point

Analytis Engine

DB

Apache Spark Job

Engine

Action

Recommender
(key-value DB)

Recommendations
refreshed
periodically

DB

Real-Time
Feedback

Recommender
(REST/gRPC API)

Requ-

customer_id,
timestamp,
what they have
browsed through in past

Recommendations generated
in real-time

**Pros**
- Pre-Computed (saved somewhere in DB)
- Decouple compute from serving
- Lower Operational Load

if spark job breaks,
we can show yesterday's
recommendations.
(stale data)

we don't
need to maintain
much for
* low latency
* high throughput
∴ system is not
up 24/7
∴ we create
recommendations
of all users.

**Pros**
- Responsive to Time-sensitive Context
- Reduced Cost for non-operative users.

**Cons:** API should have:-
- low latency
- high throughput
- expensive in infra
& developer time
(Ops)

- works well for 4* customers
- guesses mood of buyer &
updates recommendations
in home page.
- makes customers sticky on
the app.

REAL-TIME /ON-DEMAND /ONLINE

offline env. to
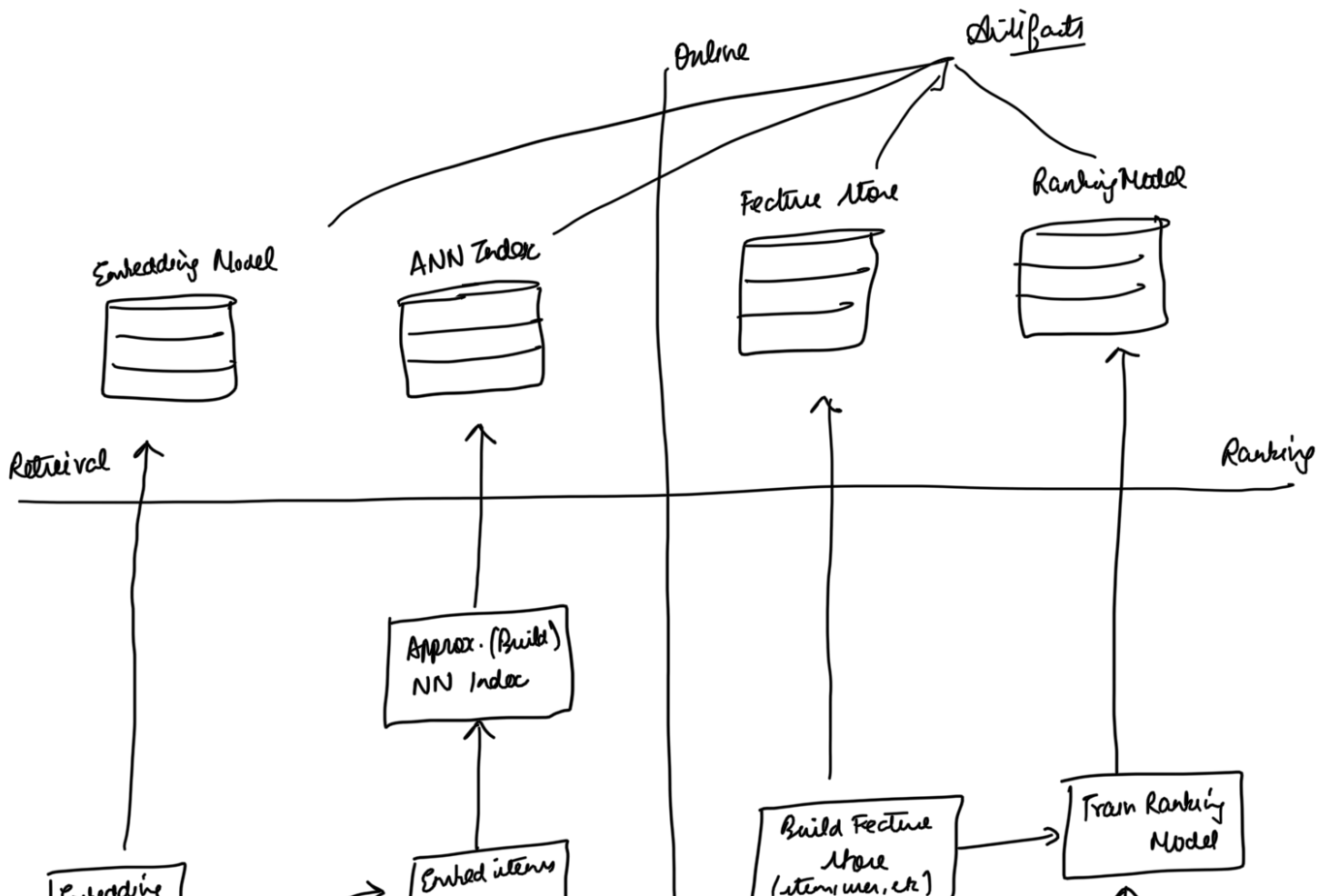
## BATCH / OFFLINE

- Most Batch Processes
  ( Training / Indexing /
  Graph Building )
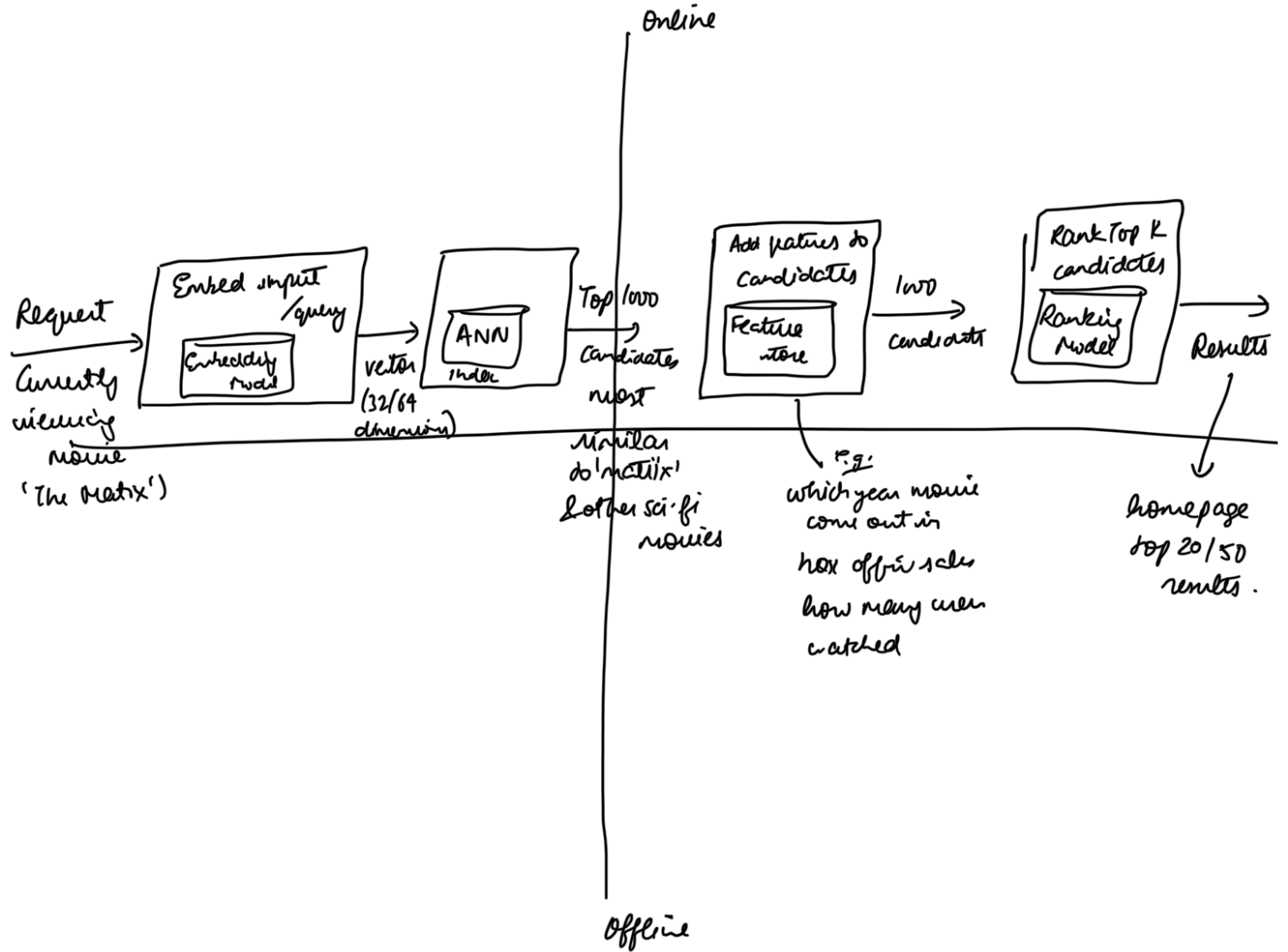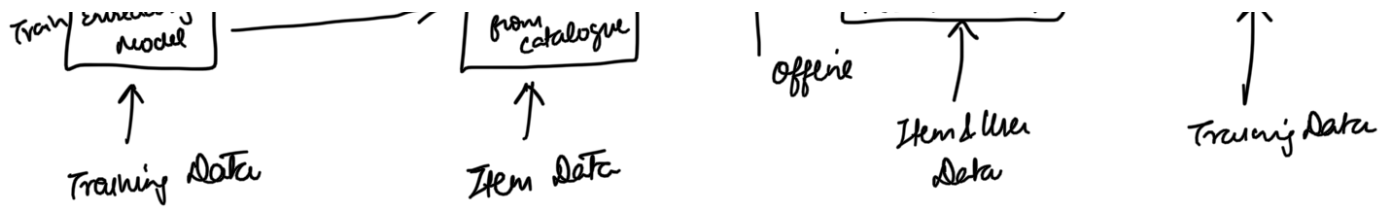
- Load data in feature
  stores.

• uses artifacts from ~~~~~~~~~
  serve requests

• — CANDIDATE RETREIVAL

  — RANKING.

  — slow but precise
  — ranks hundreds of candidates
  — adds more features
        (users, items,
         context (day/month,
         , more engagement
         with user)

Approximate
— using Nearest Neighbours (NN),
        graphs, etc

— searches millions of items to get
     hundreds of req.

— fast but coarse (not precise)



Online                              Artifacts

                    Feature Store          Ranking Model

Embedding Model      ANN Index

Retrieval ↑        ↑                  ↑                          Ranking

                  Approx. (Build)
                  NN Index

                      ↑

                              Build Feature      →    Train Ranking
                              Store                    Model
   Le ....          Embed items    (item, user, etc)

Train embedding
model

Training Data

from
catalogue

Item Data

offline

Item & User
Data

Training Data

Online

Request

Currently
viewing
movie
'The Matrix'

Embed input
/query

Embedding
model

vector
(32/64
dimension)

ANN
index

Top 1000

Candidates
most

similar
do 'matrix'
& other sci-fi
movies

Add features to
Candidates

Feature
store

e.g.
which year movie
come out in

box office sales

how many user
watched

1000
candidates

Rank Top K
candidates

Ranking
Model

Results

homepage
top 20/50
results.

Offline

Industry Examples

Offline
(a week/day)

User Behaviour

(Ruer oue

Logs

Bi-directional Item – graph

Random Walk

↓ this graph connects all items together.
( to allow random walks do items )

( based on probability / edge strength of the graph )

sequence of items

Tensorflow :-
Graph Embedding Training

→ maps each item → vector

Inner Product search ( Cosine similarity)
↓
Result

Item 2 Item Similarity Map

say item $\xrightarrow{1000}$ similar items

≡ ANN Index

③ Candidate sets ( 1 item → 600 similar Item on ANN Index)

② User Trigger (S Item)

Home Page

① Request →

Prediction Platform

④ Candidate vector list →

Ranker service Platform

⑥ Recommended List

⑤ Ranked List

Client

# Building Graphs for Query Expansion & Retrieval

**Client**

**Search - Service**

Query

KFC

Response

### Search - Service box:
- ✓ Recall
  - Spell - check
  - ↓
  - Query Understanding
  - ↓
  - Query Expansion

corrects spelling as much as possible

synonymization e.g.

KFC → Chicken
↳ Kentucky Fried Burger
+
Kentucky FC

Neo 4J
( Knowledge-
graph)

Product → Synonyms,
Restaurants,
Category etc.

∴ Query Enrichment

→ Traverse the knowledge graph finds most similar items & expands the query.

Elastic Search

✓ Precision
- ↓
- Ranker
- ↓
- Decorate

Search DB

Attributes

Docs / Sheets