

Music (Spotify | Gaana (Amazon prime Music) | Apple Music | YT Music)

↳ recommend the next few songs to a
listener (Task)

↳ 100MM Customers [scale]

10MM Songs

→ latency ...

data

ML (science)

Engg



→ lots of new users & songs - ..

D1 D2
D3

→ { Mood - specific songs ↑↑↑ .. ↑

→ language / region

→ artist / band / director ...

→ 70's / 80's / 90's songs ... (teens / childhood ...)



(Q)

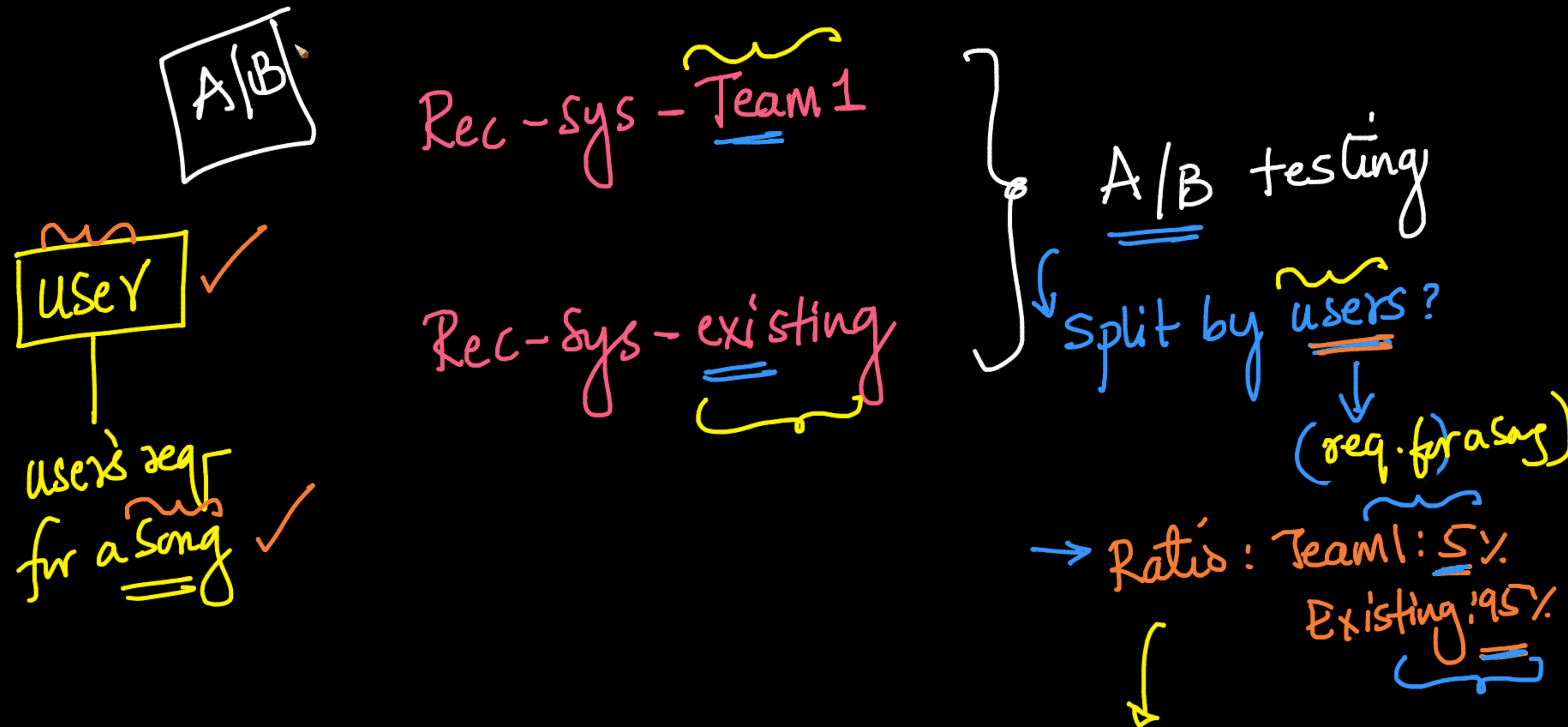
Melomics

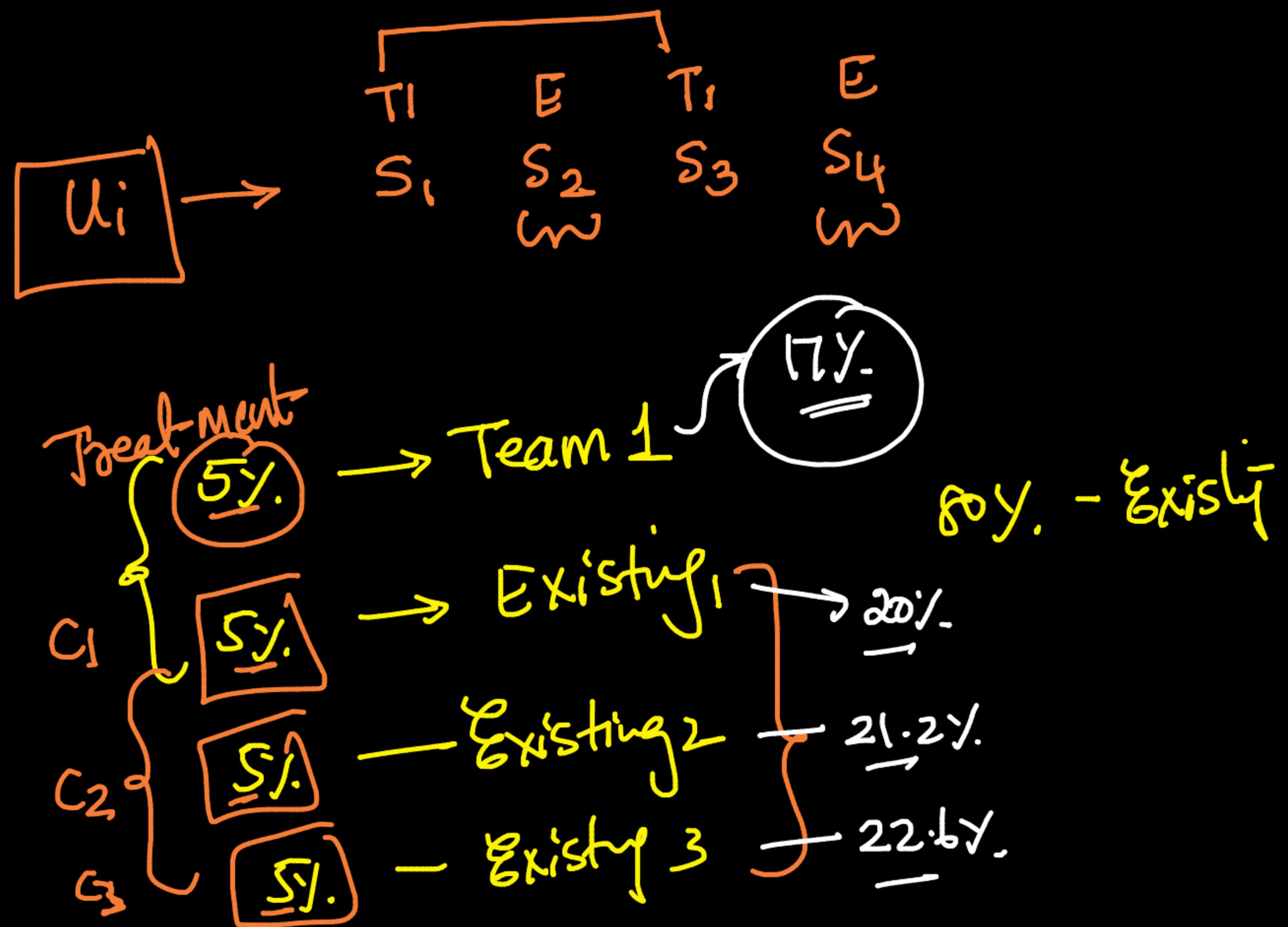
How do you measure any RecSys

(abundant)

- $\sqrt{\# \text{ skipped songs} \downarrow}$ } implicit
- $\sqrt{\text{time spent by the user} \uparrow}$
- Songs added to playlist / like → explicit
- : case

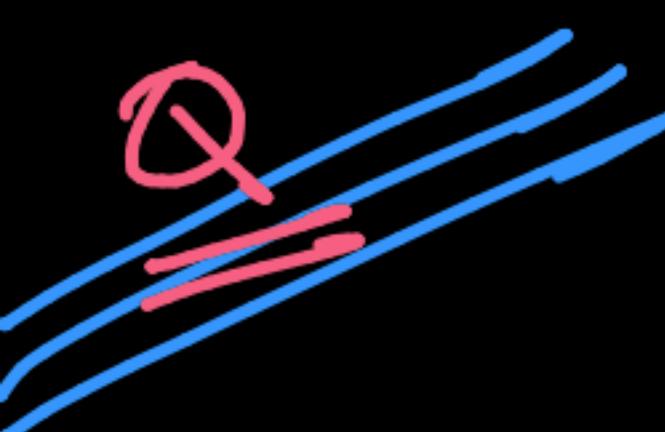
(rare)



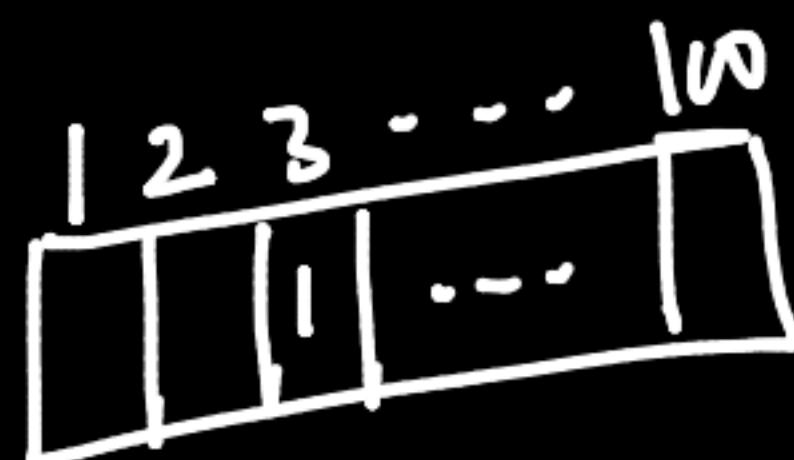


research . . .
Causal Inference
≡

A/B testing



Q ✓ { occasion → festival / ...
D: } user → features signup form, ...
age / location
User → features Title, artist, ...
Song →



{ t_k , u_i, s_j → historical data
liked / playlist (0/1)
skipped (0 - 1)

Q

{ User-features
Song-features

$$f(u, s) \rightarrow p(u_i \text{ likes } s_j)$$

ignoring t_k, u_i, s_j info

Recap:

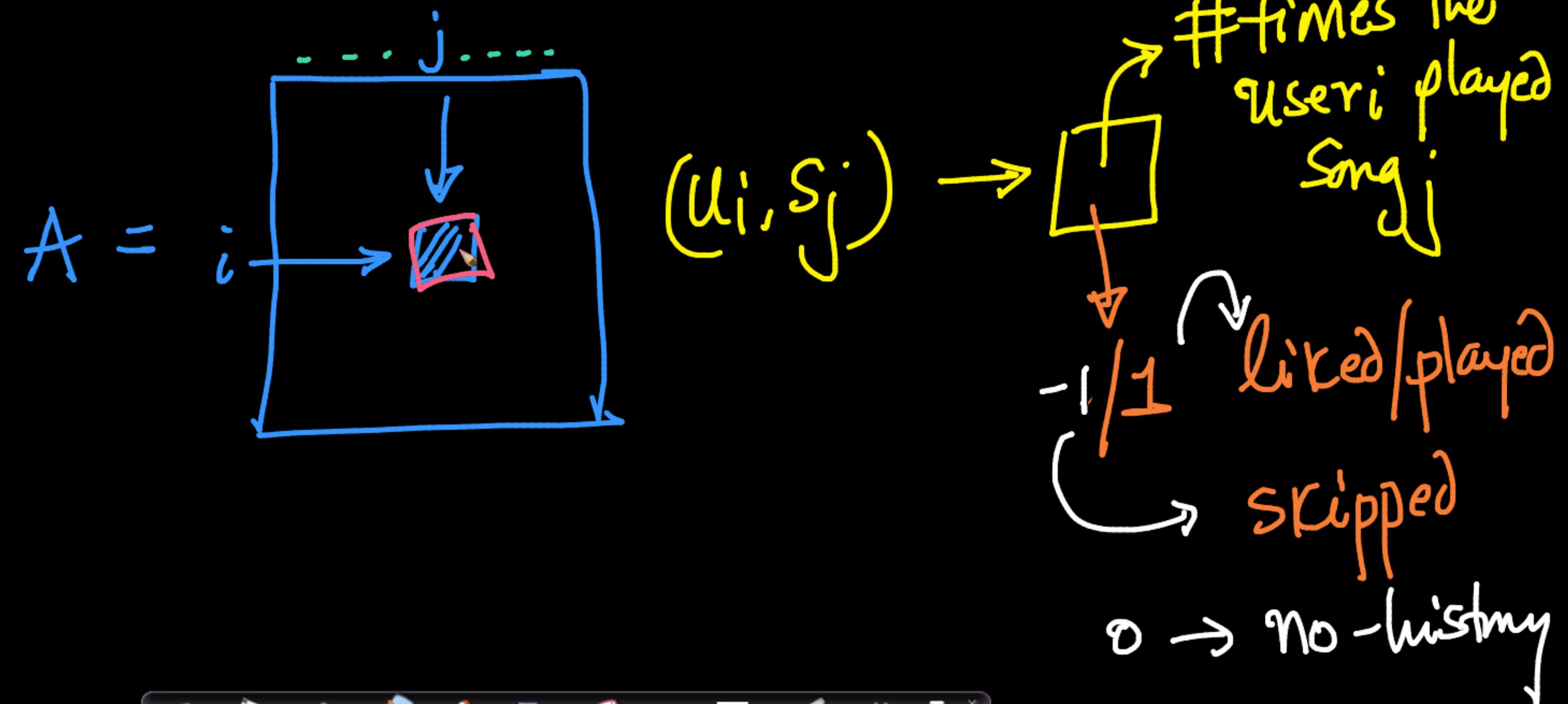
- Problem-def & challenges/nuances
- Metrics & A/B testing
- Databases

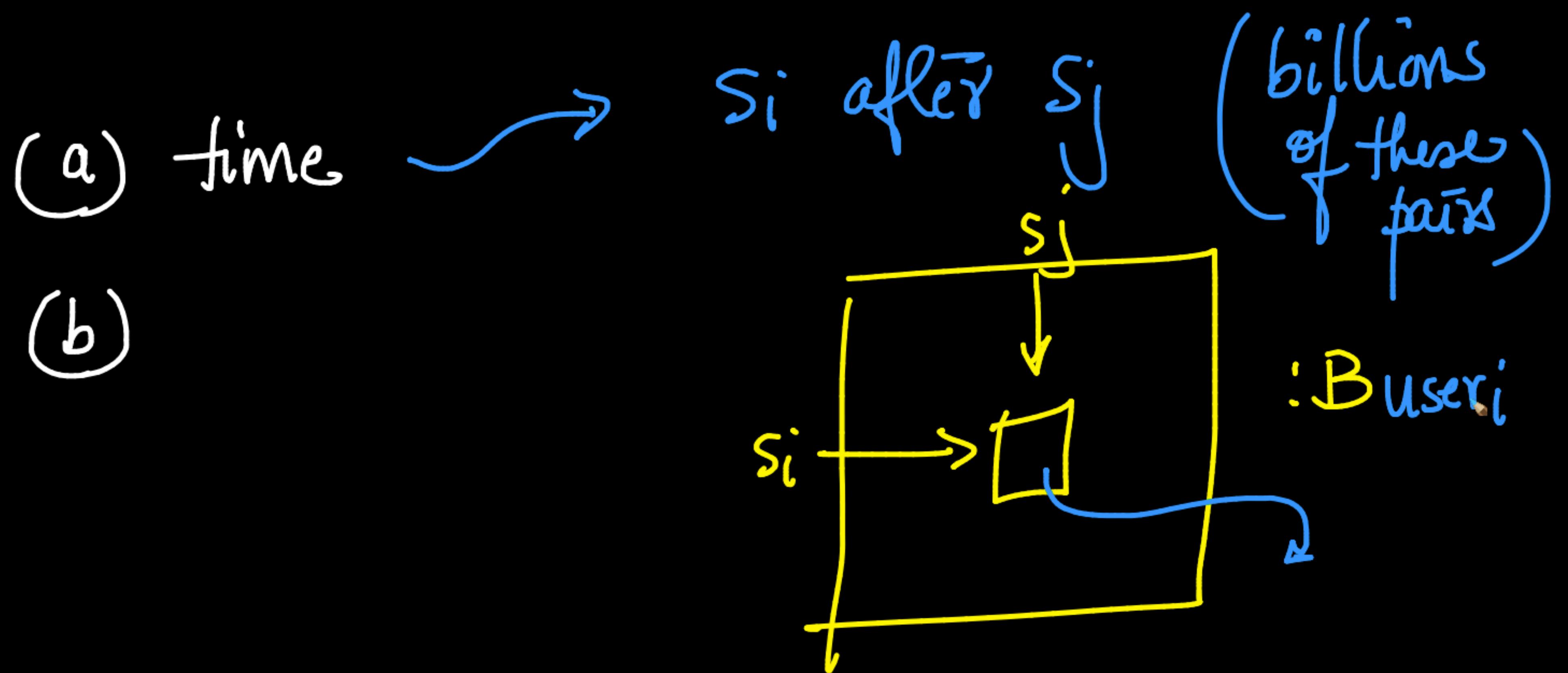
Team 1 →

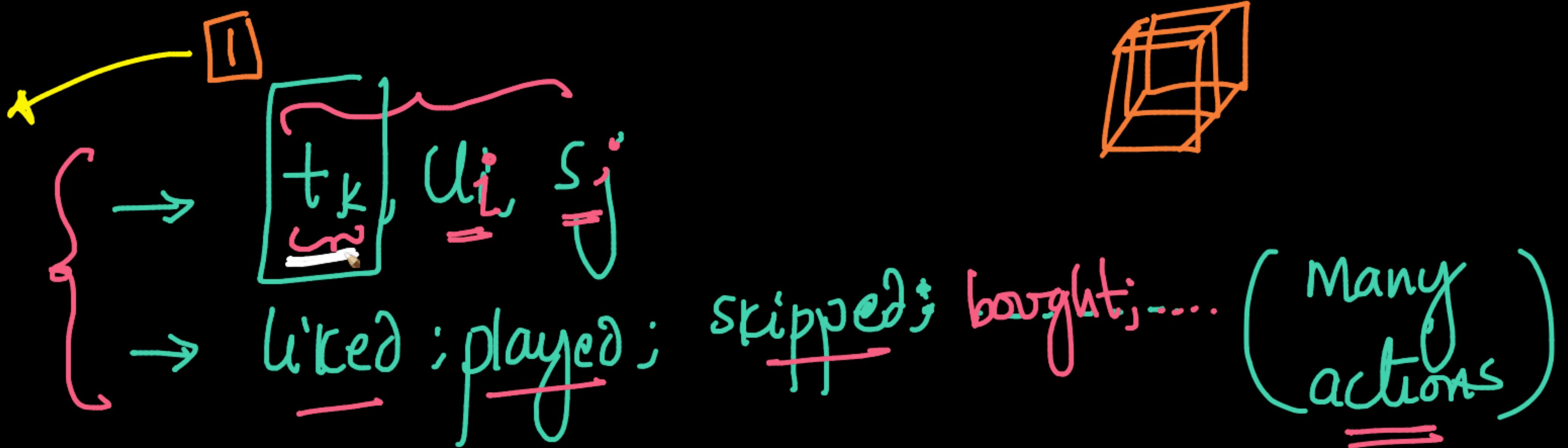
(Q)

Models

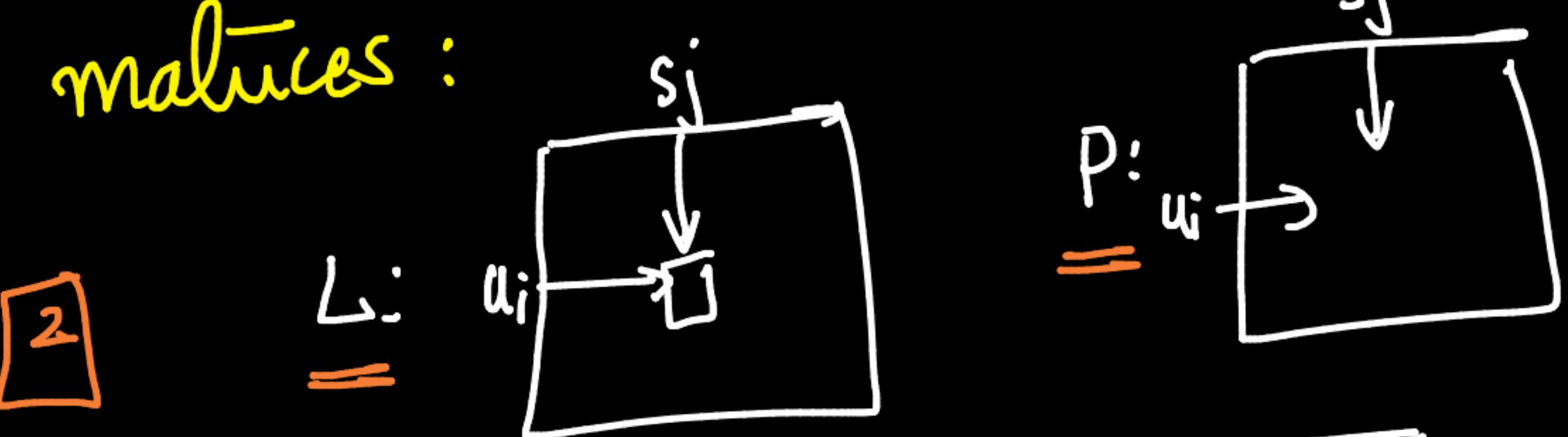
① Collaborative filtering : MF







Multiple matrices :



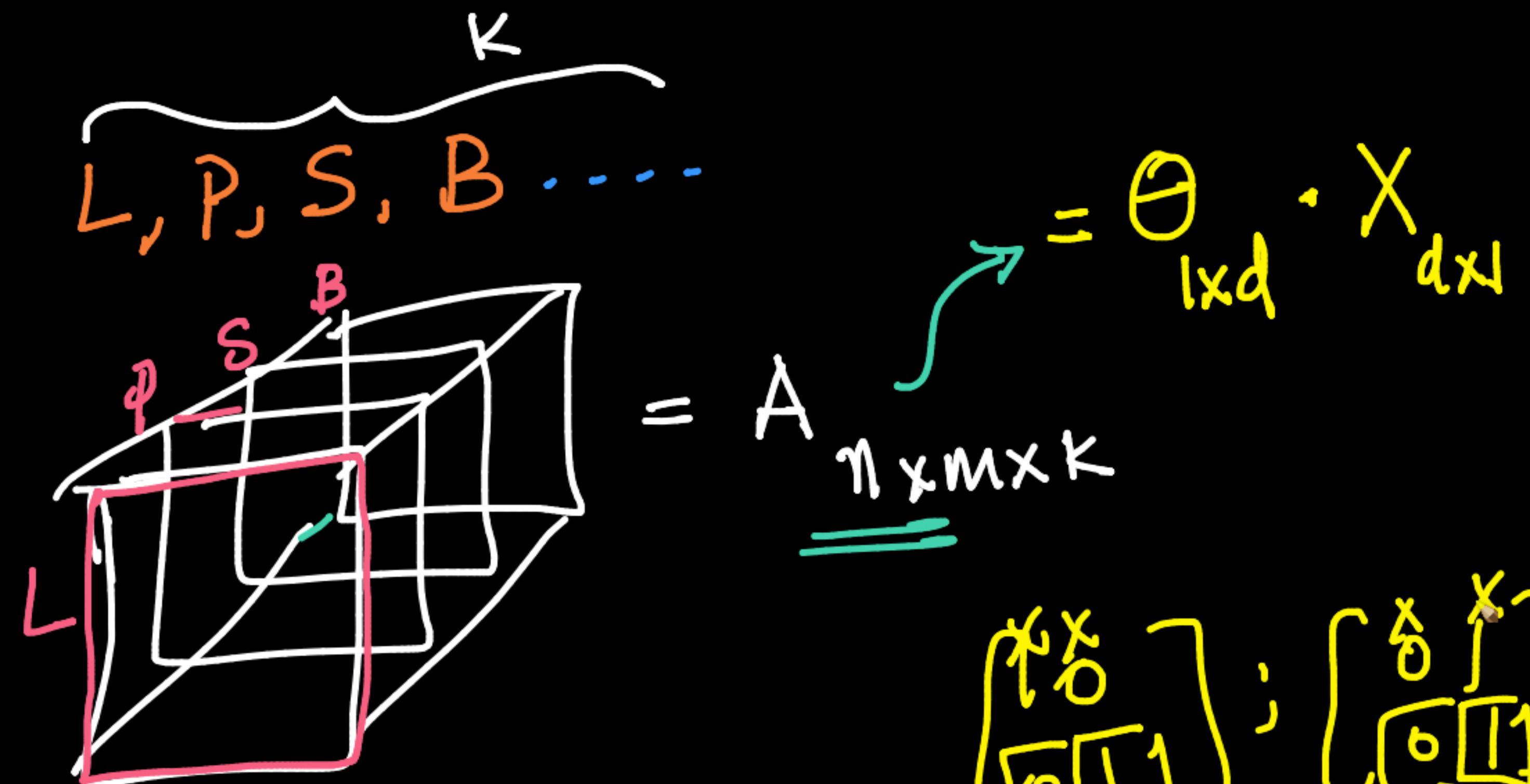
binary

$$\tilde{L} = \tilde{\Theta}_{n \times d} \cdot \tilde{X}_{d \times m}$$

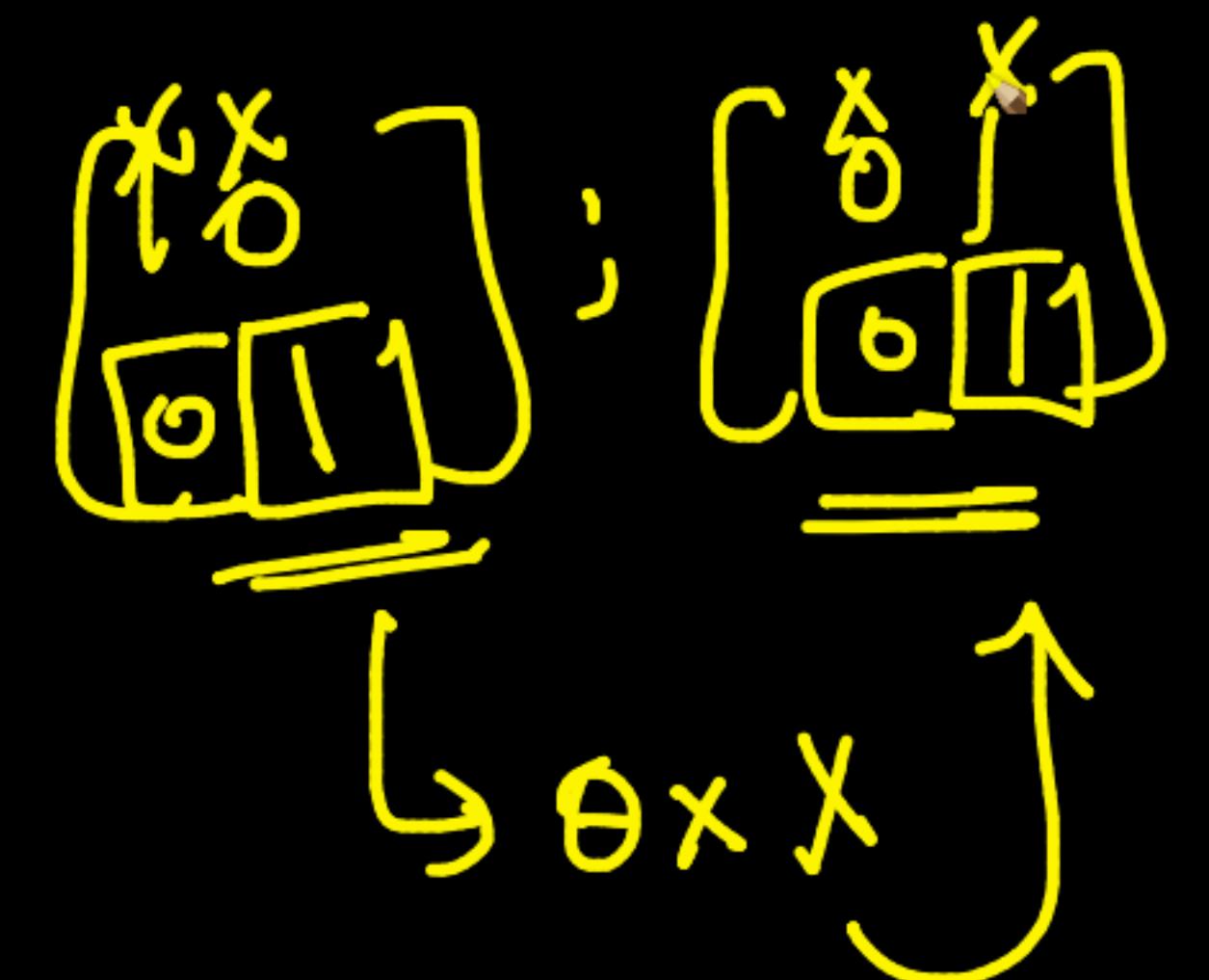
$$\tilde{P} = \tilde{\Theta}_{n \times d} \cdot \tilde{X}_{d \times m}$$

real-value
 θ

Matrices:



Typical MF: $A = B \cdot G$





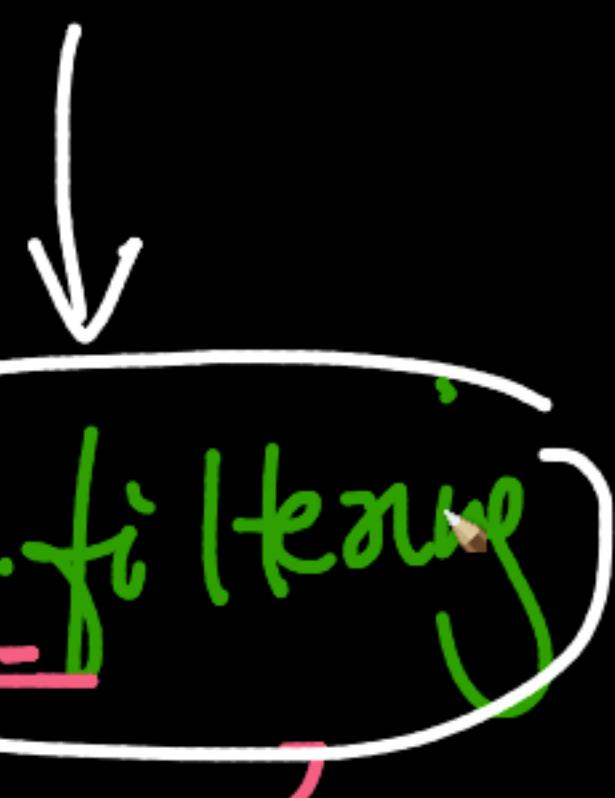
time!

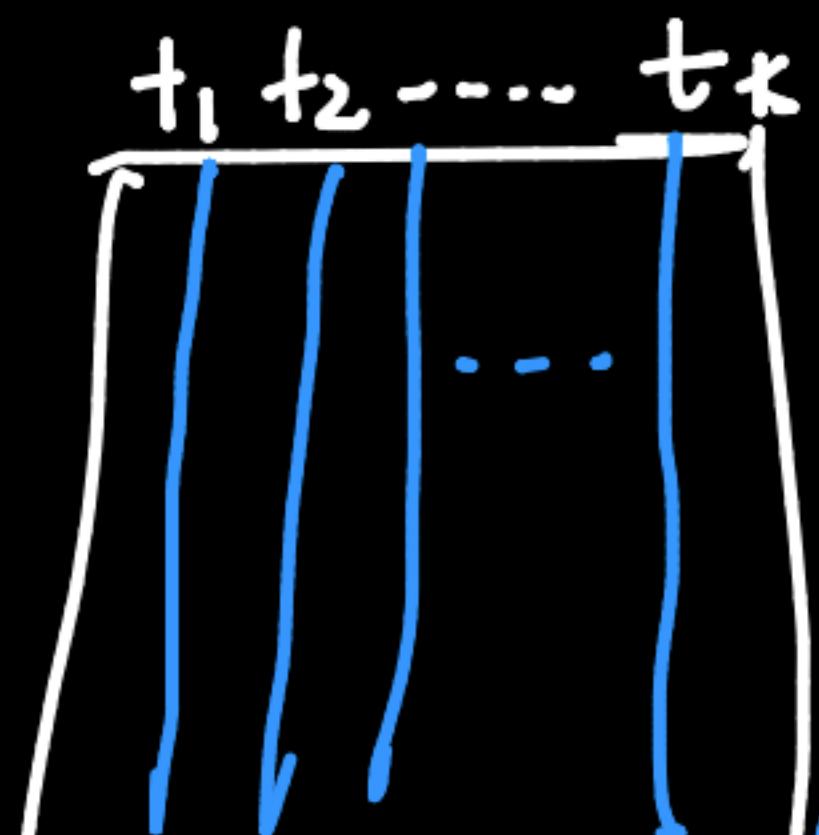
$$\lfloor \frac{n \times m}{365} \rfloor^{30}$$

Real-world

Hybrid →

Content-based + coll-filtering

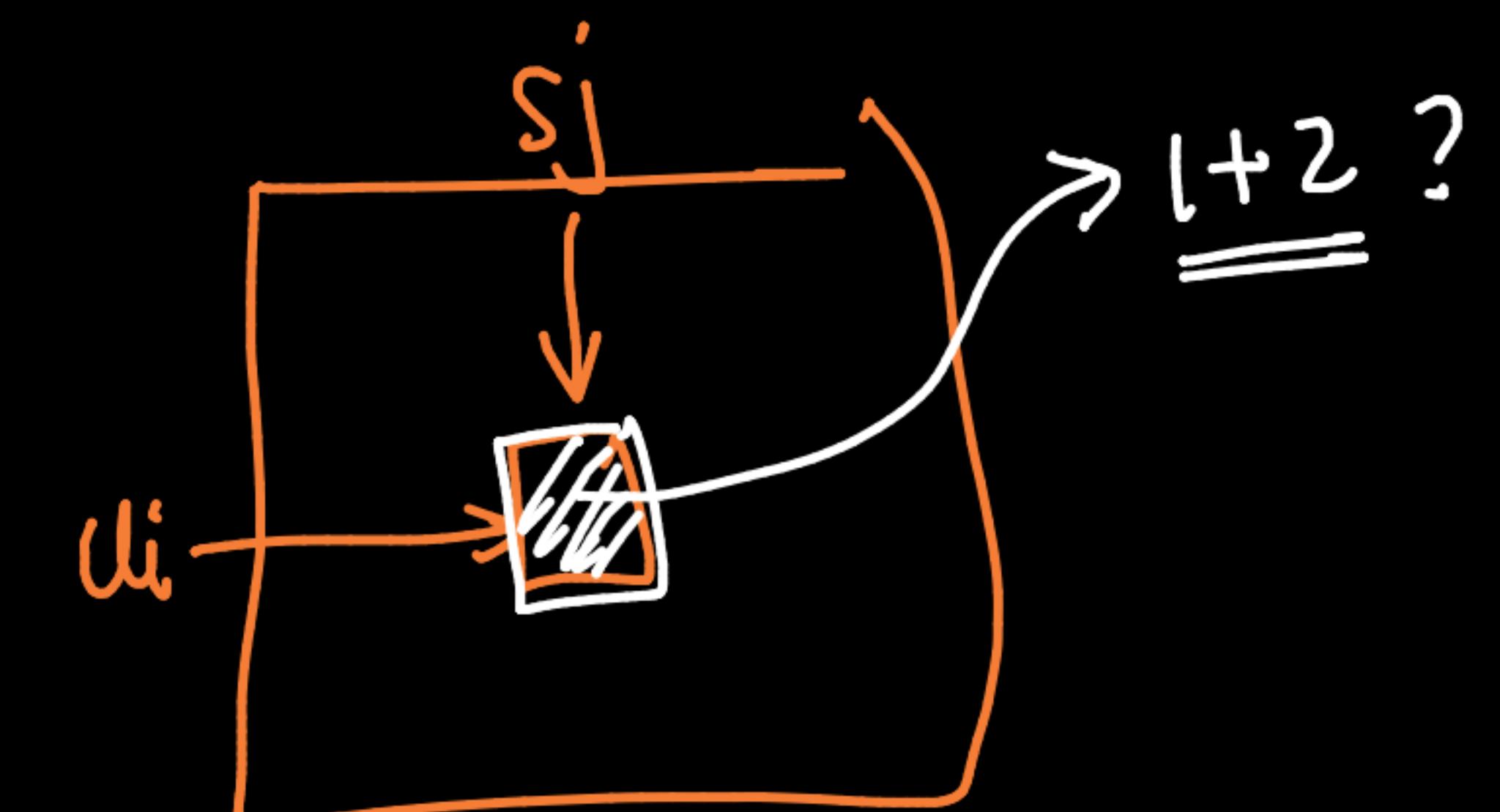


2Di time $L: n \times m \times k$  t_1, t_2, \dots, t_k  $(n \times m) \times k$

② $\tilde{L}, \tilde{R}, \tilde{S}, \dots$

Score: 1 2 0 (random)

0 0 1



$$n \times m = N$$

$$A_{ij} = \tilde{B}_i \cdot \tilde{C}_j$$

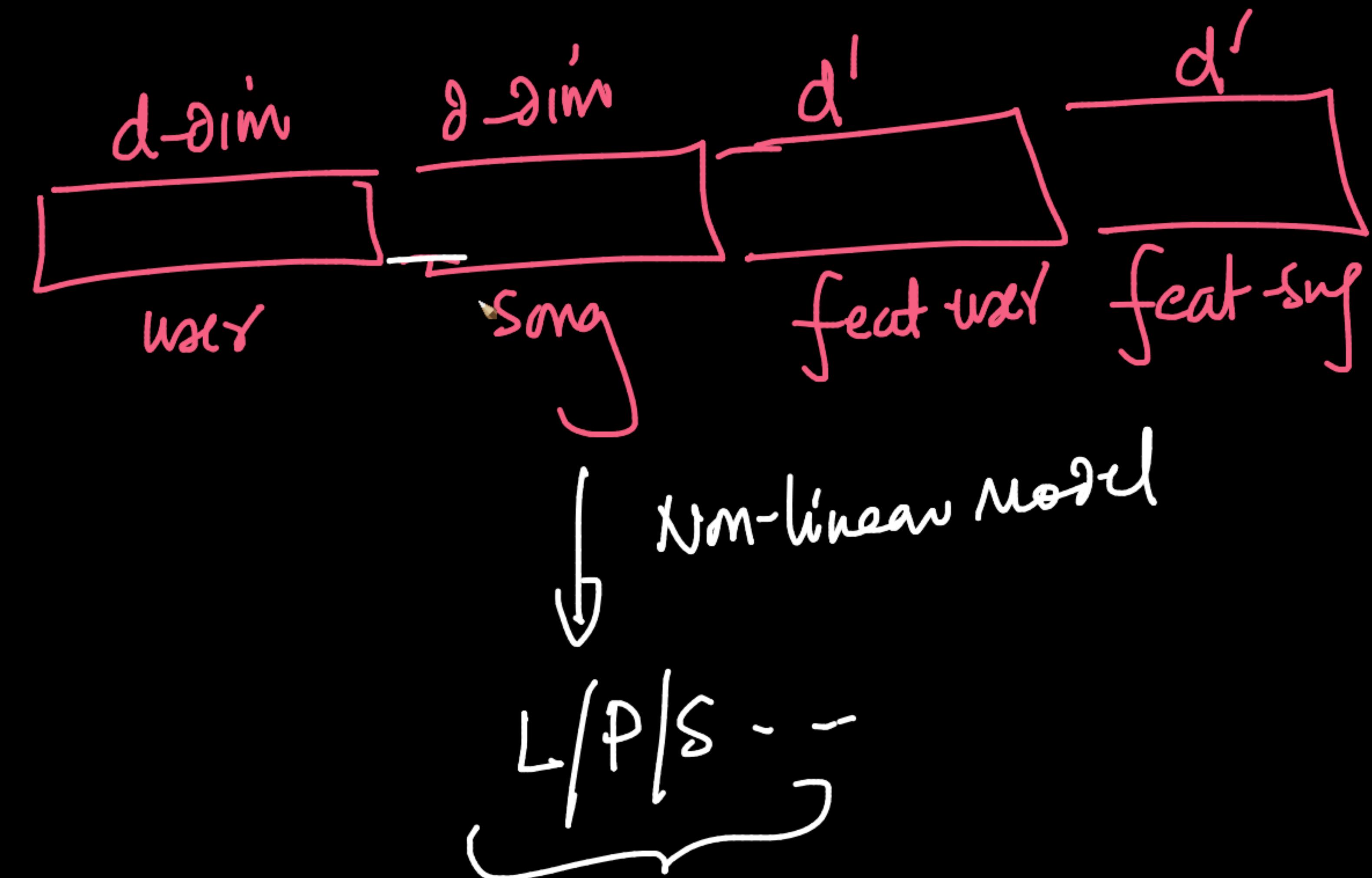
$$A = \underbrace{\begin{matrix} B & G \\ n \times d & d \times m \end{matrix}}_{\text{each row: } \text{use } \checkmark} \quad \text{each col. song}$$

each row: \checkmark

$$\tilde{L}_{N \times K} = \underbrace{\tilde{B}_{N \times d} \cdot \tilde{C}_{d \times K}}_{\downarrow}$$

each row / col represent?

Variations across time per (u_i, s_j)



"Recency" → ↓↓↓... ↓
= ↗

1 time



$L_{n \times m}$

$P_{n \times m}$

L_{ij} : Likes by u_i
on s_j in the
last 7-days

② L, P, \dots

{ ignoring all of the history }

$$L_{n \times m \times k} = \tilde{B} \cdot \tilde{C}$$



$$k = 7 / 365 / \underline{\underline{3650}}$$

[HP
Tune this]

→ One-day: · Variance } ✓

- 100MM users
- 10M Songs
=====

next-session → continuation

- Col-filtering & MF : time & multiple-signals (L, P, S ...)
 - i-i & u-u similarity ...
- scale of MF (training)
 - Interpretability
- recency & historical path
 - retention
- cold start → initial ...
 - prod volume (latency; cost)

~~Recap:~~

① Defined the problem; scale; key aspects

② Metrics; A/B-testing (briefly)

Data: $t_k, u_i, s_j \rightarrow L|P|S| \dots$ (actions)

③ MF & Col-filtering → time
many actions

④ MF & Col-filtering → many actions

Tensor fact X
Tall & skinny Matrix X
Scoring X

Good-idea

$$L_{n \times m} = A_{n \times d}^{(1)} B_{d \times m}^{(1)}$$

$$\begin{matrix} u_i^{(1)} \\ s_j \end{matrix}$$

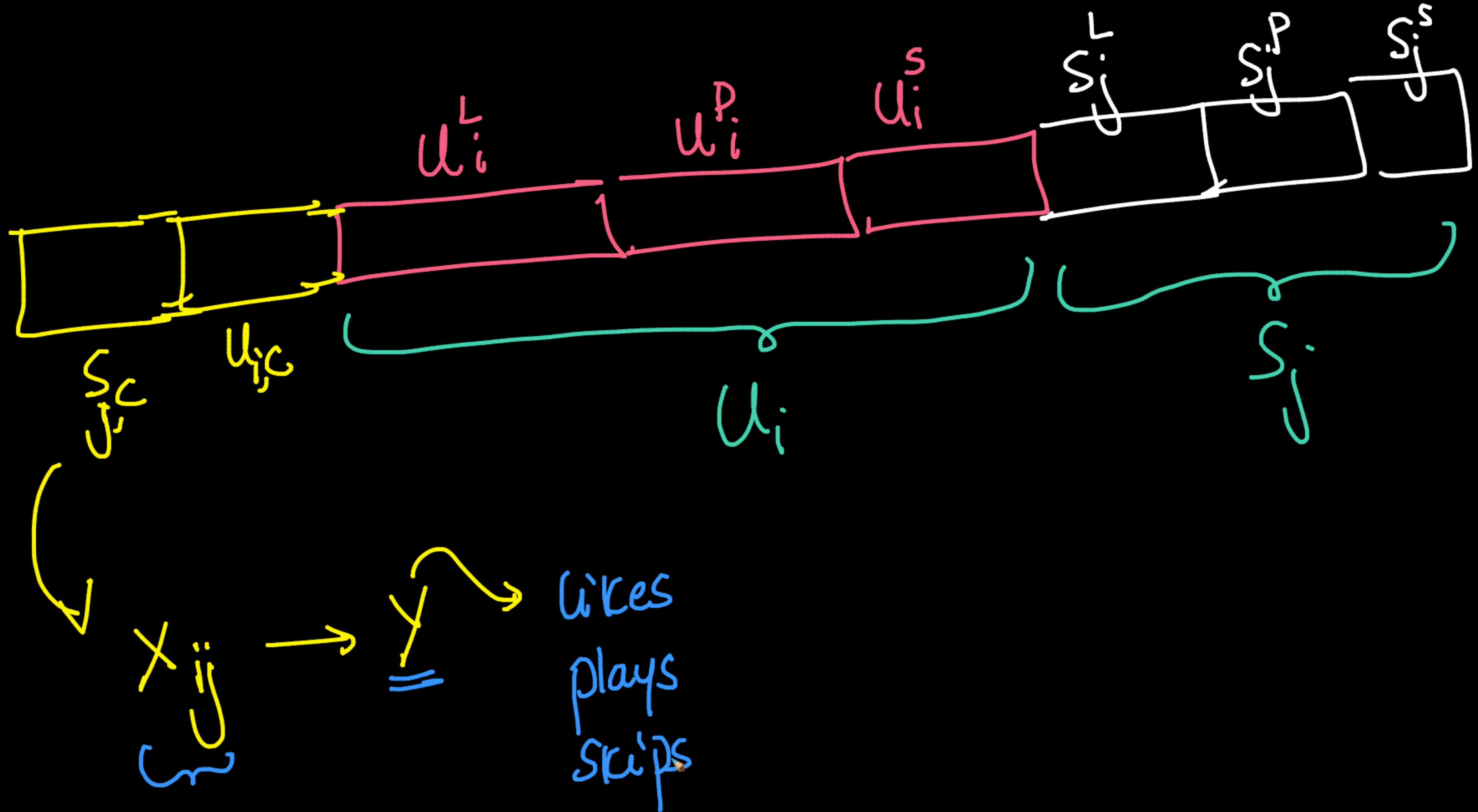
$$P_{n \times m} = A_{n \times d}^{(2)} B_{d \times m}^{(2)}$$

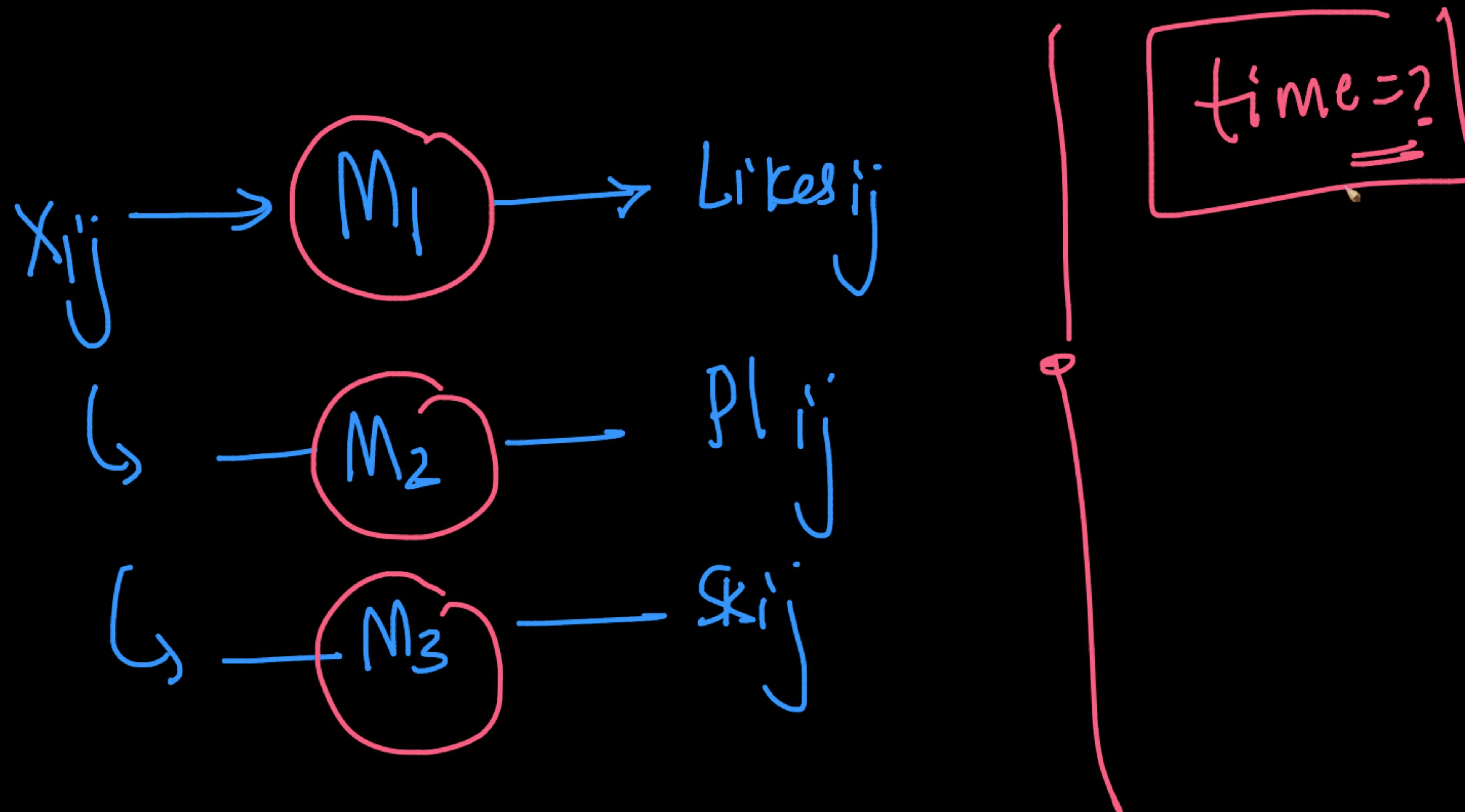
$$\begin{matrix} u_i^{(2)} \\ s_j \end{matrix} =$$

$$S_{n \times m} = A_{n \times d}^{(3)} B_{d \times m}^{(3)}$$

$$\begin{matrix} u_i^{(3)} \\ s_j \end{matrix} =$$

higher weightage to L





Time

$$\{ \text{L} \xrightarrow{\text{---}} \text{n} \times \text{m} \times \text{28} \xrightarrow{\text{---}} 7 \text{days} \times 4 (6-\text{hrs})$$

→ historical task (weighted by time)

Factorize this?

Learning with not Enough Data Part 1: Semi-Supervised Learning

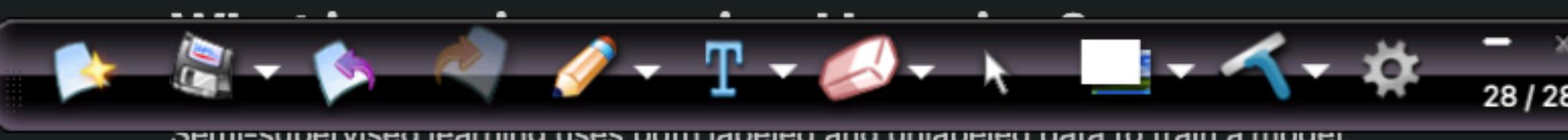
December 5, 2021 · 26 min · Lilian Weng

▶ Table of Contents

When facing a limited amount of labeled data for supervised learning tasks, four approaches are commonly discussed.

1. **Pre-training + fine-tuning:** Pre-train a powerful task-agnostic model on a large unsupervised data corpus, e.g. pre-training LMs on free text, or pre-training vision models on unlabelled images via self-supervised learning, and then fine-tune it on the downstream task with a small set of labeled samples.
2. **Semi-supervised learning:** Learn from the labelled and unlabeled samples together. A lot of research has happened on vision tasks within this approach.
3. **Active learning:** Labeling is expensive, but we still want to collect more given a cost budget. Active learning learns to select most valuable unlabeled samples to be collected next and helps us act smartly with a limited budget.
4. **Pre-training + dataset auto-generation:** Given a capable pre-trained model, we can utilize it to auto-generate a lot more labeled samples. This has been especially popular within the language domain driven by the success of few-shot learning.

I plan to write a series of posts on the topic of "Learning with not enough data". Part 1 is on *Semi-Supervised Learning*.



Interestingly most existing literature on semi-supervised learning focuses on vision tasks. And instead pre-training + fine-tuning is a more common paradigm for language tasks.