

# Recommender Systems in the Era of Large Language Models (LLMs)



Yujuan Ding<sup>1</sup>, Shijie Wang<sup>1</sup>, Liangbo Ning<sup>1</sup>, Qiaoyu Tan<sup>2</sup>, Wenqi Fan<sup>1</sup>, Qing Li<sup>1</sup>



Survey

<sup>1</sup>The Hong Kong Polytechnic University

<sup>2</sup>New York University (Shanghai)

August 4th (Day 2), Afternoon 1 & 2

IJCAI 2024, Jeju, Korea

Zoom ID: 864 7573 0054, Password: 732469



Website

# Tutorial Outline

- ◎ **Part 1: Introduction of RecSys in the era of LLMs (Dr. Wenqi Fan)**
- **Part 2: Preliminaries** of RecSys and LLMs (Dr. Yujuan Ding)
- **Part 3: Pre-training** paradigms for adopting LLMs to RecSys (Dr. Yujuan Ding)
- **Part 4: Fine-tuning** paradigms for adopting LLMs to RecSys (Liangbo Ning)
- **Part 5: Prompting** paradigms for adopting LLMs to RecSys (Shijie Wang)
- **Part 5: Future directions** of LLM-empowered RecSys (Dr. Wenqi Fan)

Website of this tutorial  
Check out the slides and more information!



# Recommender Systems (RecSys)



## Age of Information Explosion



amazon

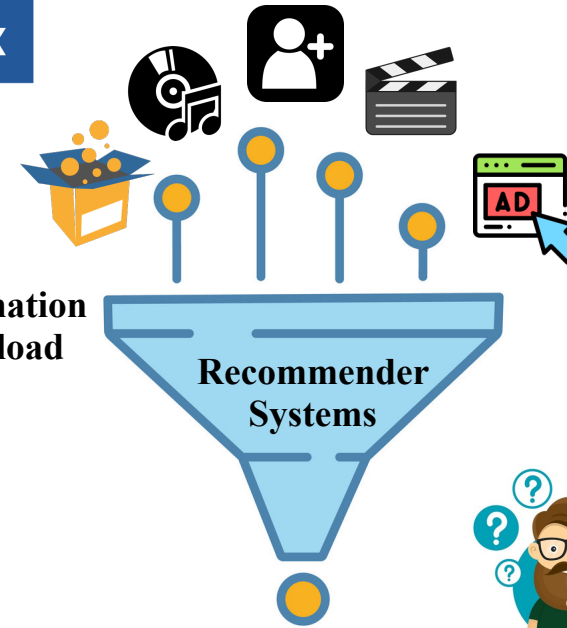


LinkedIn

facebook

淘宝网  
Taobao.com

Information  
overload



Recommend item X to user

Items can be: Products, Friends, News, Movies, Videos, etc.



# Recommender Systems (RecSys)



- Recommendation has been widely applied in online services:
  - ❖ E-commerce, Content Sharing, Social Networking ...



## Product Recommendation

Frequently bought together



Total price: \$208.9

Add all three to Cart

Add all three to List



Amazon's recommendation algorithm drives **35%** of its sales [from McKinsey, 2012]



# Recommender Systems (RecSys)



- Recommendation has been widely applied in online services:
  - ❖ E-commerce, Content Sharing, Social Networking ...

YouTube

TikTok

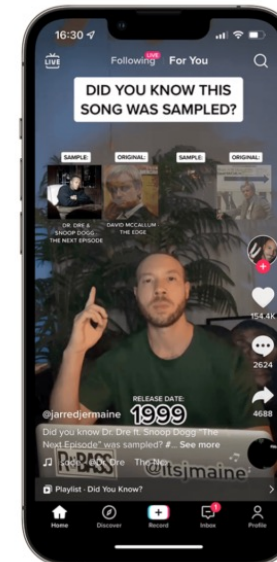
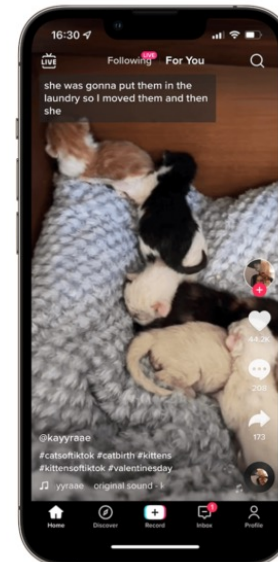
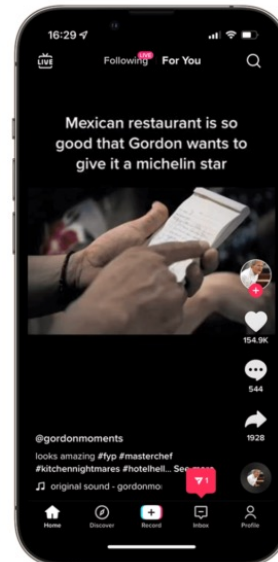


## News/Video/Image Recommendation

TikTok's recommendation algorithm

Top 10 Global Breakthrough  
Technologies in 2021

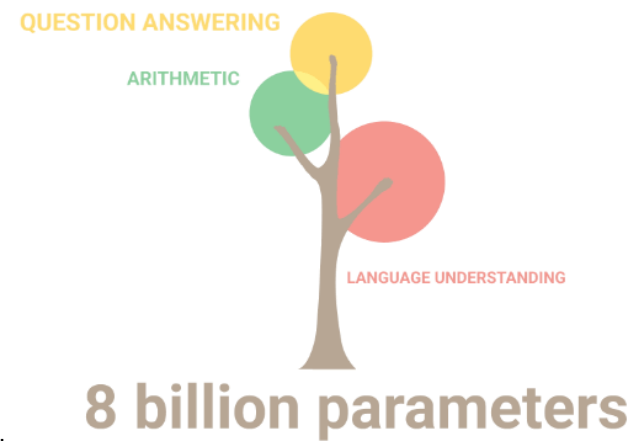
MIT  
Technology  
Review



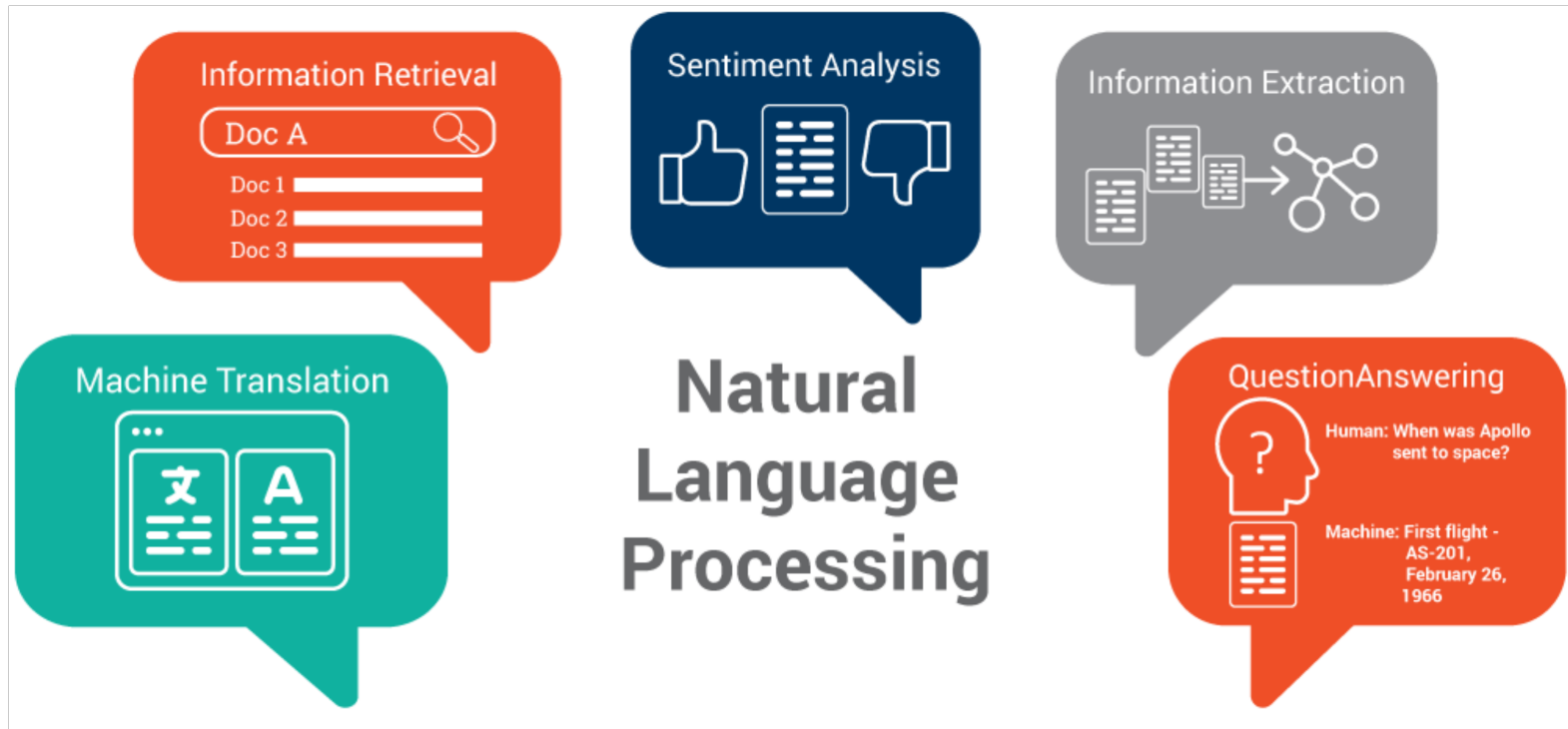
# Large Language Models (LLMs)



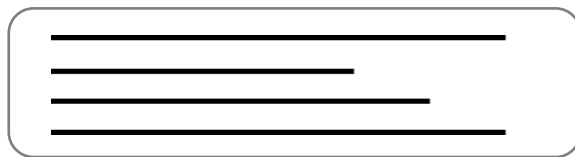
## They Are Changing Our Lives !



# LLMs in Natural Language Processing



**Input Text**



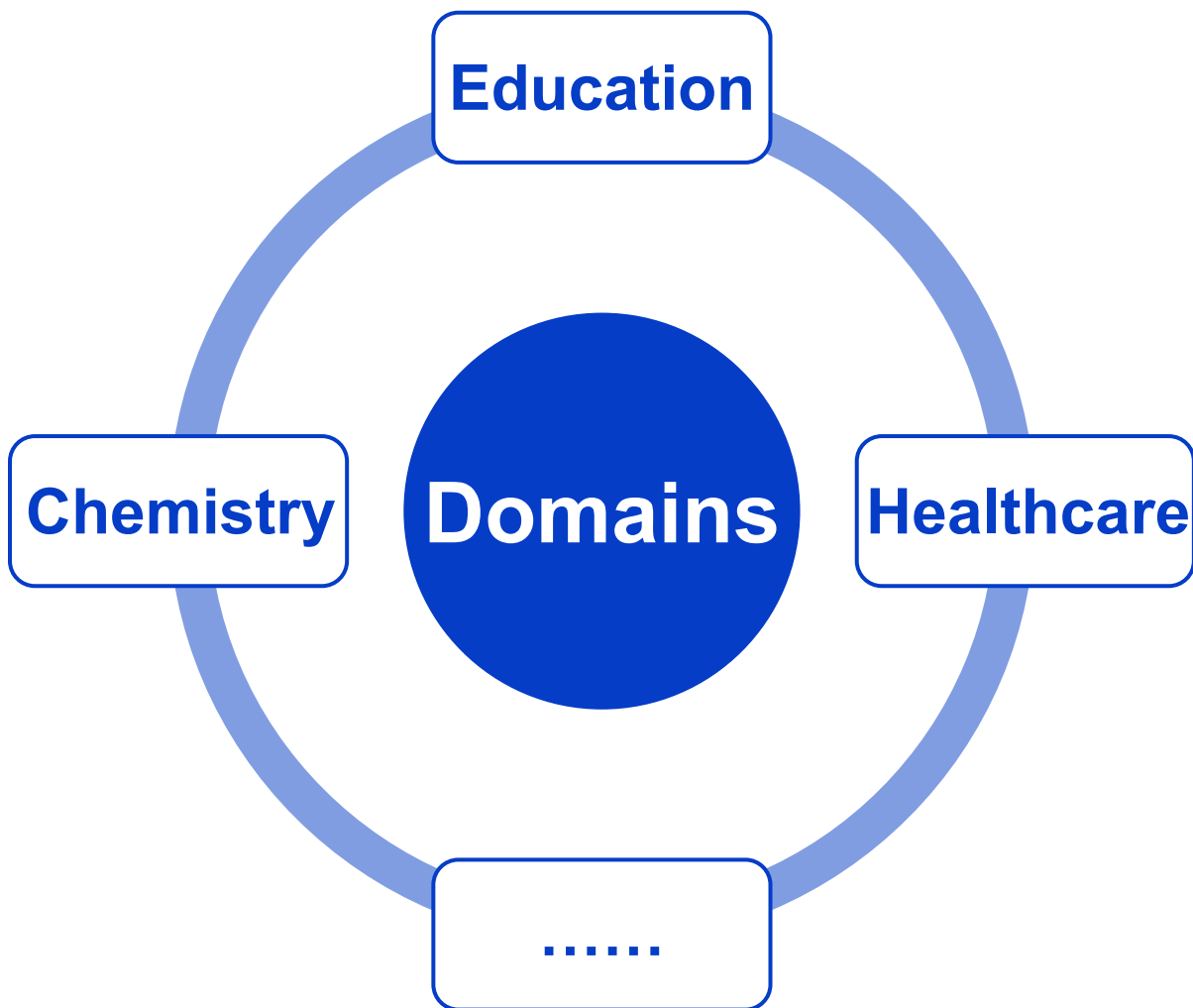
**Generated Text**



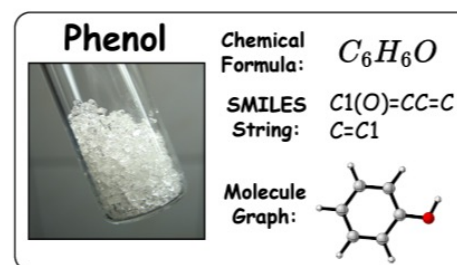
**Large Language Models (LLMs)**



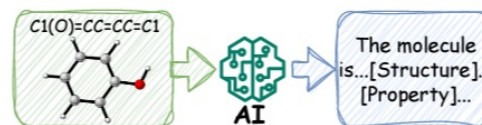
# LLMs in Downstream Domains



## □ Molecule discovery, etc.



### (a) Molecule Representations.



### (b) Molecule Captioning.



**ChatGPT**

**(a) Molecule Captioning**

Please show me a description of this molecule: "C1=CC=C(C=C1)OC2=CC=CC=C2"

The molecule is an aromatic ether in which the oxygen is attached to two phenyl substituents. It has been found in muscat grapes and vanilla. It has a role as a plant metabolite.

**(b) Text-based Molecule Generation**

Help me generate a molecule based on the given description:

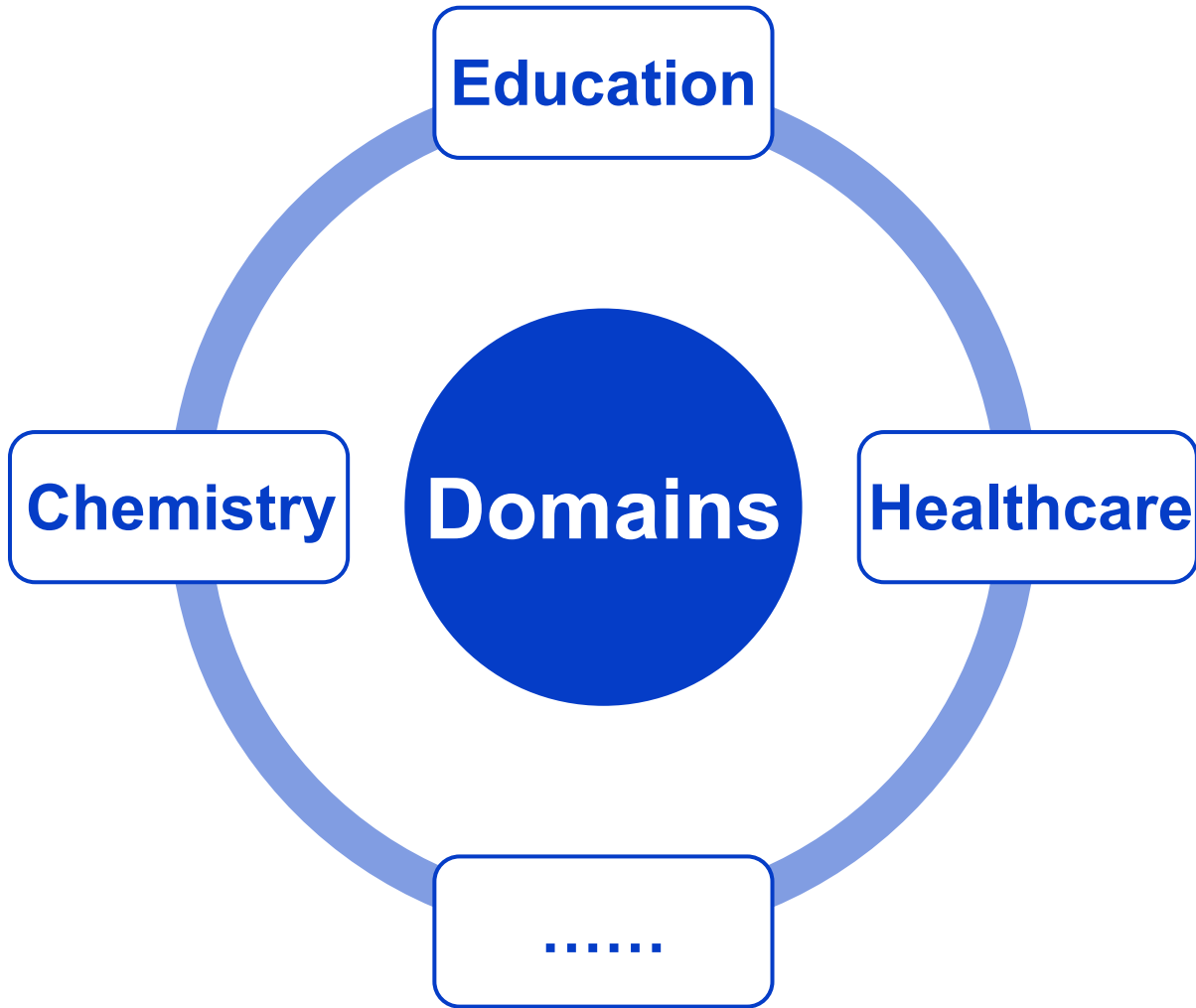
The molecule is a quinolinemonocarboxylate that is the conjugate base of xanthurenic acid, obtained by deprotonation of the carboxy group. It has a role as an animal metabolite. It is a conjugate base of a xanthurenic acid.

C1=CC2=C(C=C1)[O-]NC(=CC2=O)C(=O)O

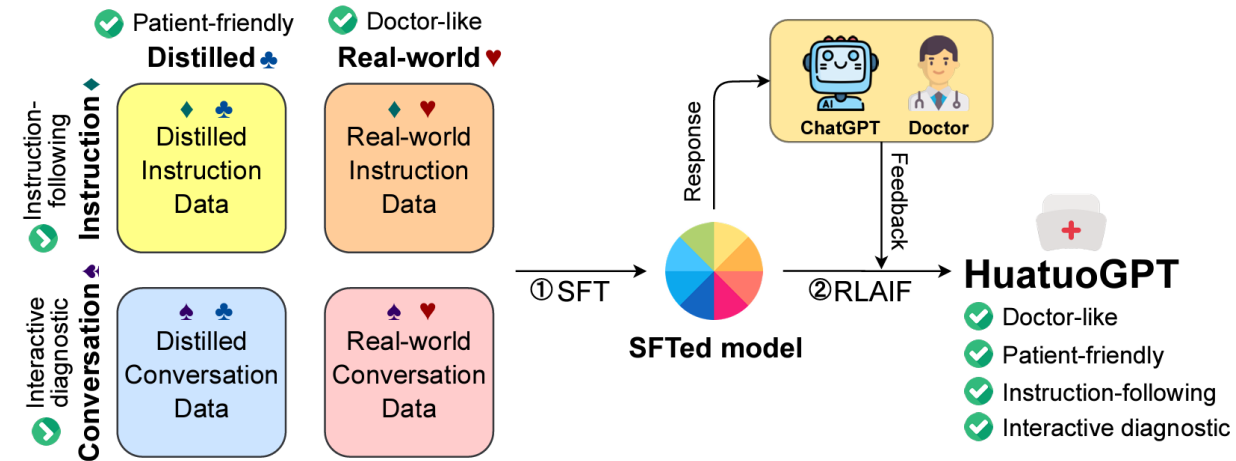




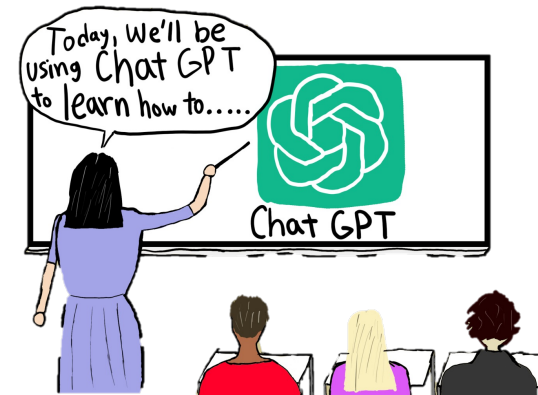
# LLMs in Downstream Domains



## ❑ Medical consultation, etc.



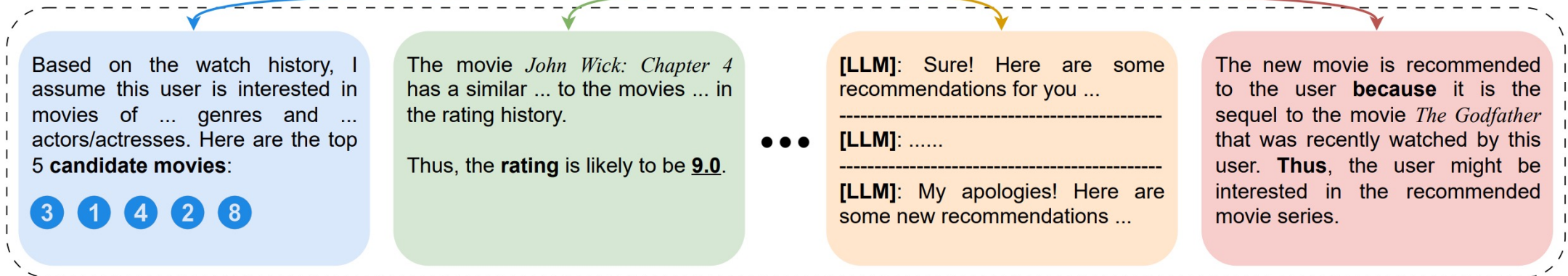
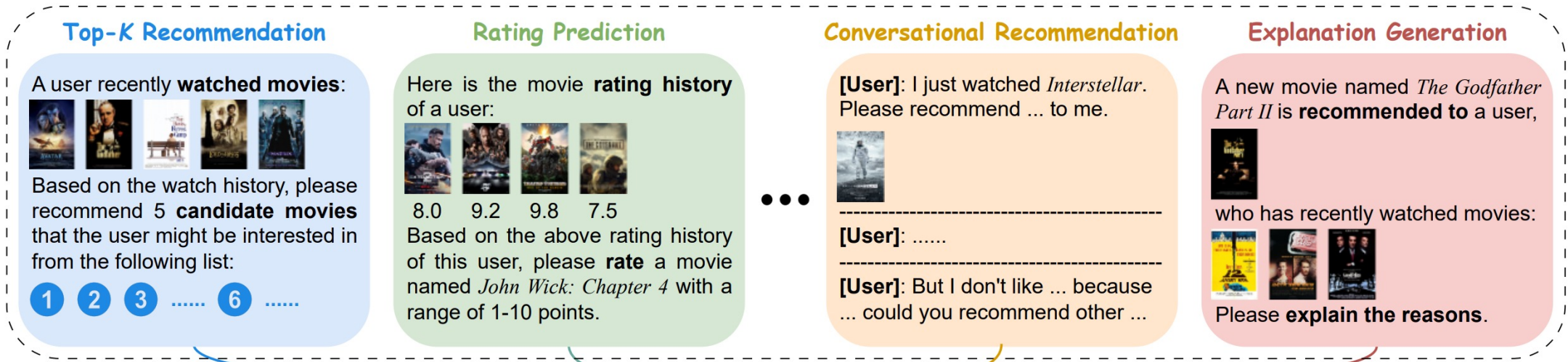
## ❑ Curriculum & Teaching, etc.



# LLMs in RecSys



## Task-specific Prompts (LLMs Inputs)



## Task-specific Recommendations (LLMs Outputs)

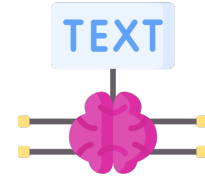


# Potentials of LLMs in RecSys



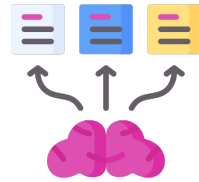
As the **parameter size** of LLMs continues to **scale up** with a larger **training corpus** ...

## ❑ Language understanding and generation ability



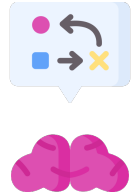
- ❖ LLMs can comprehend **human intentions** and generate **language responses** that are more human-like in nature.

## ❑ Generalization capability



- ❖ LLMs can apply their learned knowledge to **fit various downstream tasks**, even **without being fine-tuned** on specific tasks.

## ❑ Reasoning capability



- ❖ LLMs can generate the outputs with **step-by-step reasonings** to support complex **decision-making processes**.

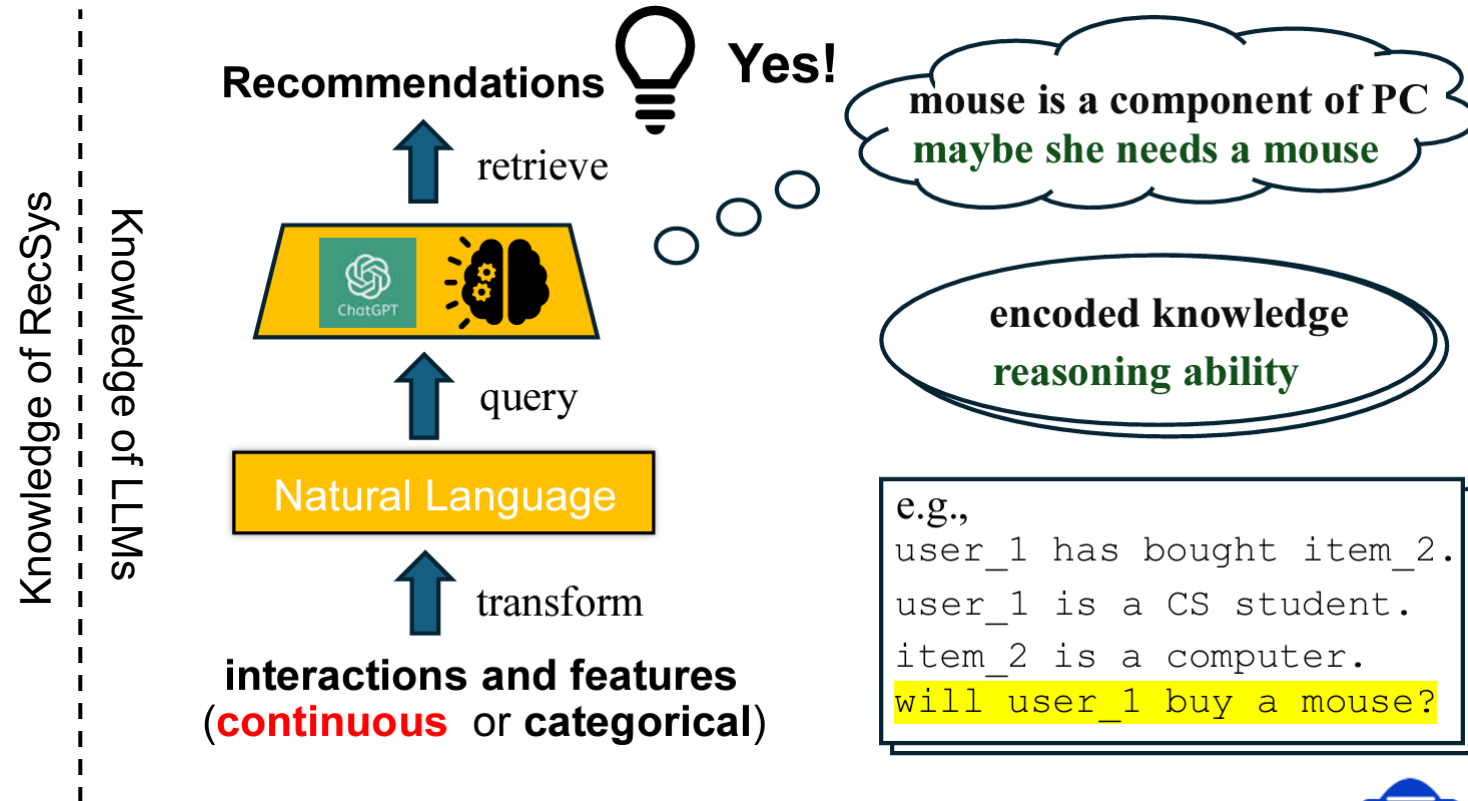


# Language Understanding & Generation



- ❑ Sufficiently capture **textual knowledge** about users and items
  - ❖ Rich **textual side information** about users and items in RecSys
  - ❖ Diverse **open-world knowledge** encoded in LLMs

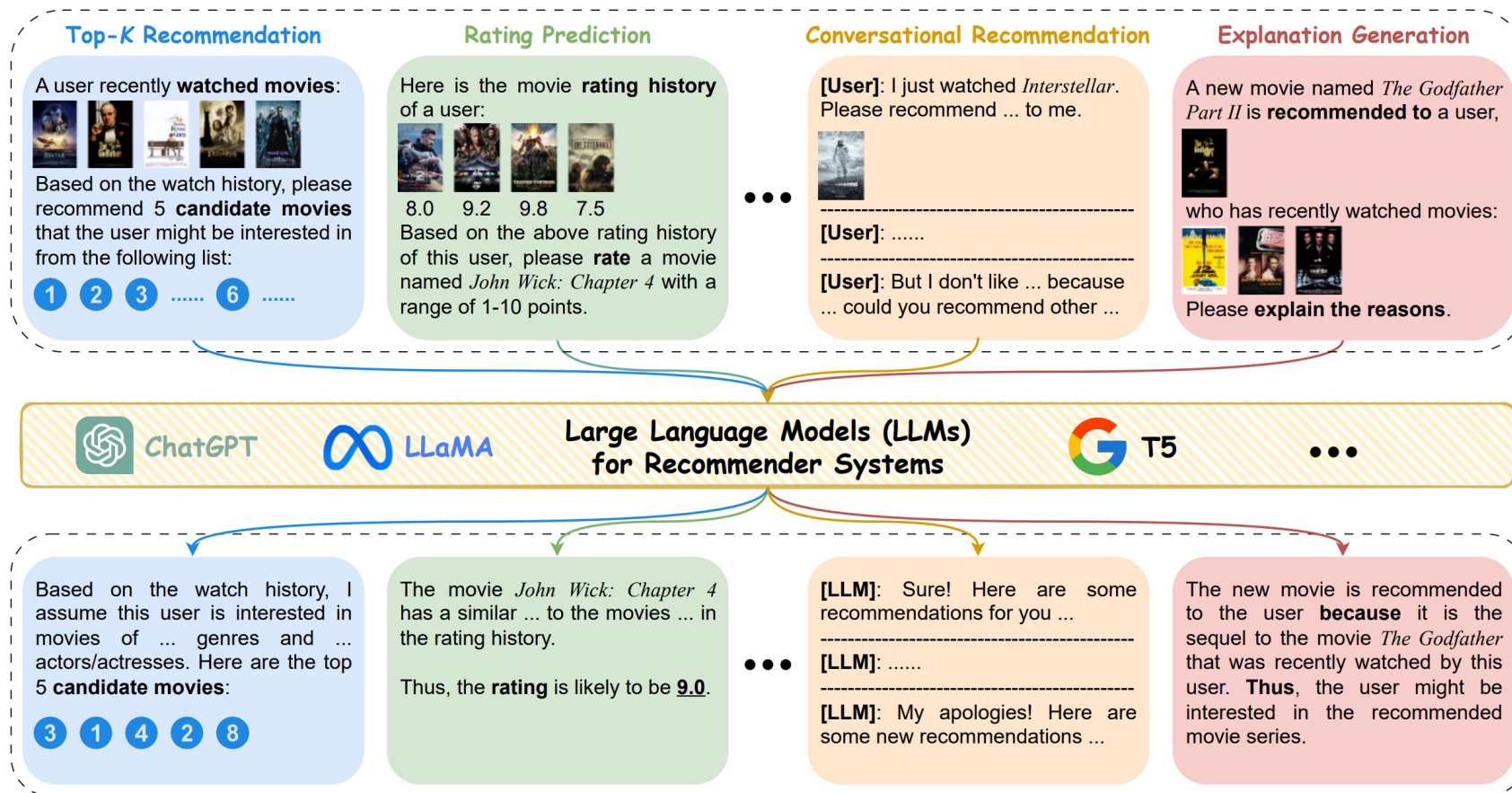
**User ID: 0057** **Item ID: 0046**  
**Item Title:** Wet n Wild Mega Last Lip Color 908C Sugar Plum Fairy  
**Review:** The color is a perfect mix of dark purple, red and pink. The only downside is the drying aspect of the lipstick, which I counteract by using lip balm before putting it on.



# Generalization



- Adapt to **various recommendation tasks** even without being fine-tuned
- ❖ LLMs can apply their **learned knowledge** to address recommendation objectives
- ❖ **Multi-task adaption** by providing appropriate task instructions or a few task demonstrations

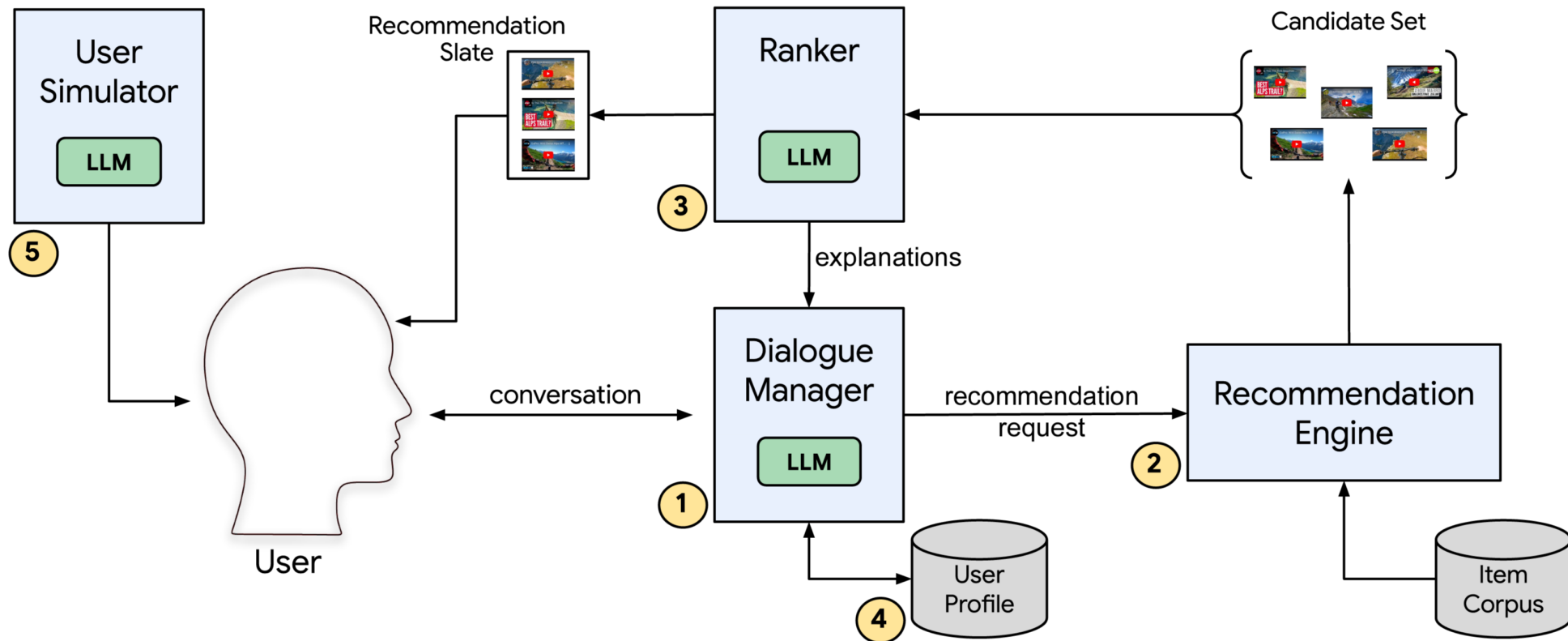


# Reasoning



□ Support complex **decision-making processes** in RecSys

- ❖ Retrieve information from **large contexts** and control **multi-step** recommendation tasks
- ❖ Generate outputs with **step-by-step reasoning** empowered by chain-of-thought prompting



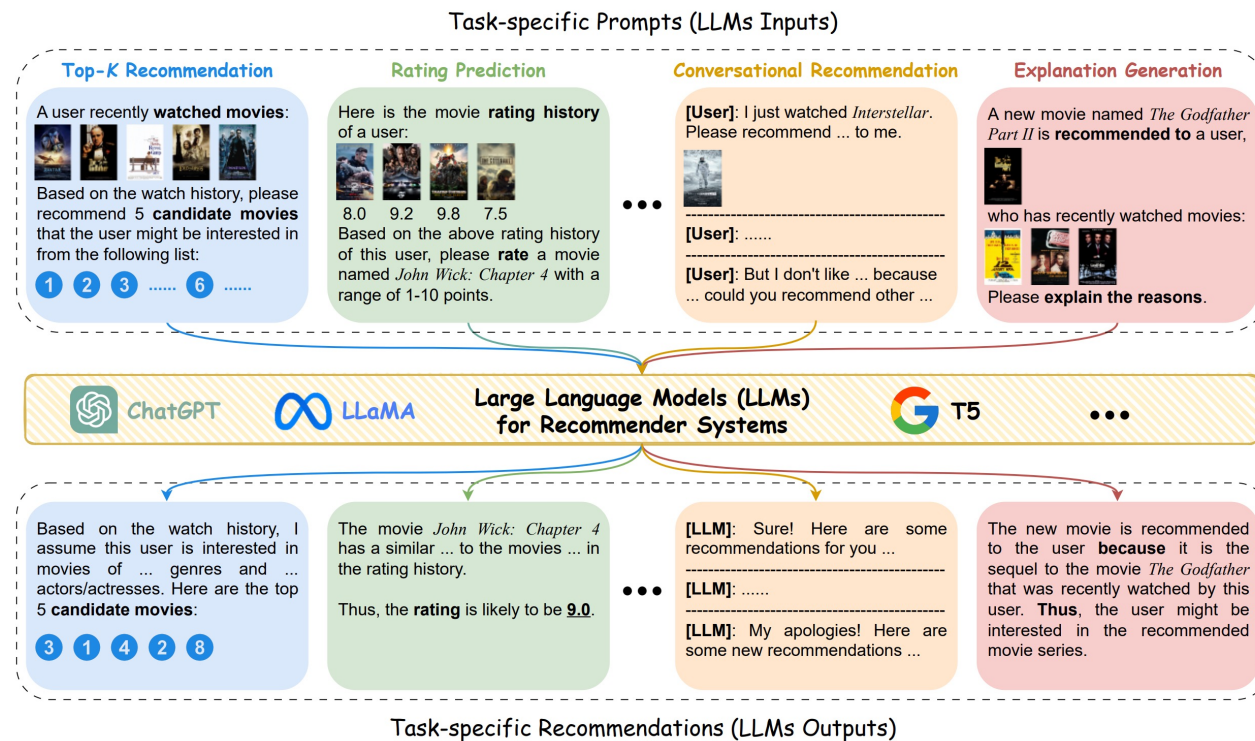
# A Comprehensive Survey Paper



## Recommender Systems in the Era of Large Language Models (LLMs)

Zihuai Zhao, Wenqi Fan, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Zhen Wen, Fei Wang, Xiangyu Zhao, Jiliang Tang, and Qing Li

<https://arxiv.org/abs/2307.02046>



# Recruitment



- ❑ Our research group are actively recruiting self-motivated **postdoc, Ph.D. students, and research assistants**, etc. **visiting scholars, interns, and self-funded students** are also welcome. Send me an email if you are interested.
  
- ❖ Research areas: machine learning (ML), data mining (DM), artificial intelligence (AI), deep learning (DNNs), large language models (LLMs), graph neural networks (GNNs), computer vision (CV), natural language processing (NLP), etc.
  
- ❖ Position details:  
<https://wenqifano3.github.io/openings.html>





# Tutorial Outline

- ⦿ **Part 1: Introduction** of RecSys in the era of LLMs (Dr. Wenqi Fan)
- ⦿ **Part 2: Preliminaries of RecSys and LLMs (Dr. Yujuan Ding)**
- **Part 3: Pre-training** paradigms for adopting LLMs to RecSys (Dr. Yujuan Ding)
- **Part 4: Fine-tuning** paradigms for adopting LLMs to RecSys (Liangbo Ning)
- **Part 5: Prompting** paradigms for adopting LLMs to RecSys (Shijie Wang)
- **Part 5: Future directions** of LLM-empowered RecSys (Dr. Wenqi Fan)

Website of this tutorial  
Check out the slides and more information!



# PART 2: Preliminaries of RecSys and LLMs



**Presenter**  
**Dr. Yujuan DING**  
**HK PolyU**

- **Recommender Systems (RecSys)**
  - Collaborative Filtering (CF)
  - Content-based Recommendation
  - Deep Recommender Systems
- **Large Language Models (LLMs)**
  - Development and Capability
  - LLM Architecture
- **LLM-based RecSys**
  - ID-based LLM RecSys
  - Text-based LLM RecSys



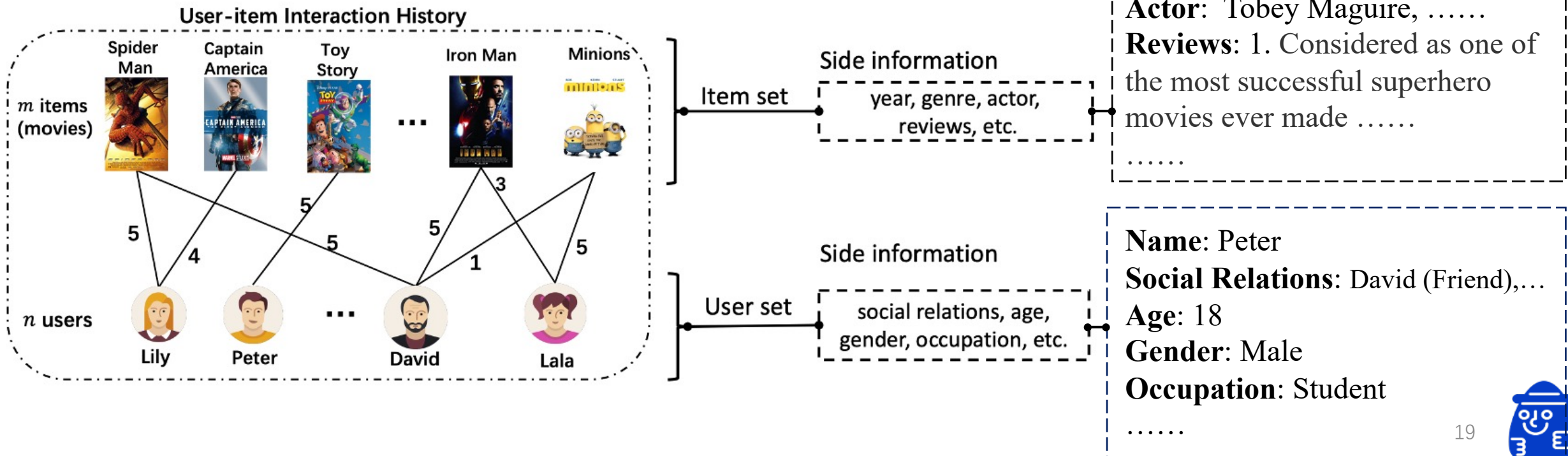
# Recommender Systems



Historical user-item interactions or additional side information (e.g., social relations, item attributes, etc)



Predictions on how likely a user would be interested or interact (click, view, purchase, etc) with a target item

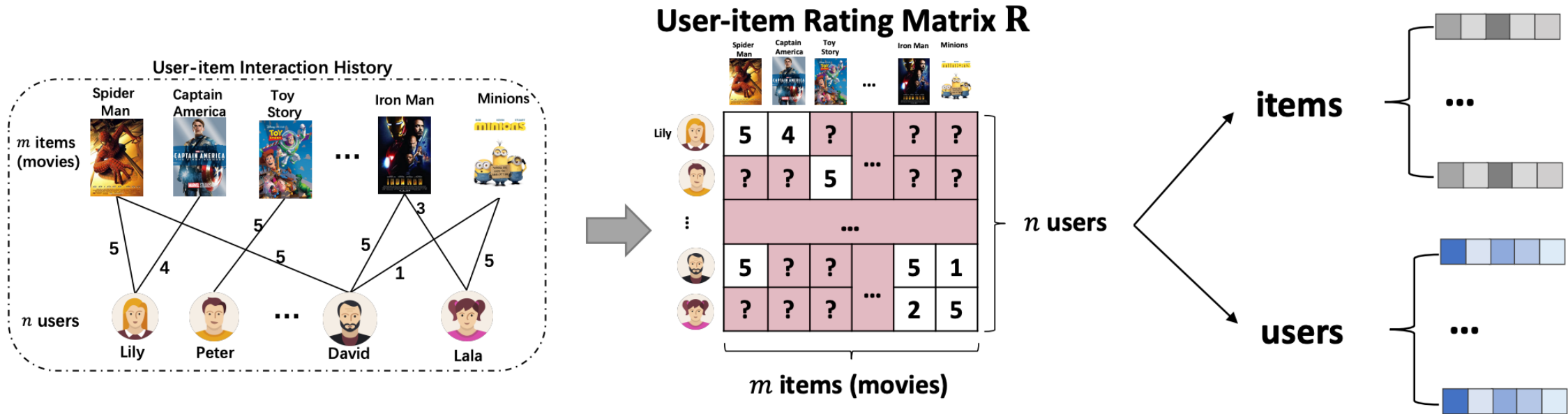


# Collaborative Filtering (CF)-based Recommendation

## CF for recommendation

- ❖ Similar users (with respect to their historical interactions) have similar preferences
- ❖ Modelling user's preferences on items based on their past interactions (e.g., ratings and clicks)

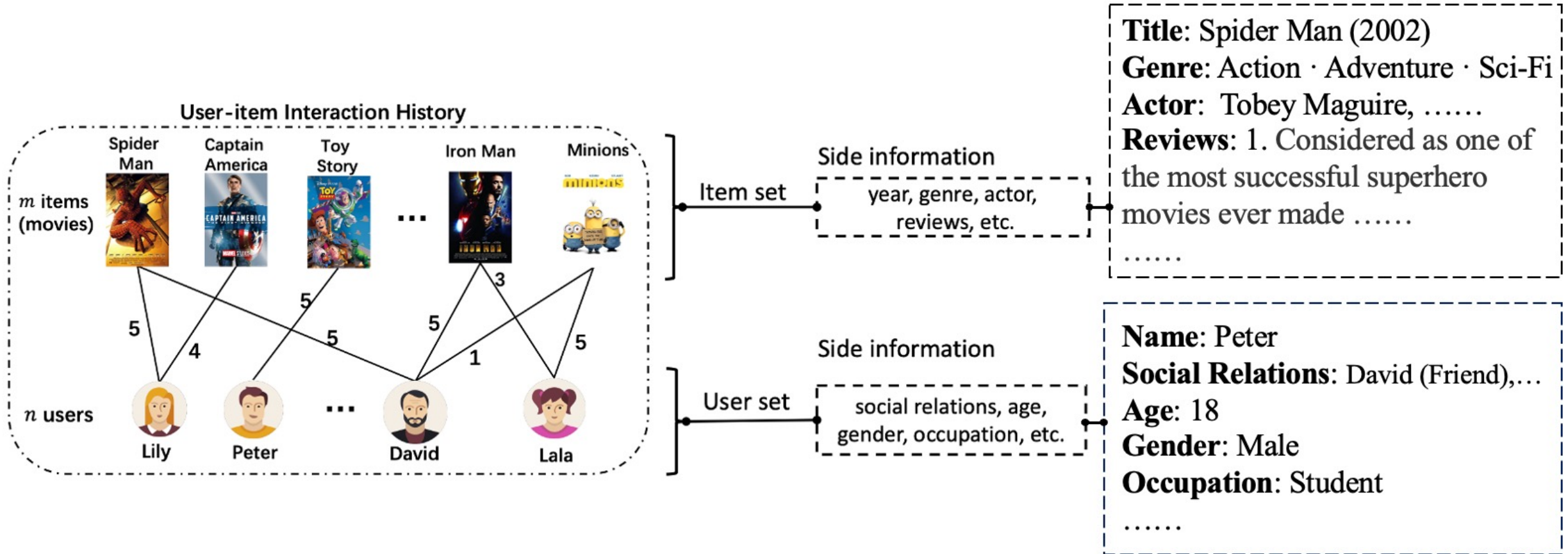
## Learning representations of users and items is the key to CF



# Content-based Recommendation



- ❑ Taking advantage of **additional knowledge/information** about users or items
- ❑ **Enhancing** user and item **representations** for improving recommendation performance



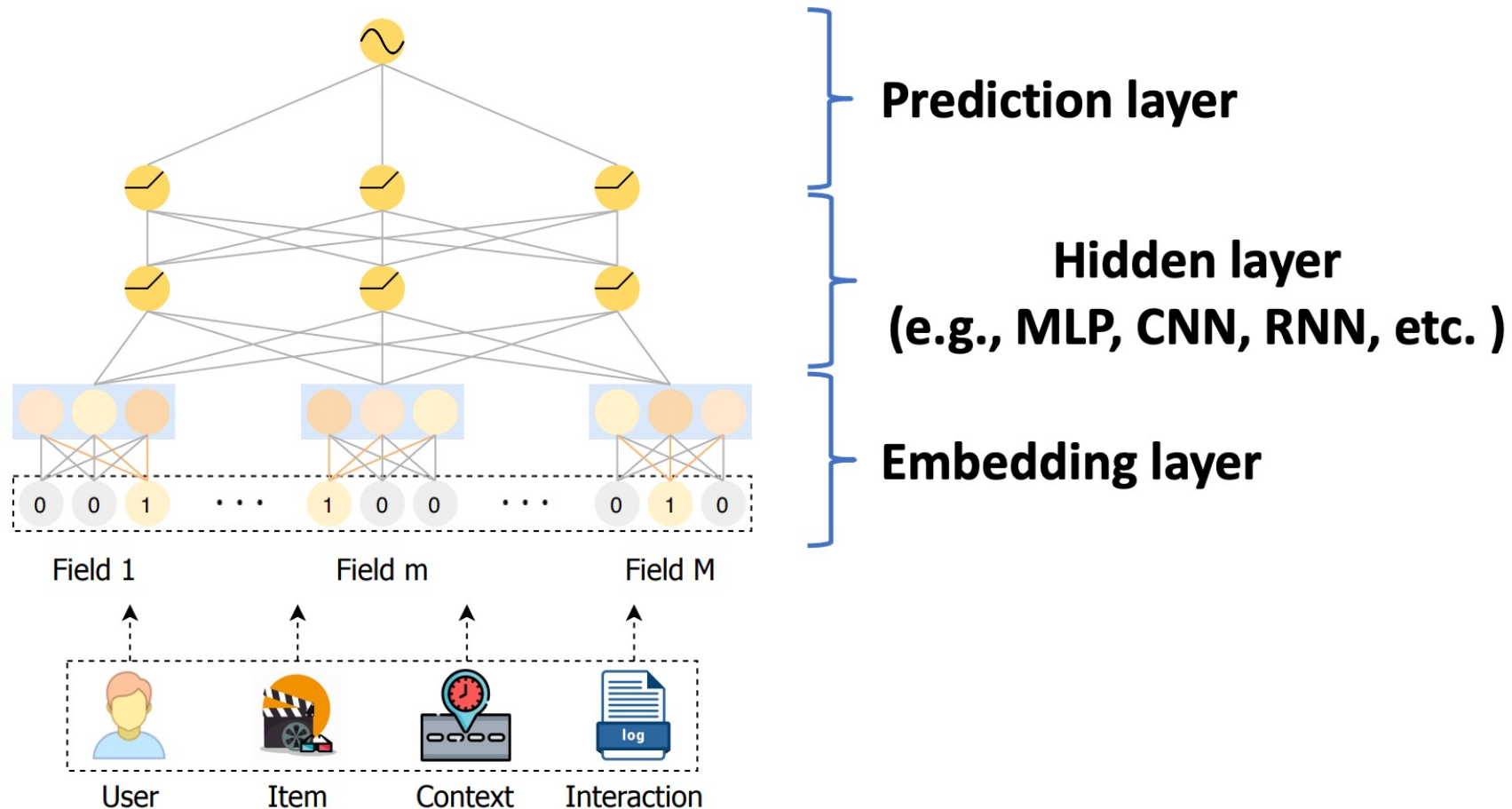
- ❑ Collaborative filtering + content == hybrid recommendation



# Deep Recommender Systems

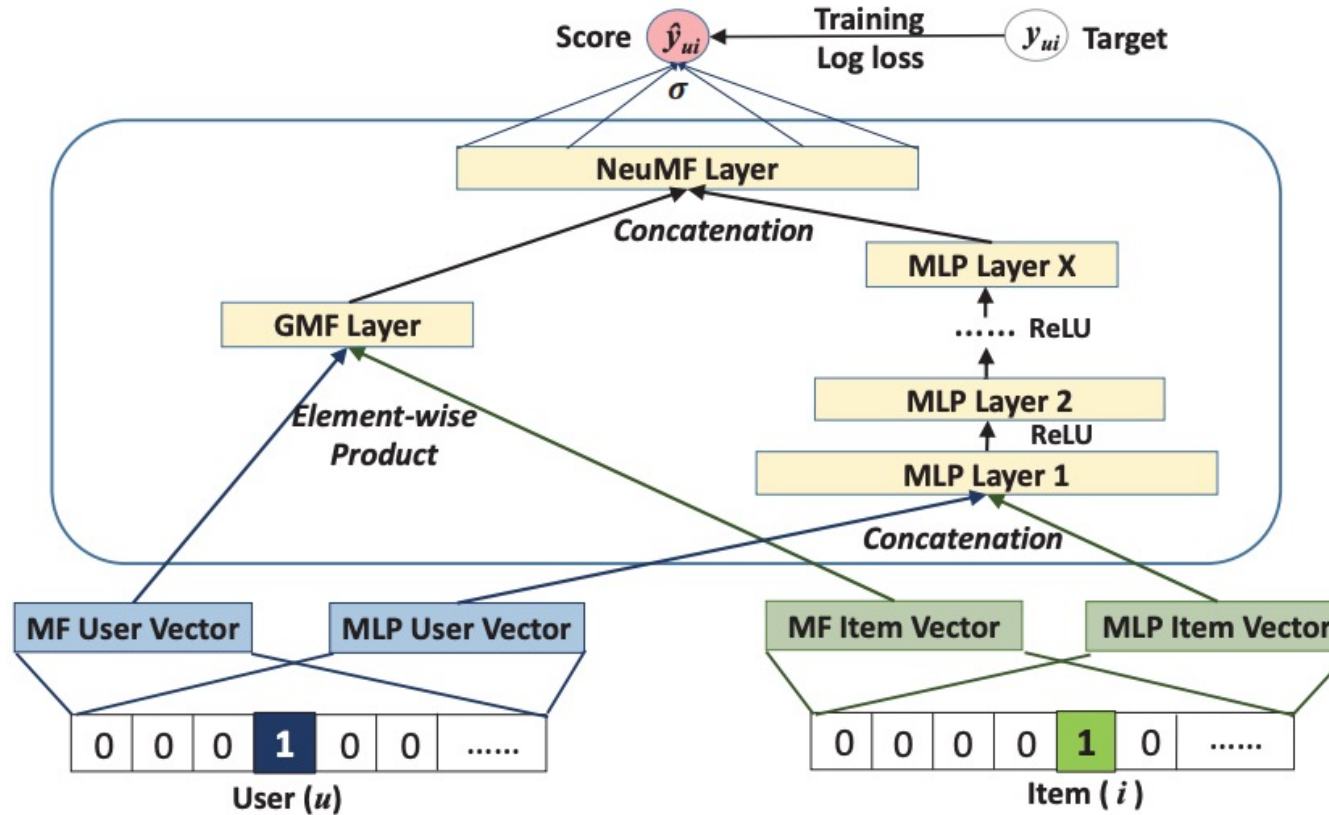


- ❑ Deep learning techniques have been effectively applied to develop recommender systems
- ❑ Remarkable representation learning capabilities



Neural Matrix Factorization (NeuMF) unifies the strengths of MF and MLP in modelling user-item interactions

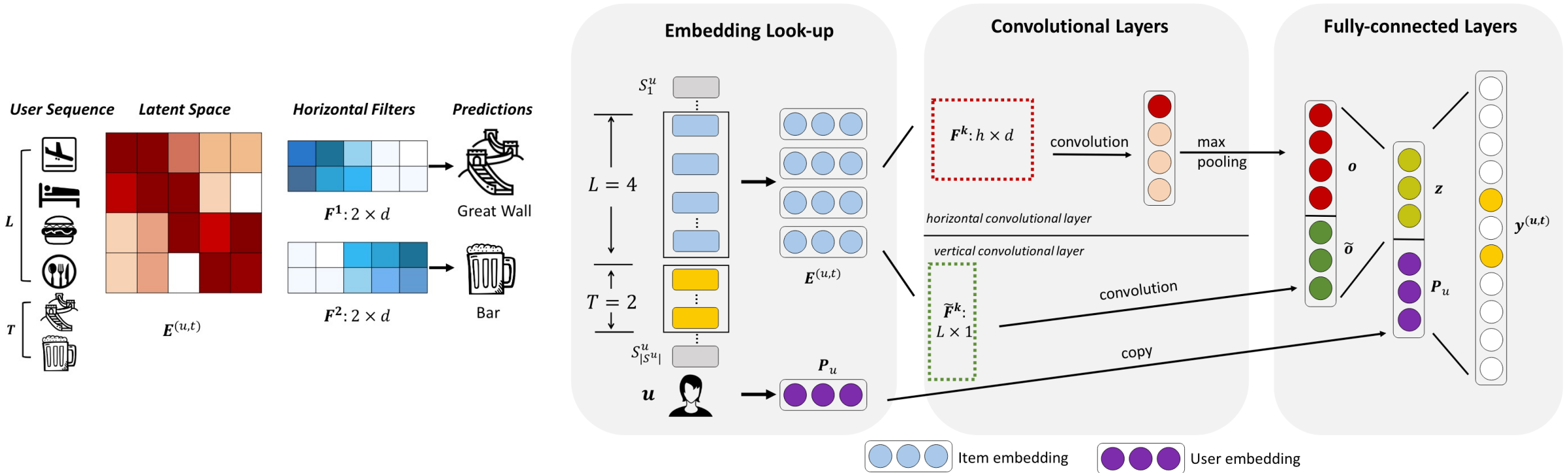
- MF uses an inner product as the interaction function
- MLP may be more capable to capture the complex structure of the interaction patterns



# Caser



- ❑ **Top-N sequential recommendation** models each user as a **sequence of items** interacted in the past and aims to **predict top-N ranked items**
- ❑ Convolutional Sequence Embedding Recommendation Model

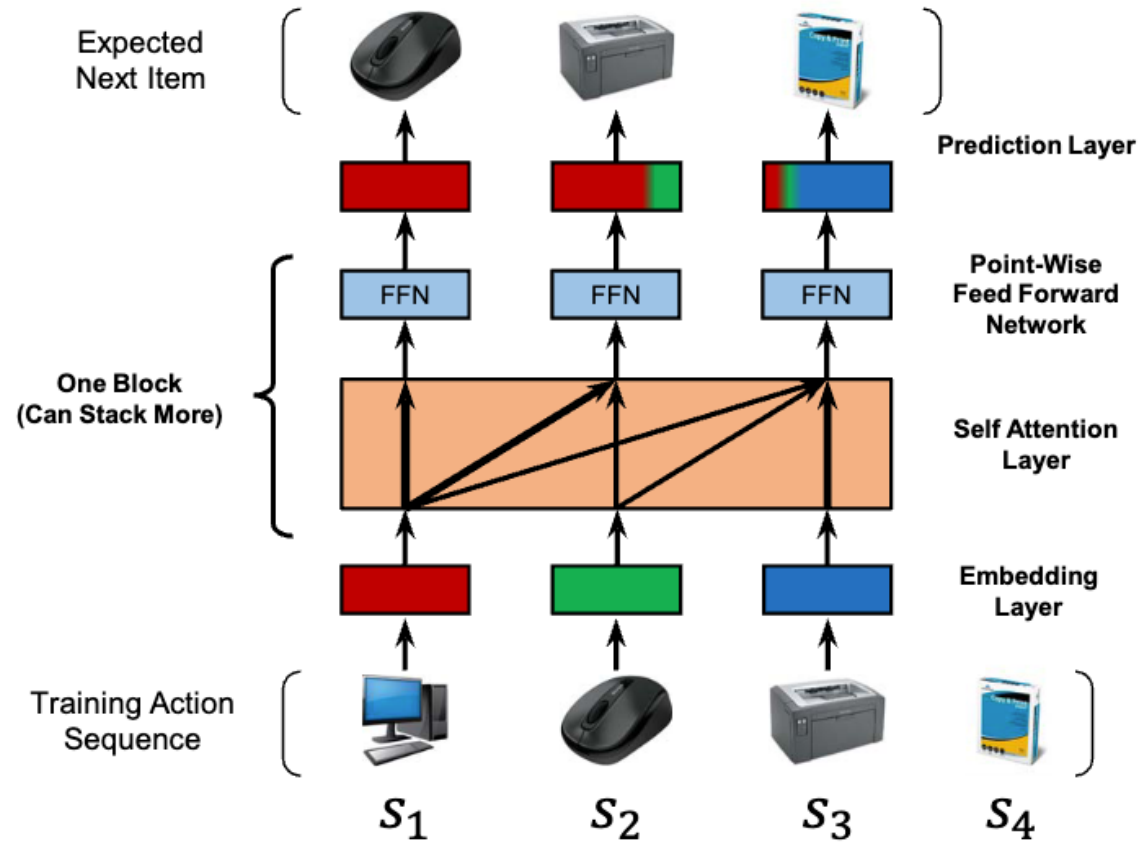




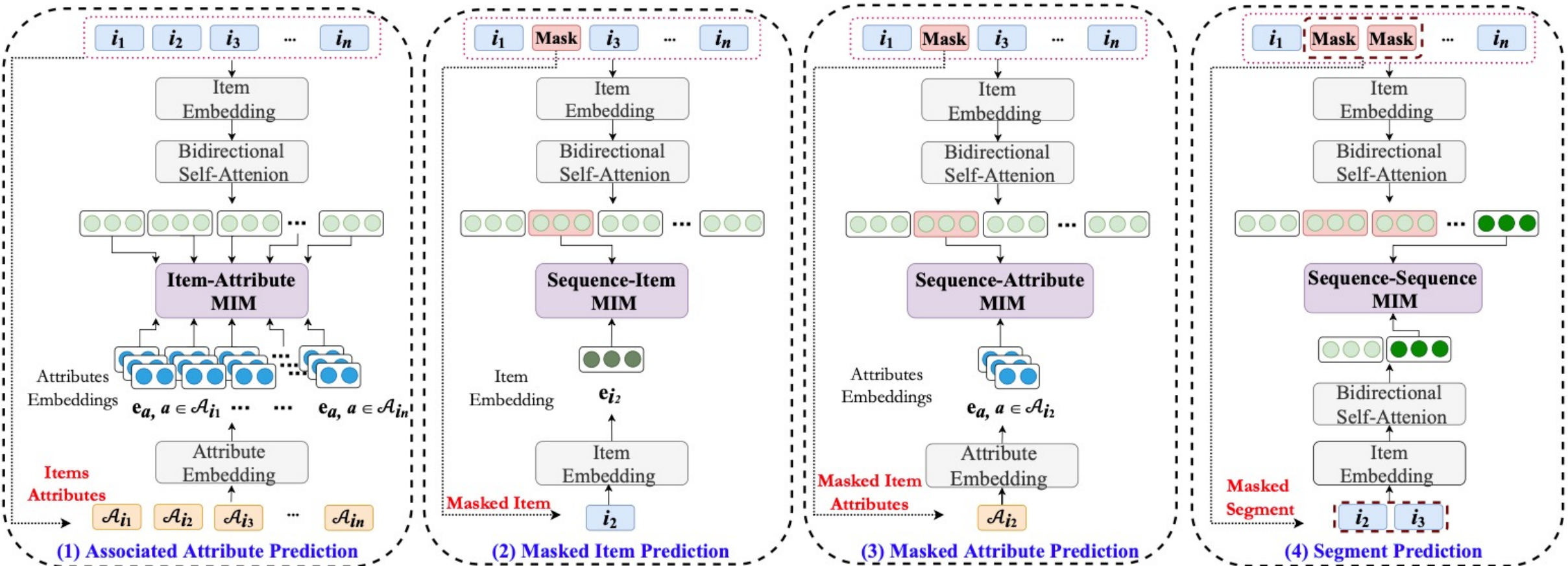
# SASRec



- ❑ Self-Attentive Sequential Recommendation
- ❑ Using an **attention** mechanism to capture **long-term semantics** and makes its **predictions** based on relatively **few actions**



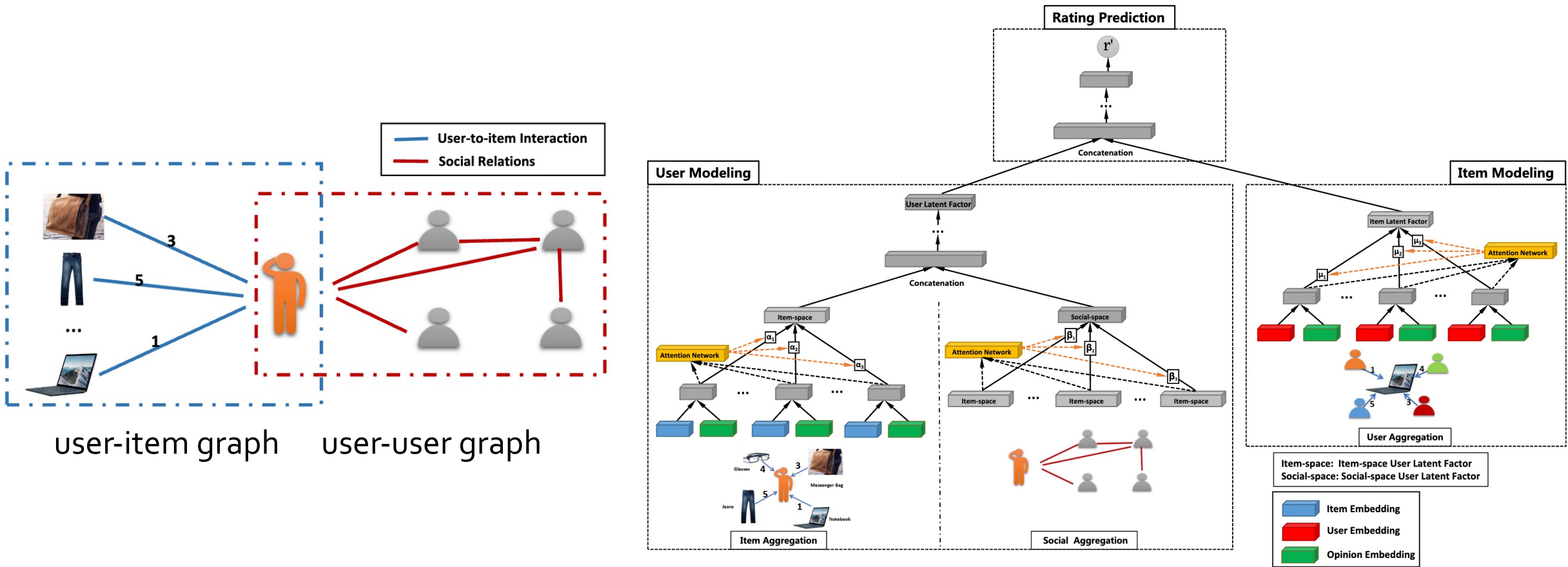
- Utilizing the intrinsic data correlation to derive **self-supervision** signals
- Enhancing the data representations via **pre-training** methods



# GraphRec



- Data in social recommender systems can be represented as **user-user** social graph and **user-item** graph



# PART 2: Preliminaries of RecSys and LLMs



Website of this tutorial

- ⊙ **Recommender Systems (RecSys)**
  - ⊙ Collaborative Filtering (CF)
  - ⊙ Content-based Recommendation
  - ⊙ Deep Recommender Systems
- **Large Language Models (LLMs)**
  - Development and Capability
  - LLM Architecture
- **LLM-based RecSys**
  - ID-based LLM RecSys
  - Text-based LLM RecSys



# Emergence of Large Language Models (LLMs)



- LLMs can be used for a variety of tasks, such as **Image Generation**

## Text to Image

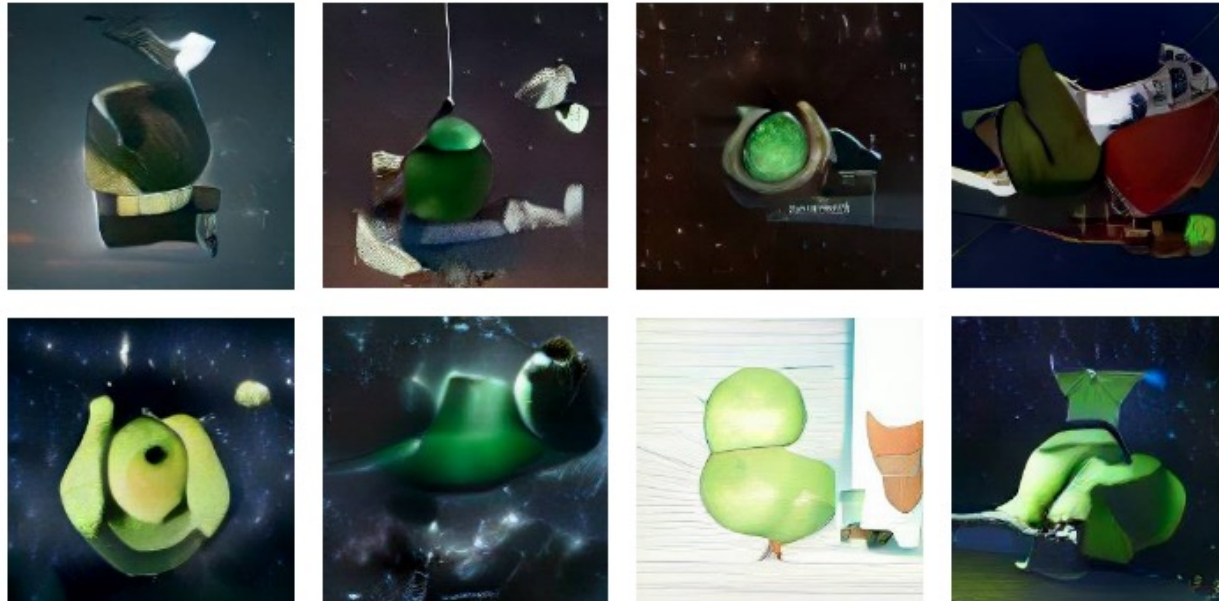
DALL·E mini

Generate images from text

What do you want to see?

an avocado armchair flying into space

an avocado armchair flying into space

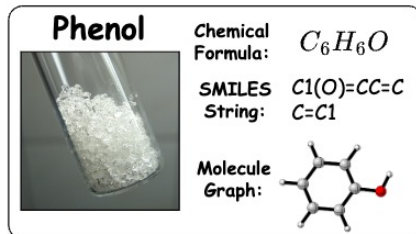


# Emergence of Large Language Models (LLMs)

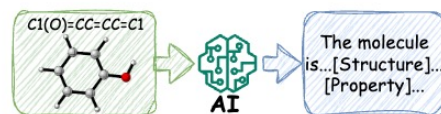


- LLMs can be used for a variety of tasks, such as **Molecule Generation**

## Text to Molecule



(a) Molecule Representations.



(b) Molecule Captioning.



(c) Text-based Molecule Generation.

**ChatGPT**

(a) Molecule Captioning

Please show me a description of this molecule:  
"C1=CC=C(C=C1)OC2=CC=CC=C2"

The molecule is an aromatic ether in which the oxygen is attached to two phenyl substituents. It has been found in muscat grapes and vanilla. It has a role as a plant metabolite.

(b) Text-based Molecule Generation

Help me generate a molecule based on the given description:  
"The molecule is a quinolinemonocarboxylate that is the conjugate base of xanthurenic acid, obtained by deprotonation of the carboxy group. It has a role as an animal metabolite. It is a conjugate base of a xanthurenic acid."

C1=CC2=C(C(=C1)[O-])NC(=CC2=O)C(=O)O

(d) Empowering ChatGPT with molecule captioning and text-based molecule generation abilities.



# Emergence of Large Language Models (LLMs)



□ LLMs can be used for a variety of tasks, such as **Recommendation**

## Text to Recommendation

### Rating Prediction

zero-shot How will user rate this product\_title: "SHANY Nail Art Set (24 Famous Colors Nail Art Polish, Nail Art Decoration)" , and product\_category: Beauty? ( 1 being lowest and 5 being highest ) Attention! Just give me back the exact number a result , and you don't need a lot of text.

few-shot Here is user rating history:  
1. Bundle Monster 100 PC 3D Designs Nail Art Nailart Manicure Fimo Canes Sticks Rods Stickers Gel Tips, 5.0;  
2. Winstonia's Double Ended Nail Art Marbling Dotting Tool Pen Set w/ 10 Different Sizes 5 Colors - Manicure Pedicure, 5.0;  
3. Nail Art Jumbo Stamp Stamping Manicure Image Plate 2 Tropical Holiday by Cheeky&reg, 5.0 ;  
4.Nail Art Jumbo Stamp Stamping Manicure Image Plate 6 Happy Holidays by Cheeky&reg, 5.0;  
Based on above rating history, please predict user's rating for the product: "SHANY Nail Art Set (24 Famouse Colors Nail Art Polish, Nail Art Decoration)", (1 being lowest and5 being highest,The output should be like: (x stars, xx%), do not explain the reason.)

### Sequential Recommendation

zero-shot Requirements: you must choose 10 items for recommendation and sort them in order of priority, from highest to lowest. Output format: a python list. Do not explain the reason or include any other words.  
The user has interacted with the following items in chronological order: ['Better Living Classic Two Chamber Dispenser, White', 'Andre Silhouettes Shampoo Cape, Metallic Black', '.....', 'John Frieda JFHA5 Hot Air Brush, 1.5 inch'].Please recommend the next item that the user might interact with.

few-shot Requirements: you must choose 10 items for recommendation and sort them in order of priority, from highest to lowest. Output format: a python list. Do not explain the reason or include any other words.  
Given the user's interaction history in chronological order: ['Avalon Biotin B-Complex Thickening Conditioner, 14 Ounce', 'Conair 1600 Watt Folding Handle Hair Dryer', '.....', 'RoC Multi-Correxion 4-Zone Daily Moisturizer, SPF 30, 1.7 Ounce'], the next interacted item is ['Le Edge Full Body Exfoliator - Pink']. Now, if the interaction history is updated to ['Avalon Biotin B-Complex Thickening Conditioner, 14 Ounce', 'Conair 1600 Watt Folding Handle Hair Dryer', '.....', 'RoC Multi-Correxion 4-Zone Daily Moisturizer, SPF 30, 1.7 Ounce', 'Le Edge Full Body Exfoliator - Pink'] and the user is likely to interact again, recommend the next item.



# What are Language Models?



## □ Narrow Sense

- ❖ A **probabilistic model** that assigns a probability to every **finite sequence** (grammatical or not)

Sentence: “the cat sat on the mat”

$$P(\text{the cat sat on the mat}) = P(\text{the}) * P(\text{cat}|\text{the}) * P(\text{sat}|\text{the cat}) \\ * P(\text{on}|\text{the cat sat}) * P(\text{the}|\text{the cat sat on}) \\ * P(\text{mat}|\text{the cat sat on the})$$

Implicit order



## □ Broad Sense

- ❖ Encoder-only models (BERT, RoBERTa, ELECTRA)
- ❖ Decoder-only models (GPT-X, OPT, LLaMa, PaLM)
- ❖ Encoder-decoder models (T5, BART)

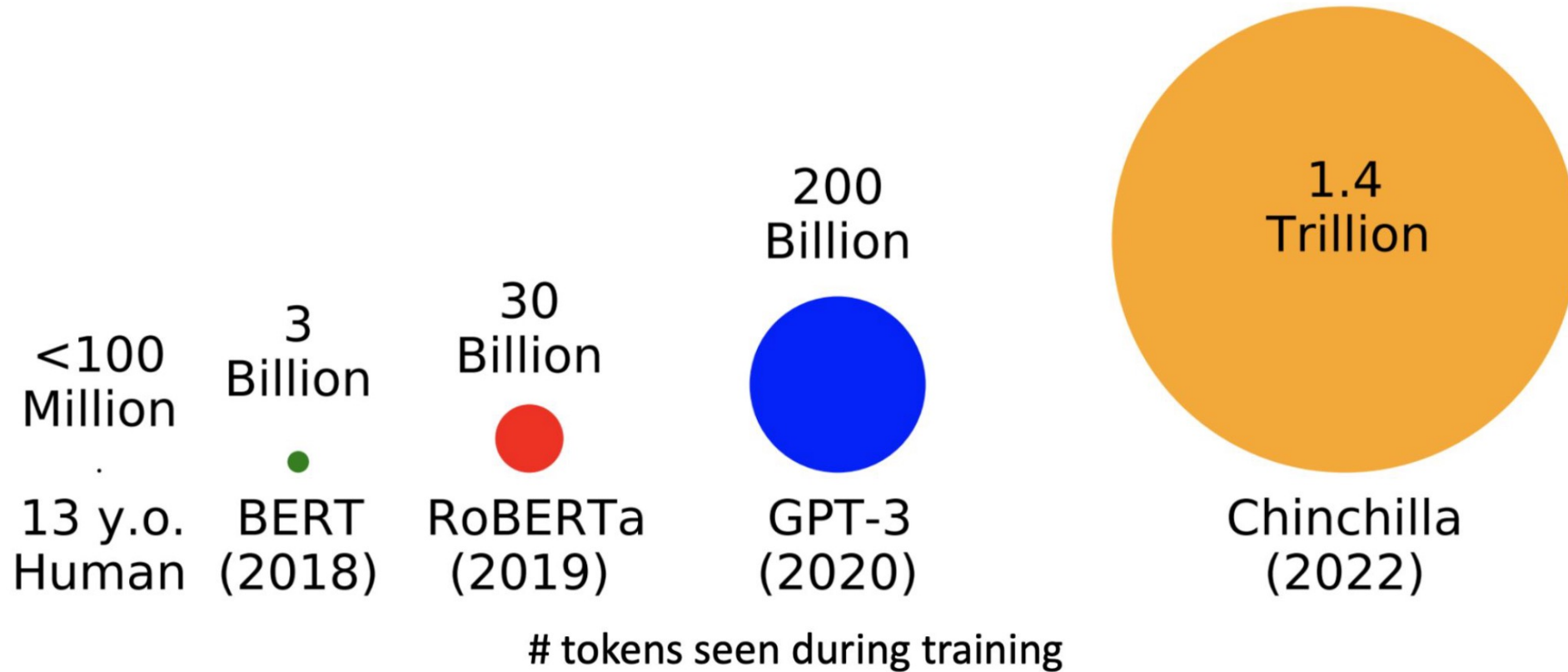




# Large Language Model Development



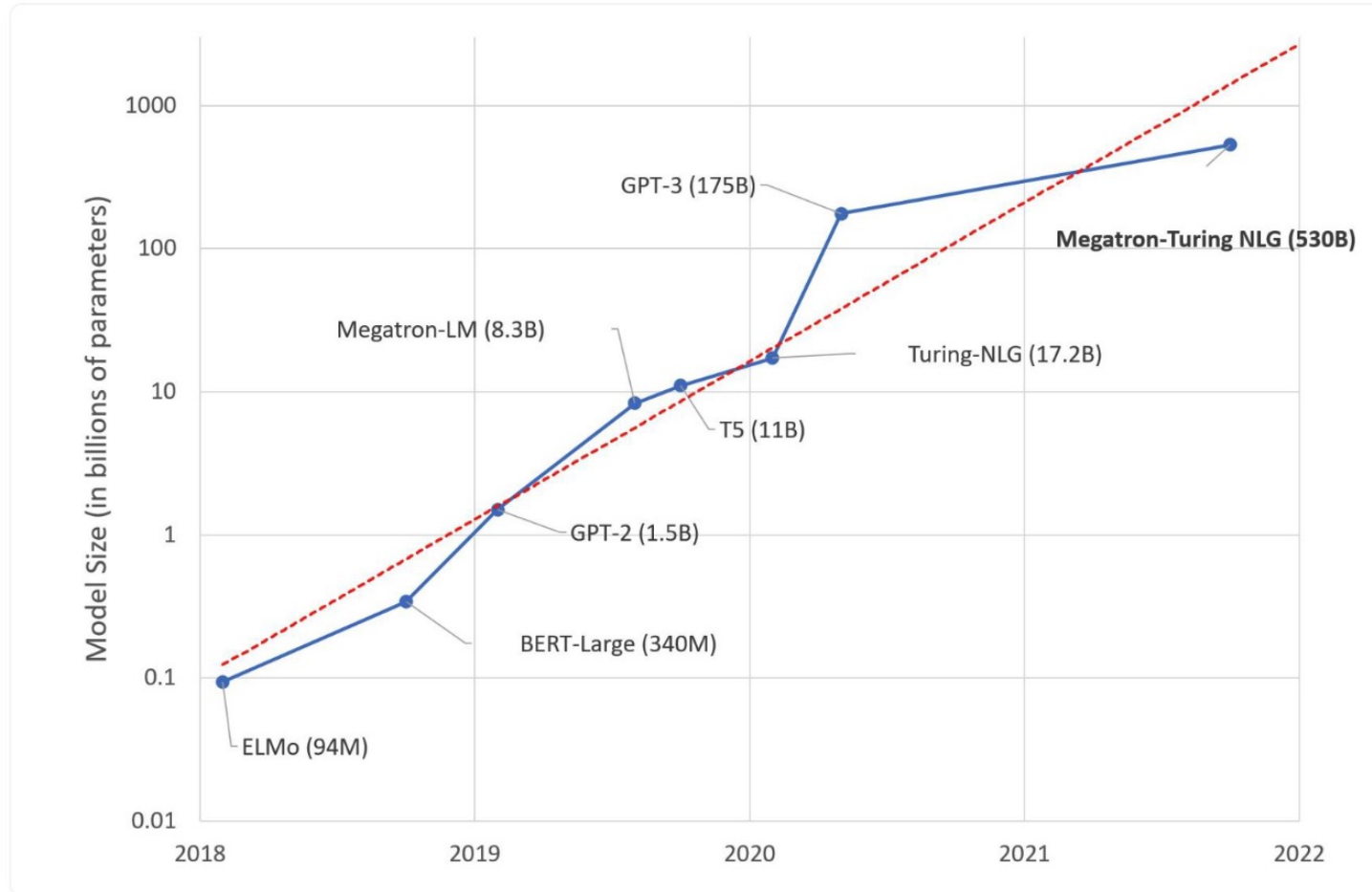
- ☐ Trained on more and more data – **Hundreds of Billions of Tokens**



# Large Language Model Development



- ❑ Larger and larger models – **Billions of Parameters**

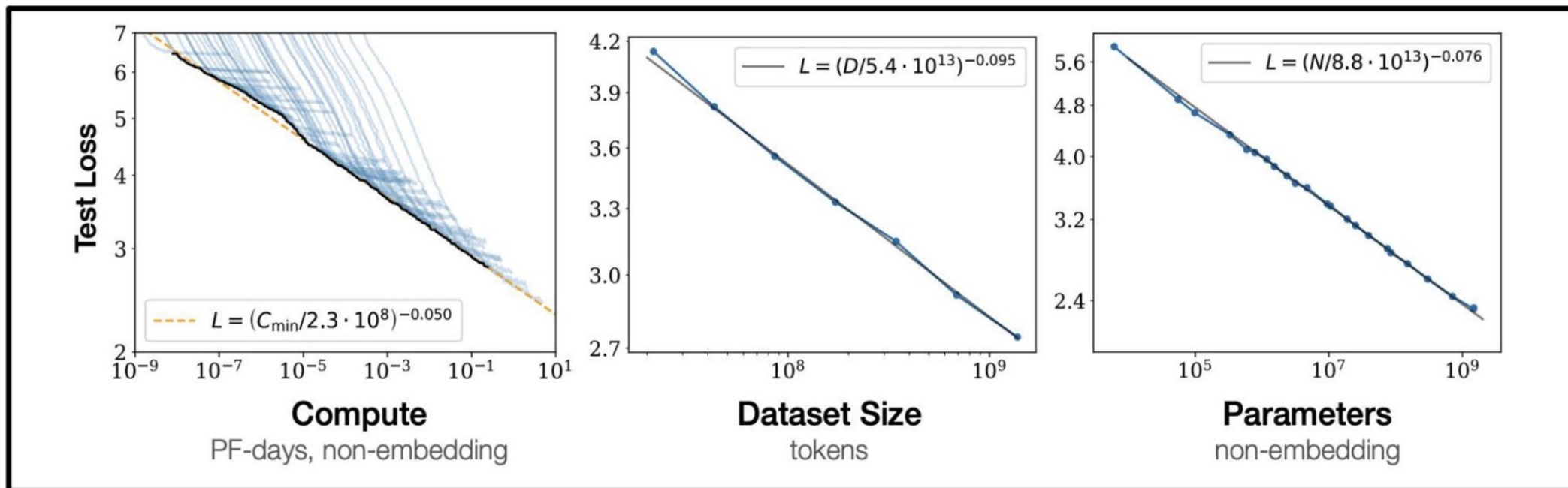


# Why Large Language Models?



## Scaling Law for Neural Language Models

- ❖ Performance depends strongly on scale! We keep getting better performance as we scale the model, data, and compute up!

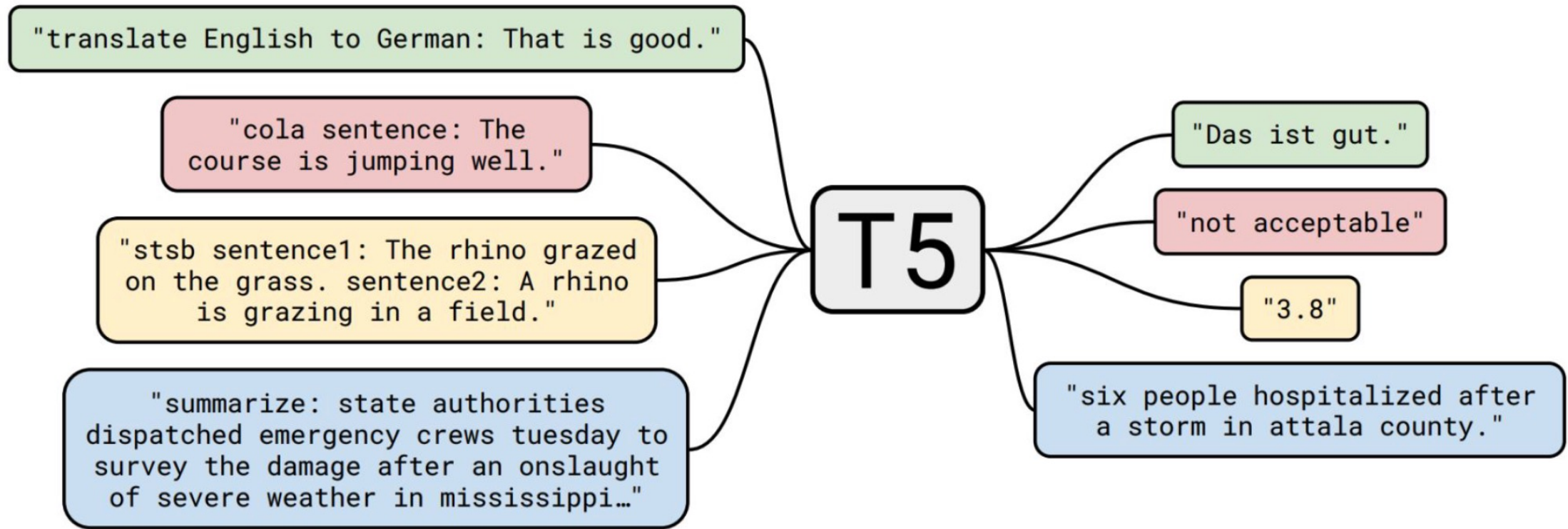


# Why Large Language Models?



## □ Generalization

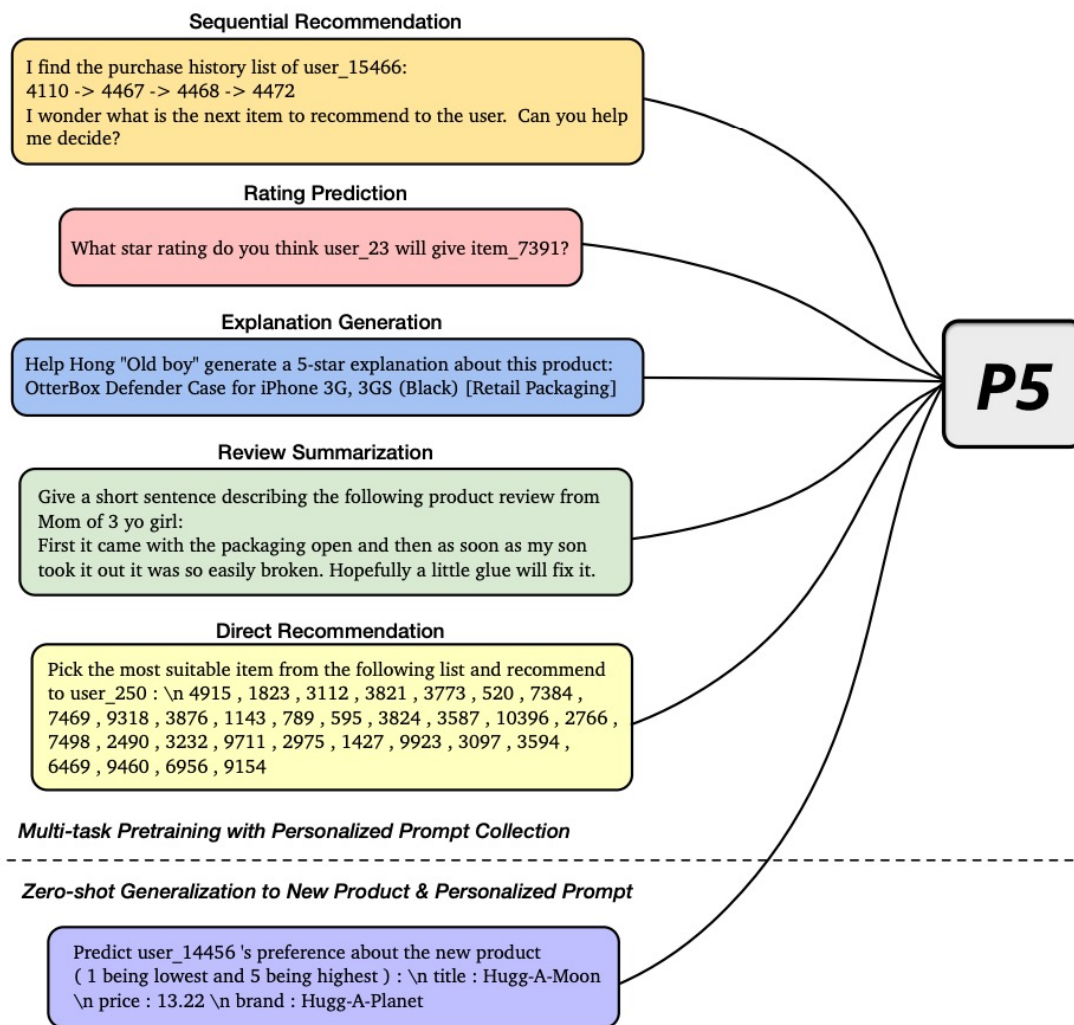
- ❖ We can now use one single model to solve many NLP tasks.



# Why Large Language Models?



## Strong Zero-shot/Few-shot Ability



## Multiple Tasks in One Model

- Sequential recommendation
- Rating prediction
- Explain generation
- Review summarization
- Direct recommendation



# Large Language Model Structure



## Encoder-Only Models

- ❖ BERT, RoBERTa, ELECTRA

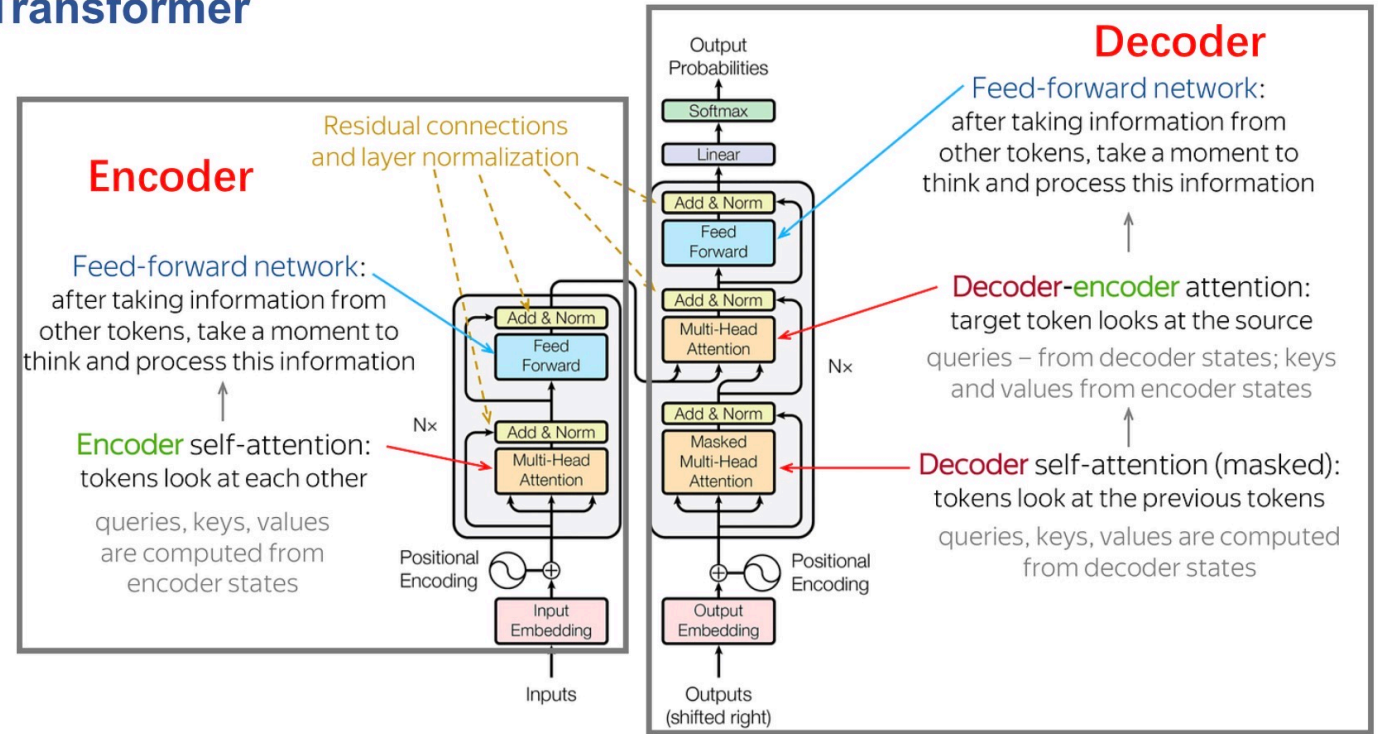
## Decoder-Only Models

- ❖ GPT-X, OPT, LLaMa, PaLM

## Encoder-Decoder Models

- ❖ T5, BART

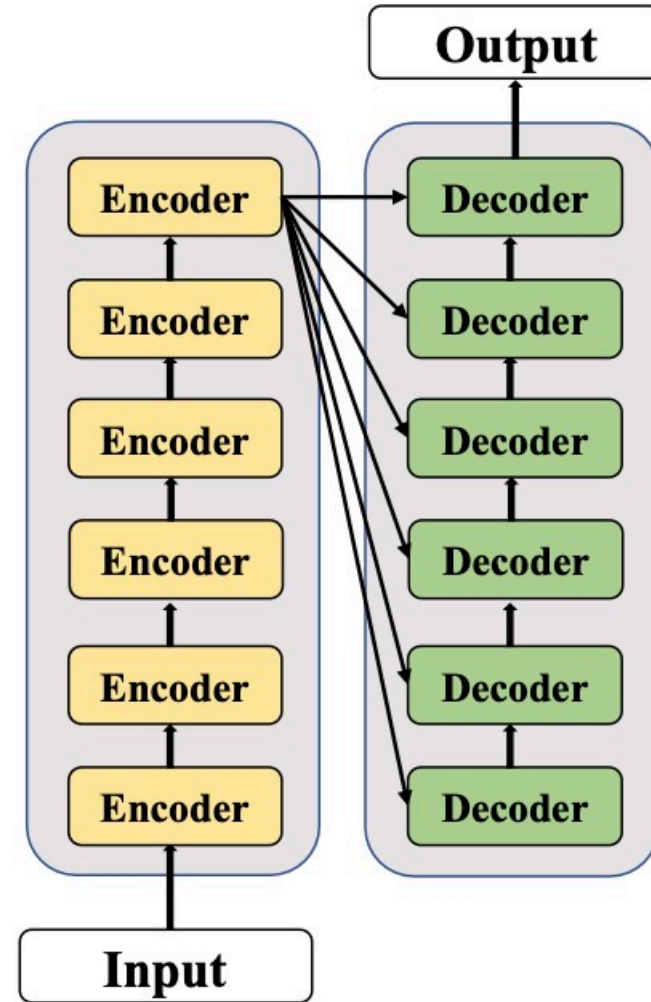
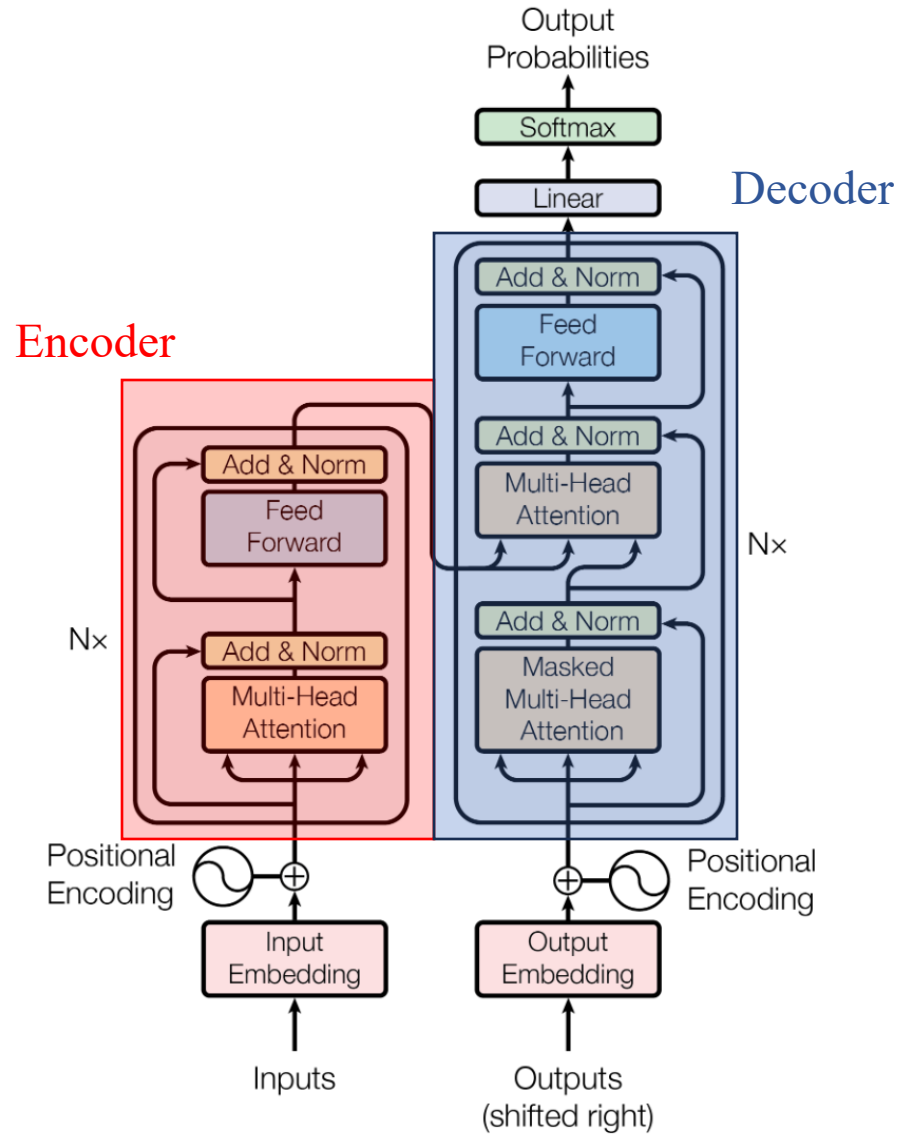
## Transformer



The Transformer – model architecture



# Transformer

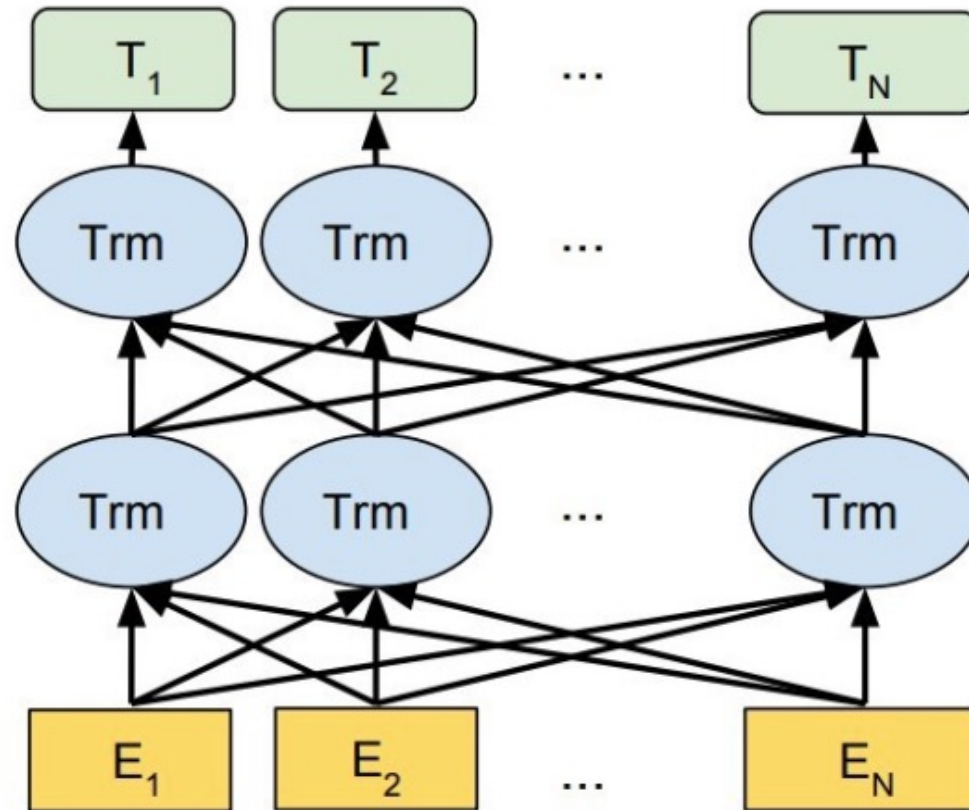


# Encoder-Only Models: BERT



- ❑ BERT uses a **bidirectional** Transformer

## BERT (Ours)

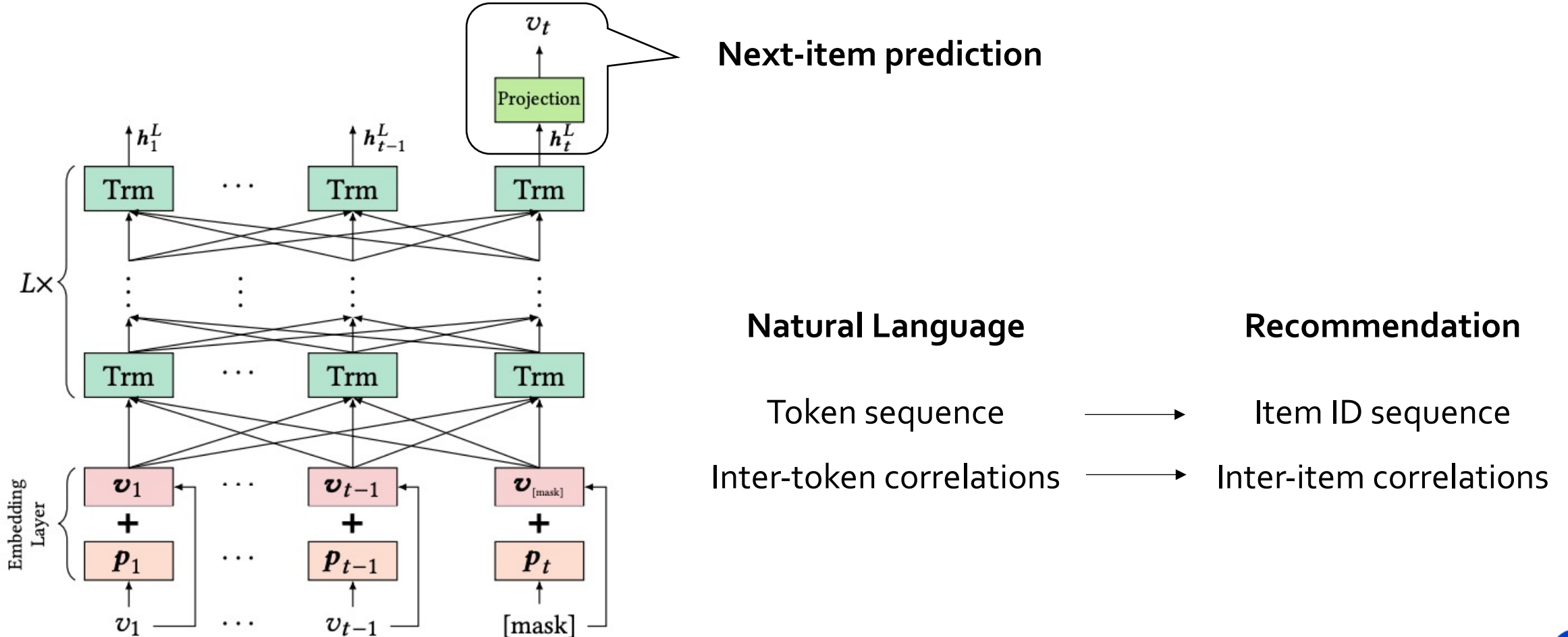




# Encoder-Only Models for Rec: BERT4Rec



- Adopt **Bidirectional Encoder Representations** from Transformers to model the sequential nature of user behaviors

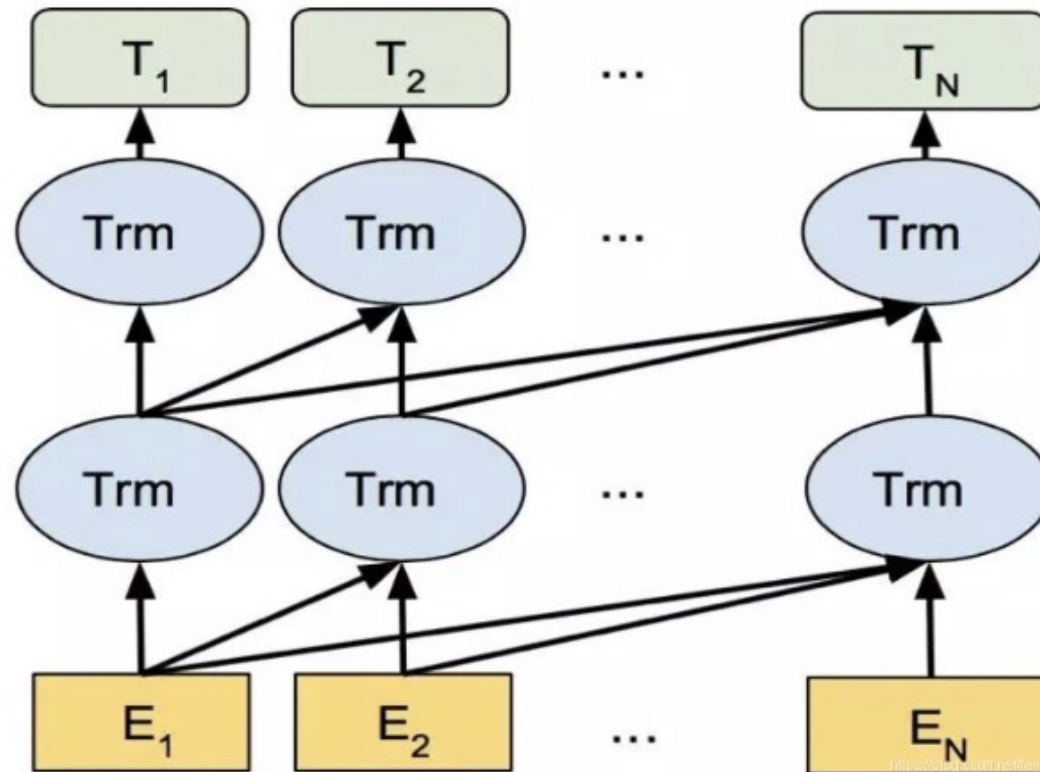


# Decoder-Only Models: GPT



- OpenAI GPT uses a **left-to-right** Transformer

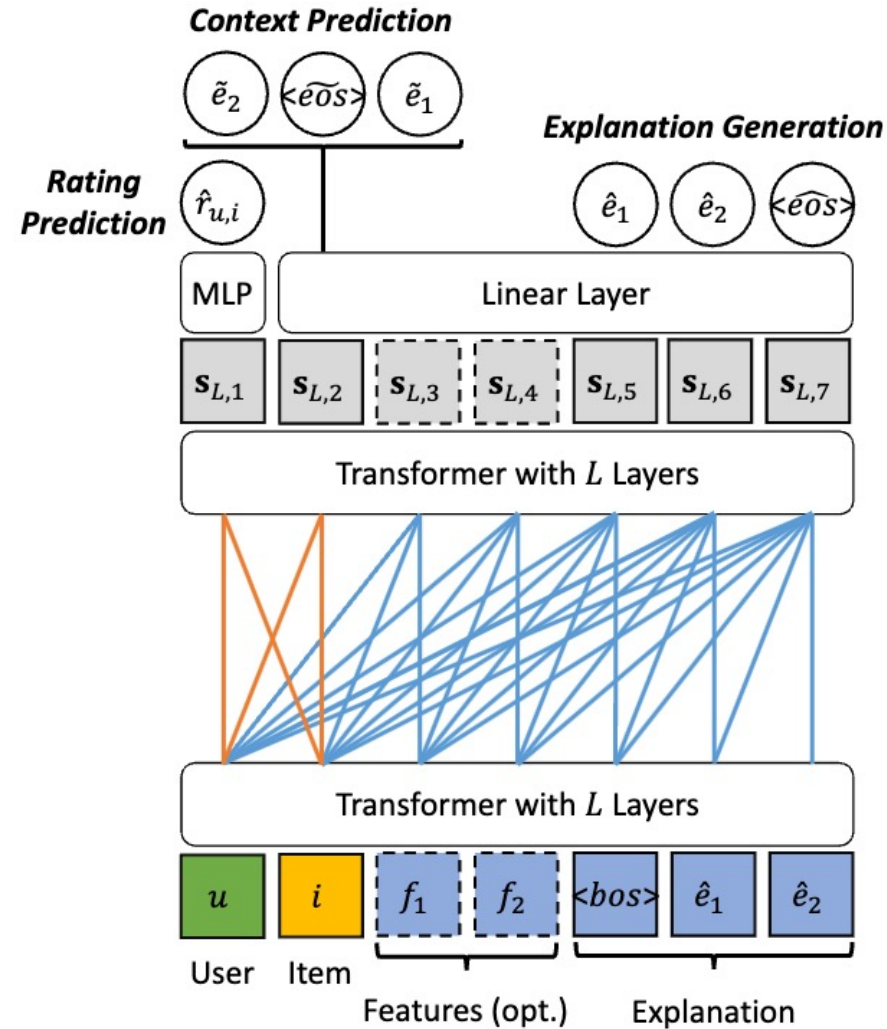
## OpenAI GPT



# Decoder-Only Models for Rec: PETER



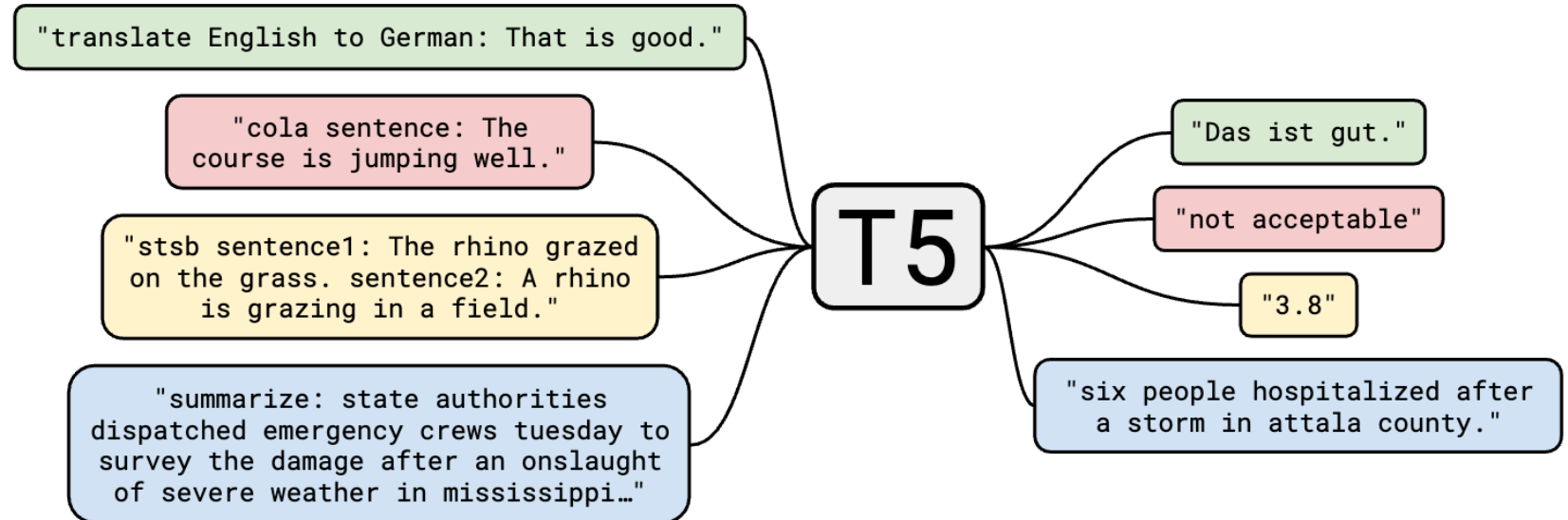
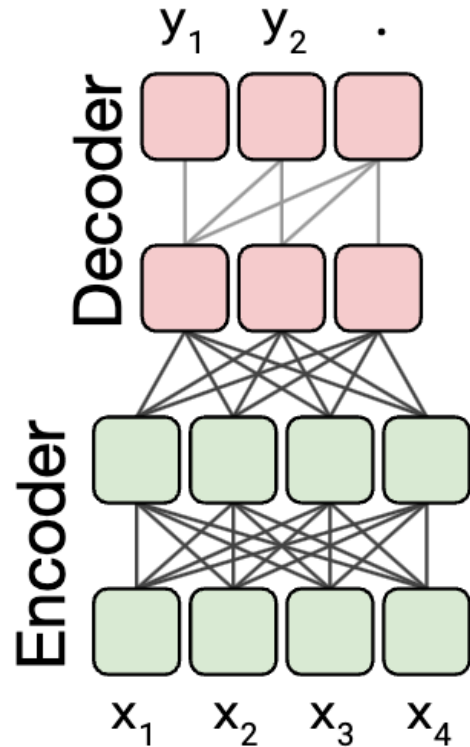
- Utilizing the IDs to predict the words in the target explanation



# Encoder-Decoder Models: T5



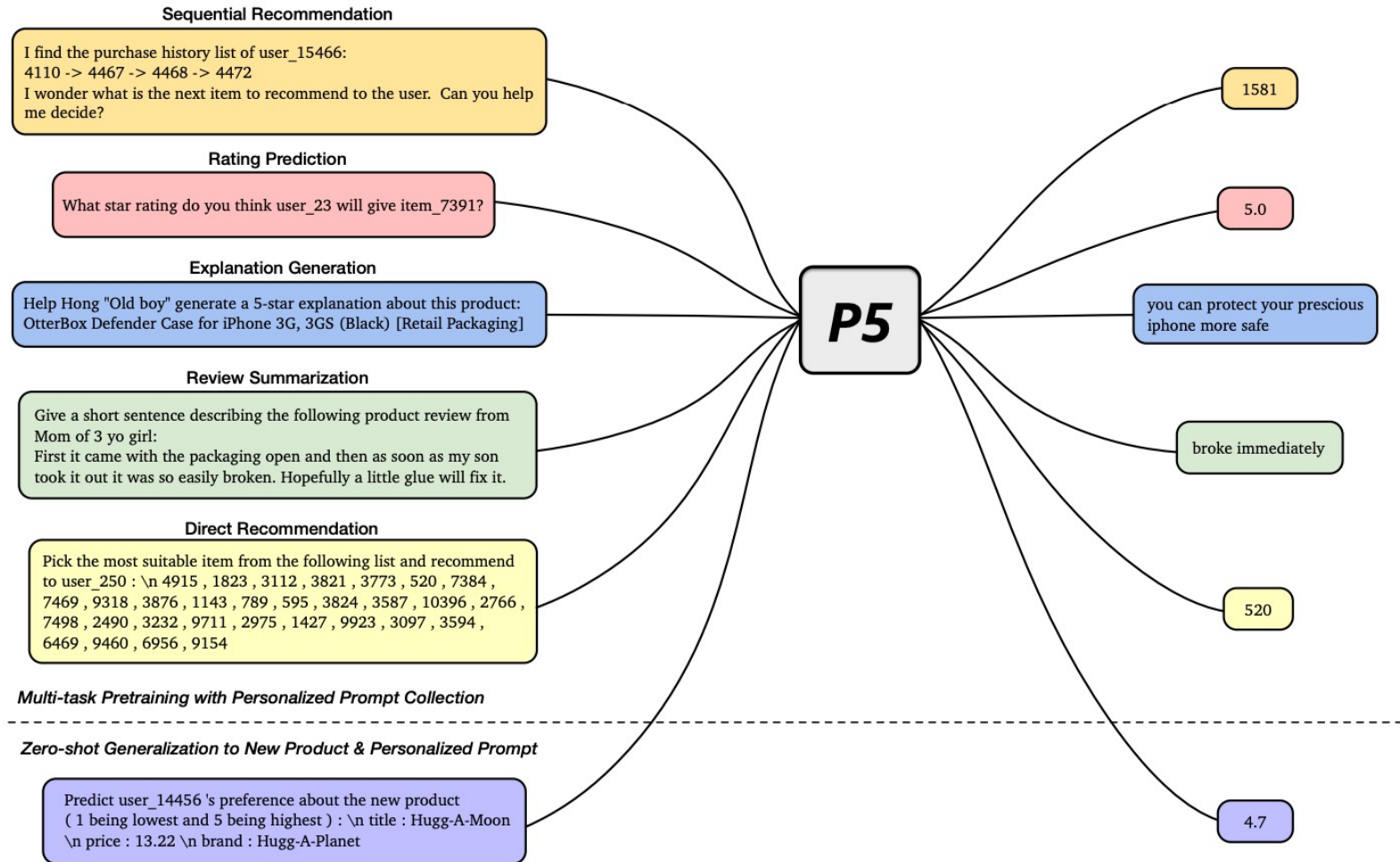
- ❑ T5 handles any text-to-text task by converting every natural language processing problem into a text generation problem.



# Encoder-Decoder Models for Rec: P5



- Text-to-text paradigm - “Pretrain, Personalized Prompt, and Predict Paradigm” (P5) for recommendation: converting five problems into a text generation problem.



# PART 2: Preliminaries of RecSys and LLMs



Website of this tutorial

- ⊙ **Recommender Systems (RecSys)**
  - ⊙ Collaborative Filtering (CF)
  - ⊙ Content-based Recommendation
  - ⊙ Deep Recommender Systems
- ⊙ **Large Language Models (LLMs)**
  - Development and Capability
  - LLM Architecture
- **LLM-based RecSys**
  - ID-based LLM RecSys
  - Text-based LLM RecSys



# LLM-based RecSys: ID-based & Text-based



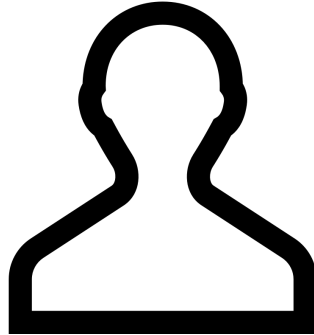
- ❑ Users and Items can be represented in various ways






Index or Content ?

## ○ User ID

U8189cf6745fc0d808977bdb0b9f22995

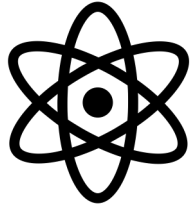
Username: Jack0513



Poster	Movie Name	Numeric ID
	In Broad Daylight	1697292155
	The Marvels	1699436461
	TAYLOR SWIFT   THE ERAS TOUR	1695730583
	The Dark Knight Rises	1699611567
	Oppenheimer	1687513232



# User & Item Representation in LLMs



ID-based LLM RecSys



Text-based LLM RecSys

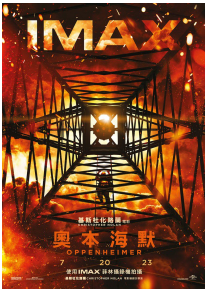
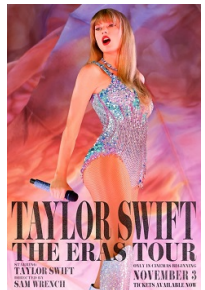





# ID-based LLM RecSys



- Various ways of assigning IDs

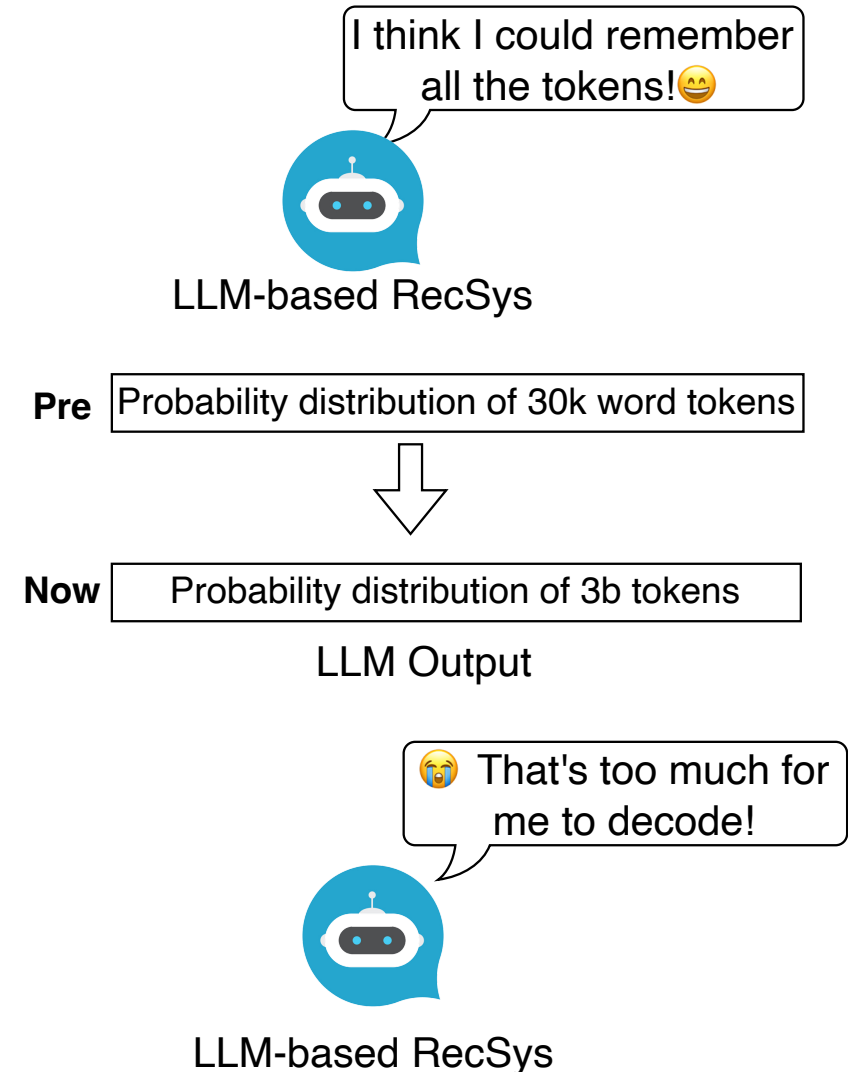
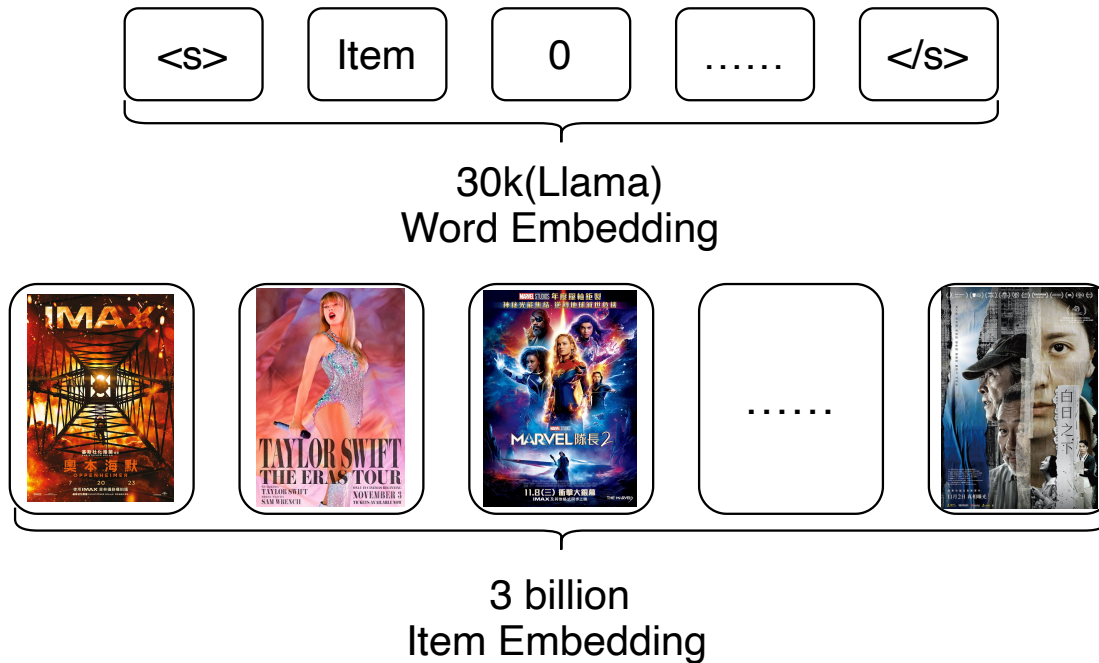
	Randomly	Based on Popularity	Based on Time
	AXGGWD027	01	1687513232
	XJSGDG0881	02	1695730583
	BXGW2UD803	03	1699436461



# ID-based LLM RecSys



- ❑ IDs are originally for unique identification
- ❑ However, the embedding of LLMs cannot hold millions of items and users



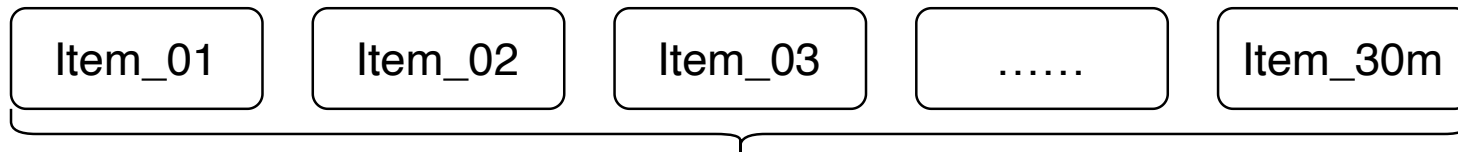
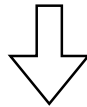
# ID-based LLM RecSys



- Normally, we can represent users and items with a span of tokens.
- The format is like “[Prefix]\_[ID]”. Examples:
  - ❖ User\_0123 : [“User”, “\_”, “0”, “1”, “2”, “3”]
  - ❖ Item\_5471 : [“Item”, “\_”, “5”, “4”, “7”, “1”]

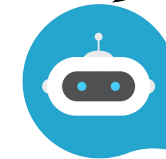


3 billion



Token Spans

😊 No extra decoding cost!  
I can handle that!



LLM-based RecSys

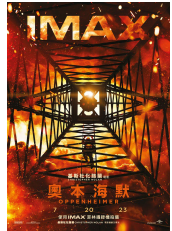
- ❖ However, for Item\_1003, it could be [“Item”, “\_”, “100”, “3”], which might be confusing for LLMs!



# ID-based LLM RecSys



- Indexing methods might affect the performance of RecSys



01

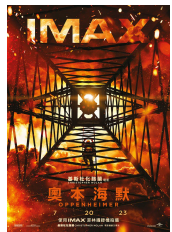
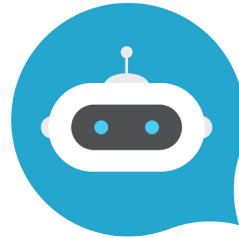


02



03

Is Item "04" still a movie?

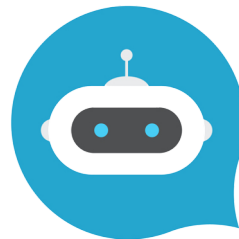


01



02

Do Item "01" and Item "02" share similar characteristics?

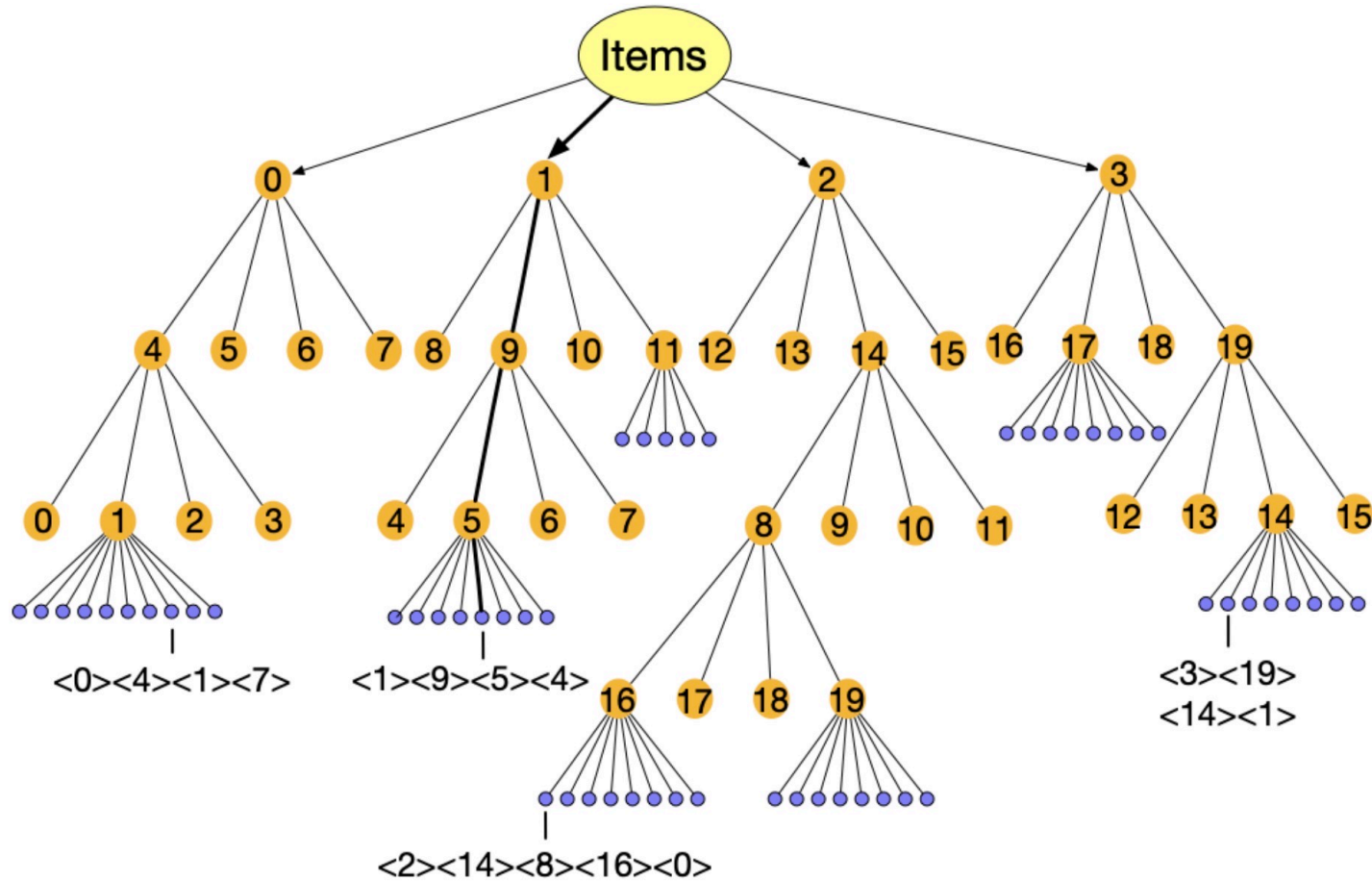


# ID-based LLM RecSys



□ Introducing more Information to enhance the ID representation

❖ Collaborative Indexing

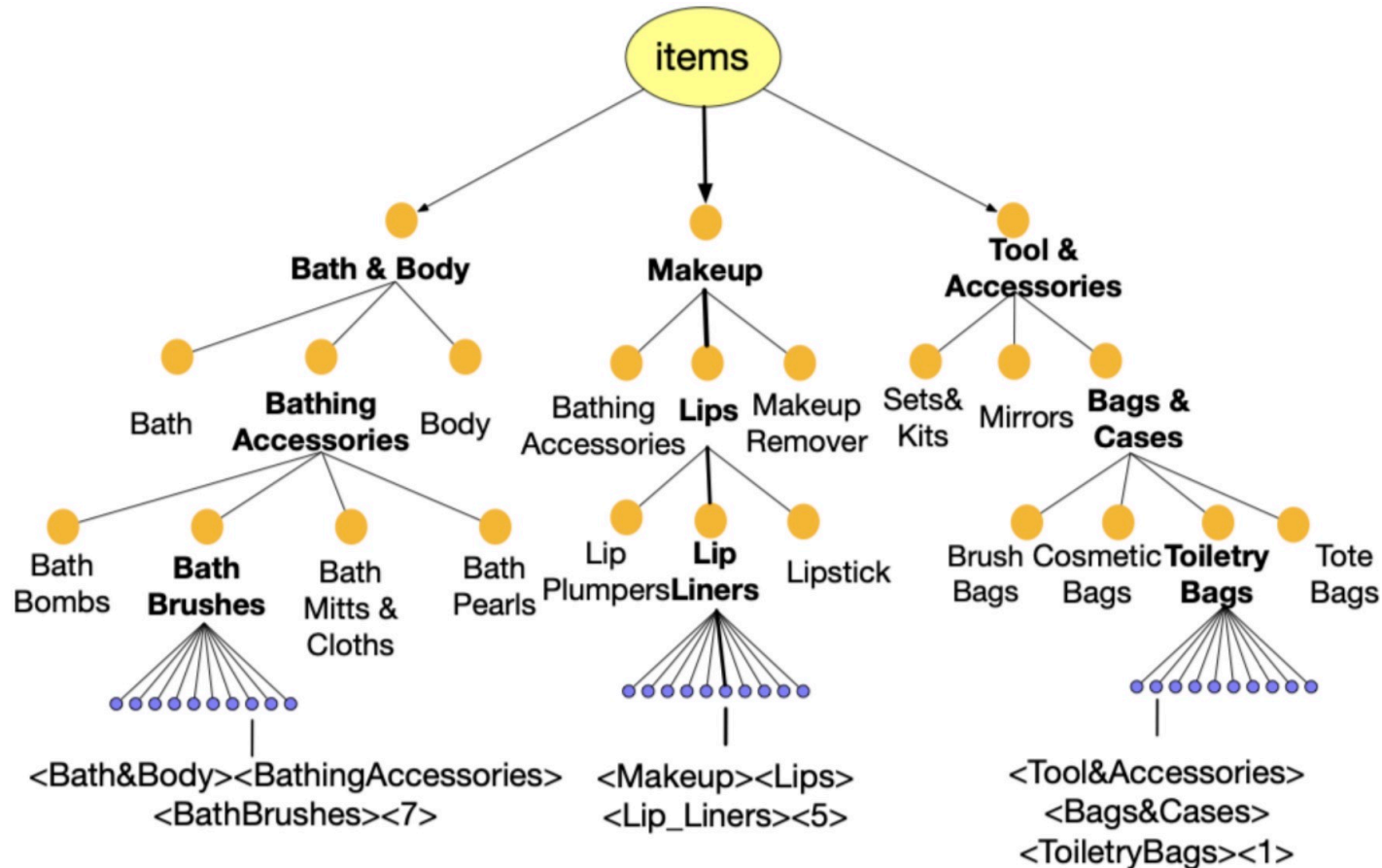


# ID-based LLM RecSys



□ Introducing more Information to enhance the ID representation


❖ Semantic Indexing



# ID-based LLM RecSys

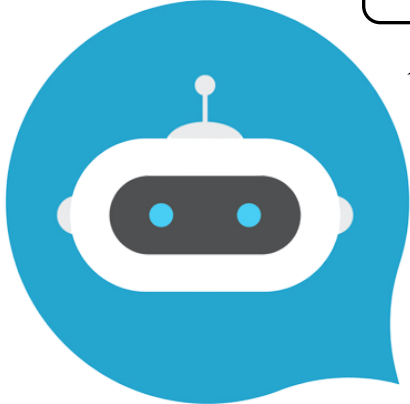


- Modeling user interaction history with Markov chain



Question

User\_4782 has bought  
3472, 7653, 0192, 4271.  
What will he buy next?



RecSys

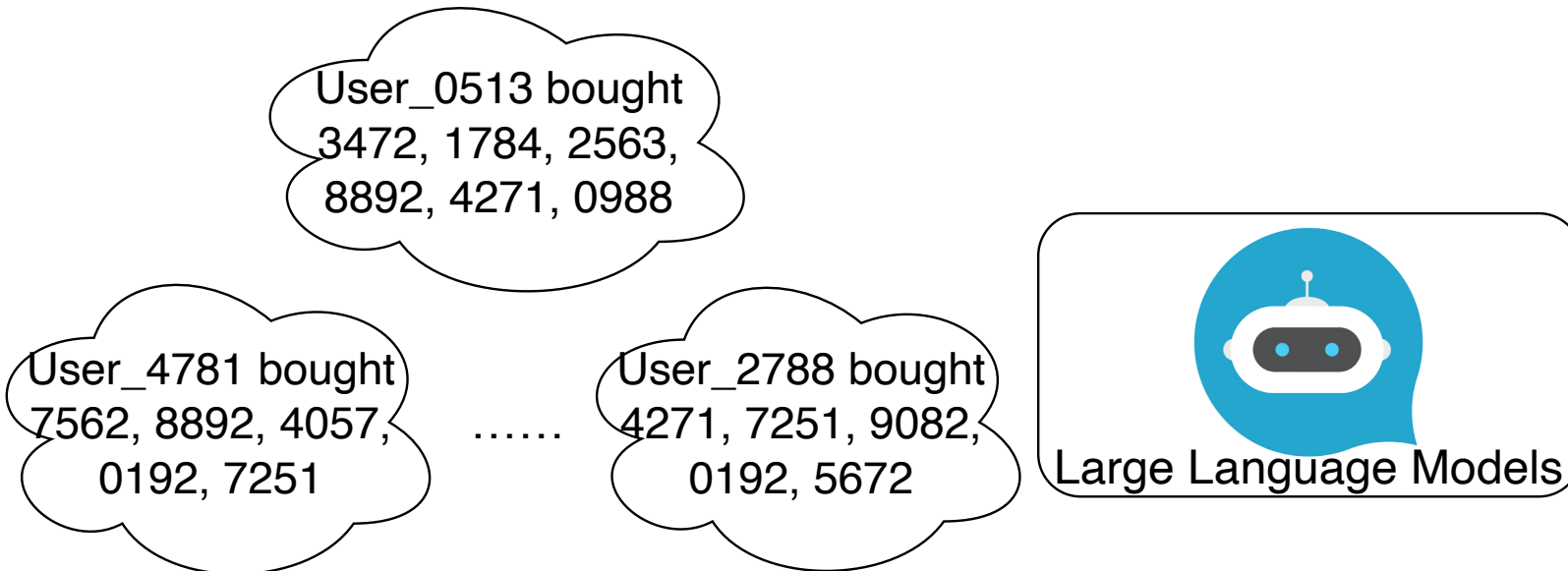
User\_4782 will buy 7251.



# ID-based LLM RecSys



- Modelling user interaction history with Markov chain



Pre-training & Fine-tuning

$$v_{i+1} = \arg \max_v P(v_{i+1} | v_1, v_2, \dots, v_i)$$



User\_4782 will buy 7251.

Modelling the probability of the next item





# ID-based LLM RecSys



## □ The N-gram probability in NLP

### ❖ Unigram

$$P("3472") = \frac{1}{16}$$

$$P("2563") = \frac{1}{16}$$

$$P("4271") = \frac{2}{16}$$

$$P("7562") = \frac{1}{16}$$

$$P("0192") = \frac{2}{16}$$

$$P("9082") = \frac{1}{16}$$

$$P("1784") = \frac{1}{16}$$

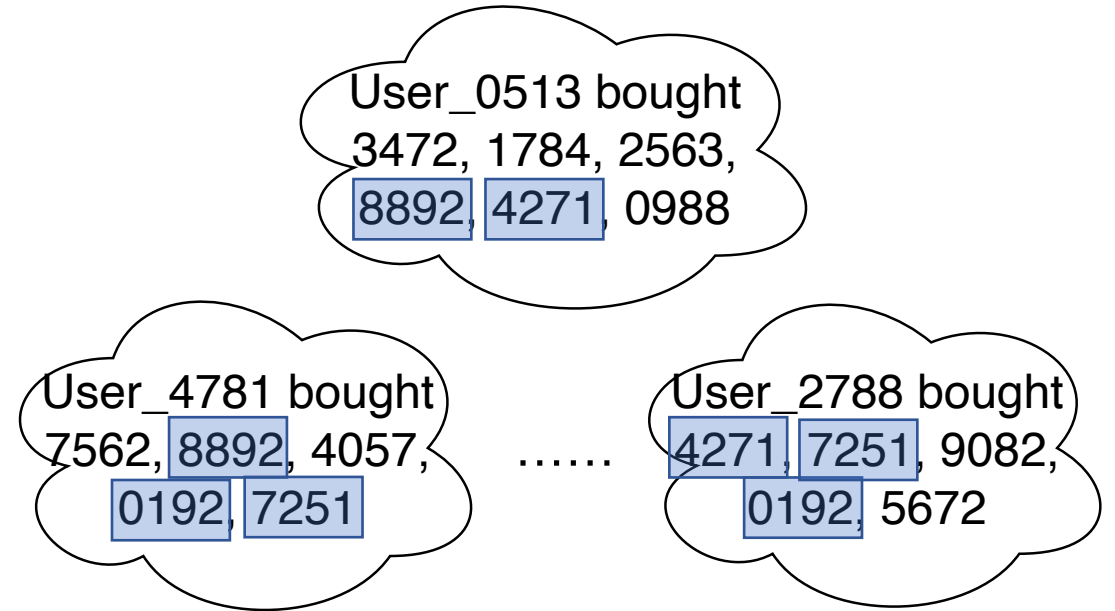
$$P("8892") = \frac{2}{16}$$

$$P("0988") = \frac{1}{16}$$

$$P("4057") = \frac{1}{16}$$

$$P("7251") = \frac{2}{16}$$

$$P("5672") = \frac{1}{16}$$



User\_4782 has bought 3472, 7653, 0192, 4271. What will he buy next?



Question



# ID-based LLM RecSys



□ The N-gram probability in NLP

❖ Bigram

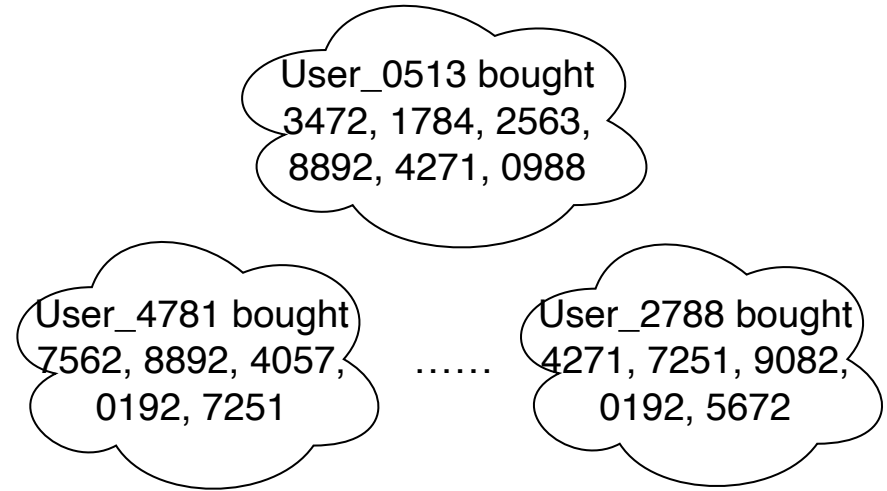
$$P("0988" \mid "4271") = \frac{1}{2}$$

$$P("7251" \mid "4271") = \frac{1}{2}$$

❖ Which one to choose?

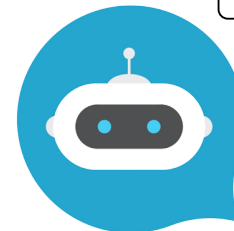
$$P("0988") = \frac{1}{16}$$

$$P("7251") = \frac{2}{16}$$



Question

User\_4782 has bought 3472, 7653, 0192, 4271. What will he buy next?



RecSys

User\_4782 will buy 7251.



# ID-based LLM RecSys

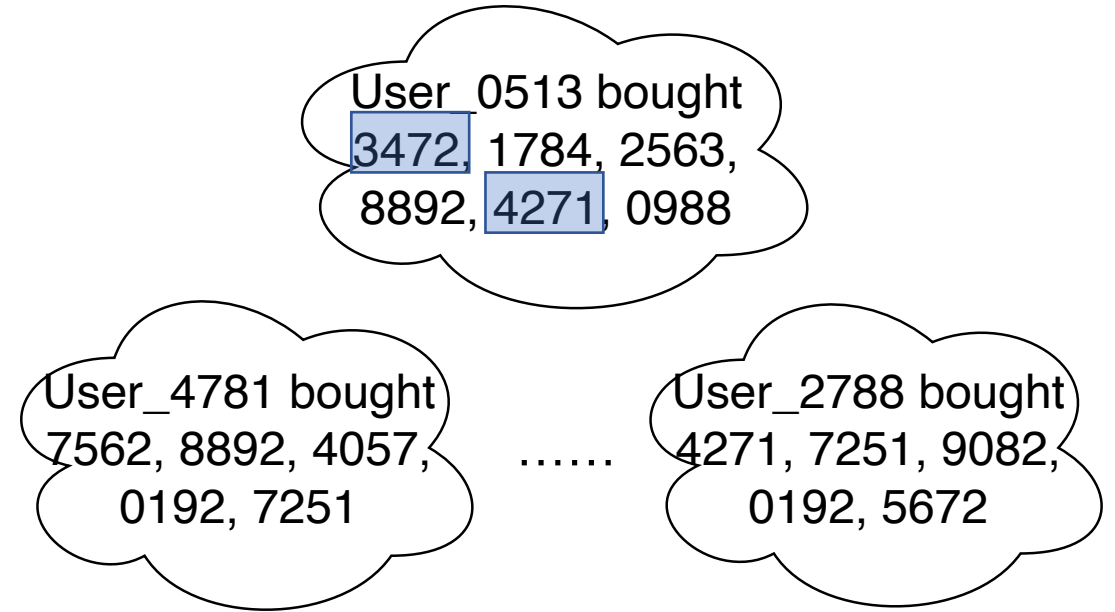


## ❑ The N-gram probability in NLP

- ❖ The co-occurrence of item IDs
- ❖ User\_0513 bought 3472, ..., 4271, 0988
- ❖ User\_4782 bought 3472, ..., 4271, ?



- ❖ Is "0988" a better answer than "7251"?



Question

User\_4782 has bought 3472, 7653, 0192, 4271. What will he buy next?



# ID-based LLM RecSys

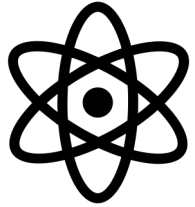


## Contextual representations of words in LLMs

- ❖ User\_0513 bought 3472, 1784, 2563, 8892, 4271, 0988
- ❖ User\_4782 bought 3472, 7653, 0192, 4271, ?
- ❖ The item representations can vary for different contexts



# User & Item Representation in LLMs



ID-based LLM RecSys



Text-based LLM RecSys



# Text-based LLM RecSys



## ❑ GPT4Rec

- ❖ Item **title** contains rich semantic information
- ❖ It's a natural way to use **text** to describe items

Previously, the customer has bought

Ben Nye Banana Luxury Face Powder 3.0 oz Makeup Kim Kardashian NEW!!!.  
Rosallini Women Stainless Steel Extension Eyelash Applicator Tool Fish Tail Clip.  
Beauty Flawless Makeup Blender Sponge Puff (size 1). Fruit Of The Earth 100%  
Aloe Vera 24oz Gel Pump.

In the future, the customer wants to buy

Fine-tuned GPT-2



Ben Nye Luxury Powders - Banana 1.5oz.  
Beautyblender Solid Blendercleanser 1 oz.  
Professional 15 Color Concealer Camouflage Makeup Palette.  
Pro Beauty Makeup Sponge Blender Flawless Smooth Shaped Water Droplets Puff (Random Color).  
L'Oreal Paris True Match Super Blendable Makeup, Natural Buff, 1.0 Ounces.



GPT4Rec

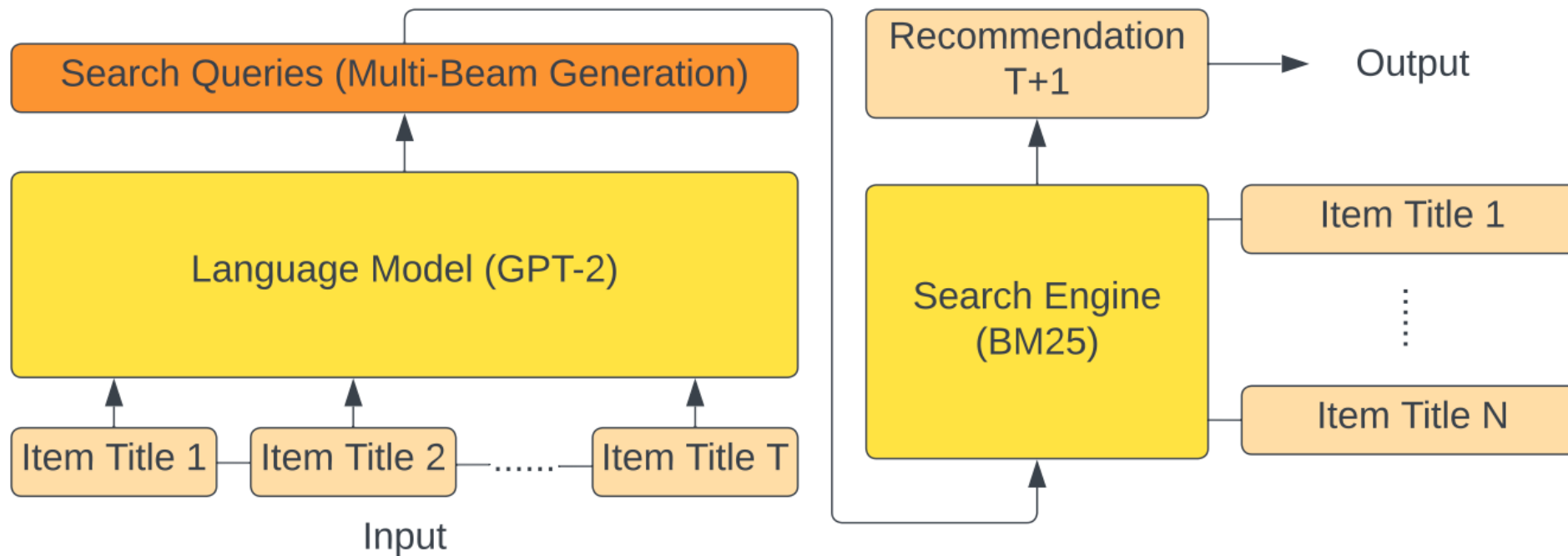


# Text-based LLM RecSys



## □ GPT<sub>4</sub>Rec

- ❖ In the era of LLMs, **Retrieval-Augmented Generation (RAG)** could be a way to improve the capability of LLMs
- ❖ RAG also enhances the explainability of LLM-based RecSys



# Text-based LLM RecSys



## □ TF-DCon

- ❖ Content-level condensation for recommendation
- ❖ Condense Item title and description to refine item representation

Enhance item titles based on given contents in the following format:

[title] {title}, [abstract] {abstract}, [category] {category}

You should rephrase the title to be clear, complete, objective, and neutral. Only provide the new title in the following format:

[newtitle] {newtitle}



[title] {Health Weightloss Watch},

[abstract] {Man Shares Time-Lapse Video of Six-Month Weight-Loss Journey We're big fans of weight-loss stories, but we usually only get to see the before and after photos. Very rarely do we get to see someone's physique transform right before our very eyes.},

[category] {Health}



[newtitle] {A Six-Month Weight-Loss Journey Captured in Time-Lapse Video},



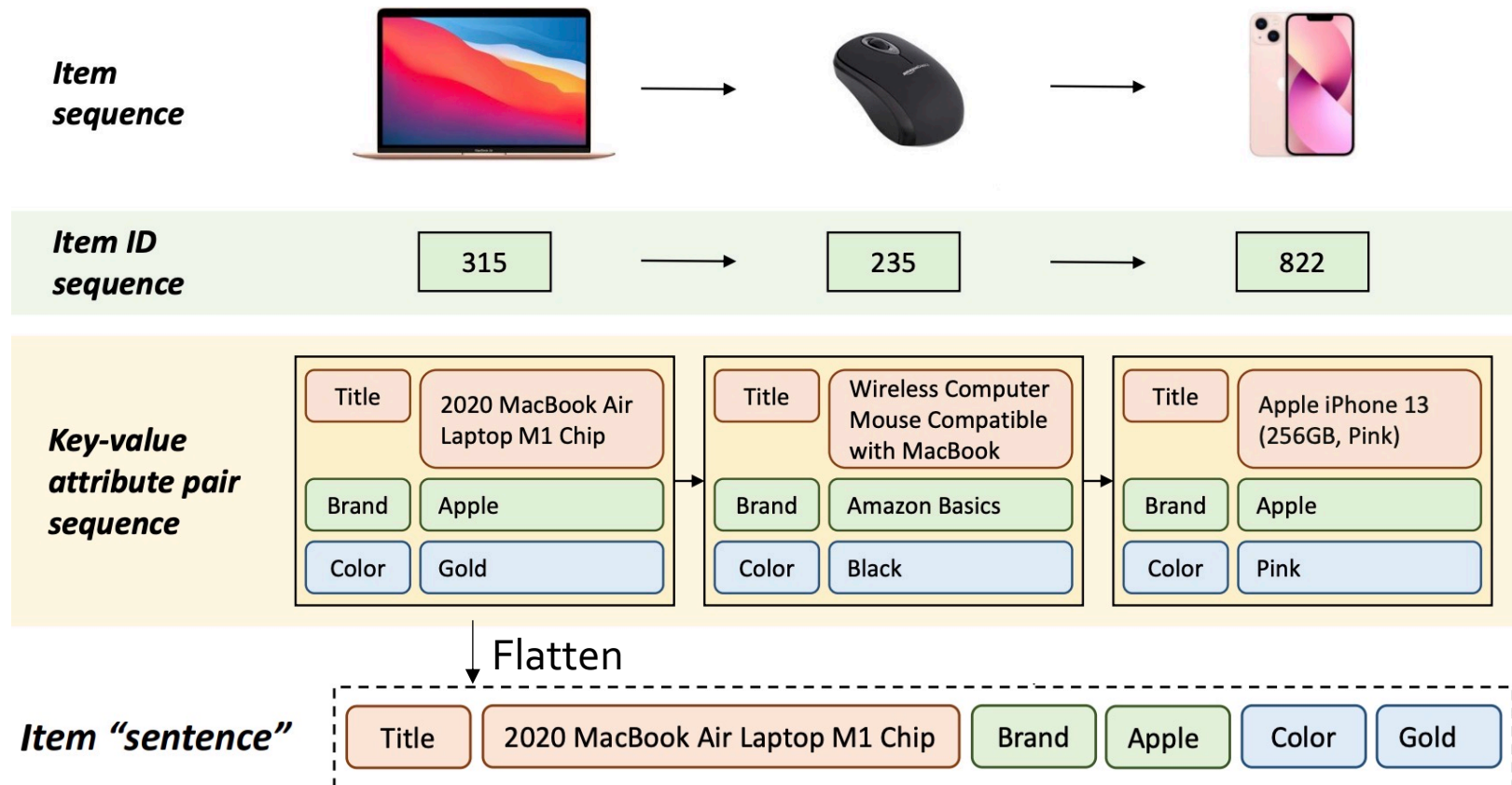


# Text-based LLM RecSys



## □ Recformer

- ❖ Use key-value attribute pairs to represent items



# Tutorial Outline

- ⦿ **Part 1: Introduction** of RecSys in the era of LLMs (Dr. Wenqi Fan)
- ⦿ **Part 2: Preliminaries** of RecSys and LLMs (Dr. Yujuan Ding)
- ⦿ **Part 3: Pre-training paradigms for adopting LLMs to RecSys (Dr. Yujuan Ding)**
- **Part 4: Fine-tuning** paradigms for adopting LLMs to RecSys (Liangbo Ning)
- **Part 5: Prompting** paradigms for adopting LLMs to RecSys (Shijie Wang)
- **Part 6: Future directions** of LLM-empowered RecSys (Dr. Wenqi Fan)

Website of this tutorial  
Check out the slides and more information!



# PART 3: RecSys Pre-training



**Presenter**  
**Dr. Yujuan DING**  
**HK PolyU**

- **Pre-training in NLP**
  - What is pre-training?
  - Why is pre-training needed?
  - NLP pre-training methods
- Pre-training LLM-based RecSys
  - What and why?
  - RecSys pre-training methods

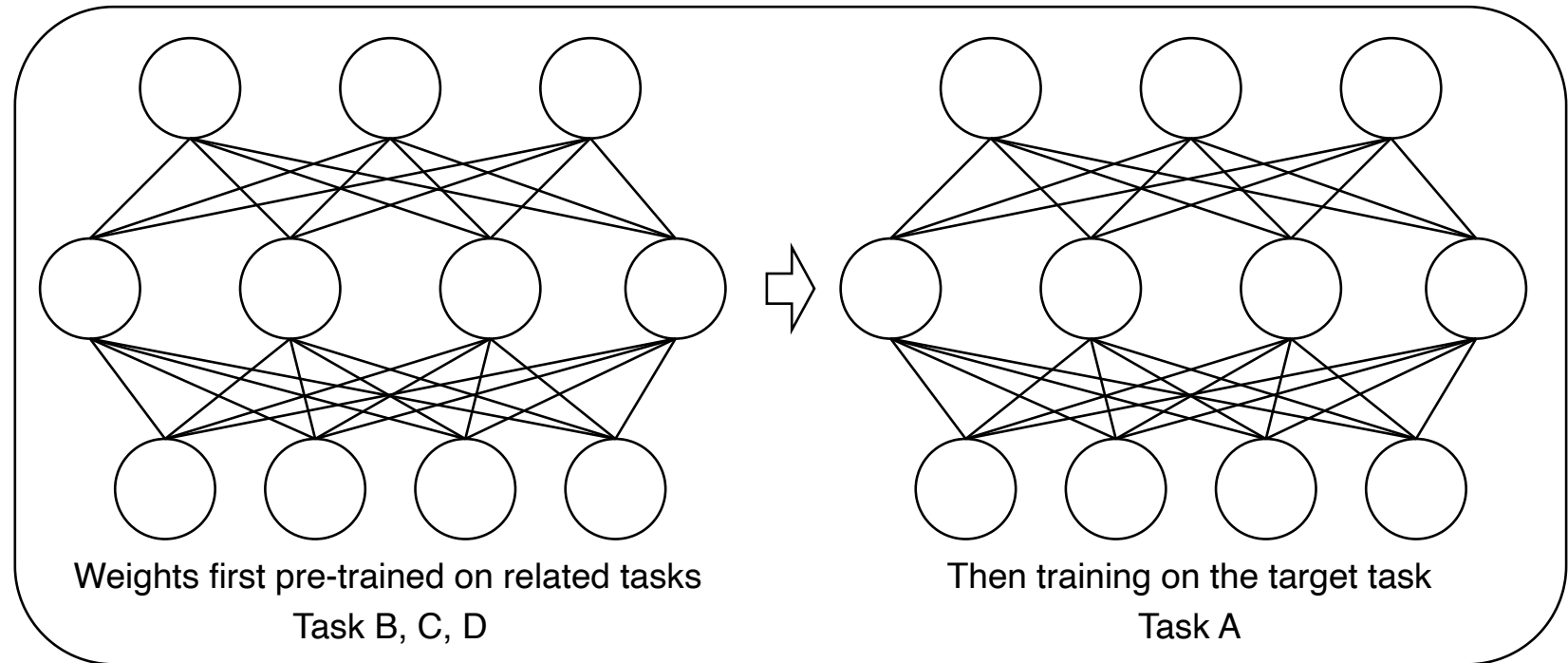
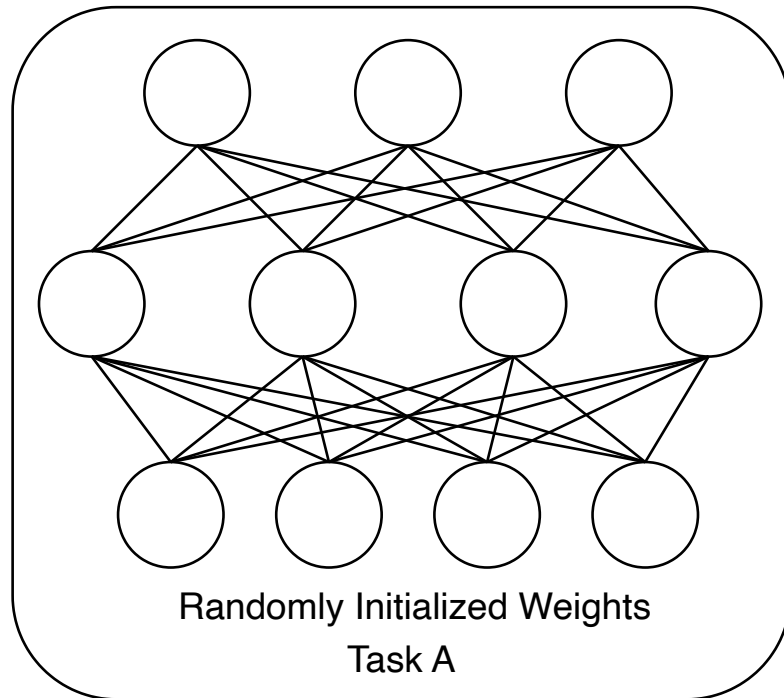


# Pre-training in NLP



## □ What is pre-training?

- ❖ Core Idea: knowledge transfer
- ❖ Technically

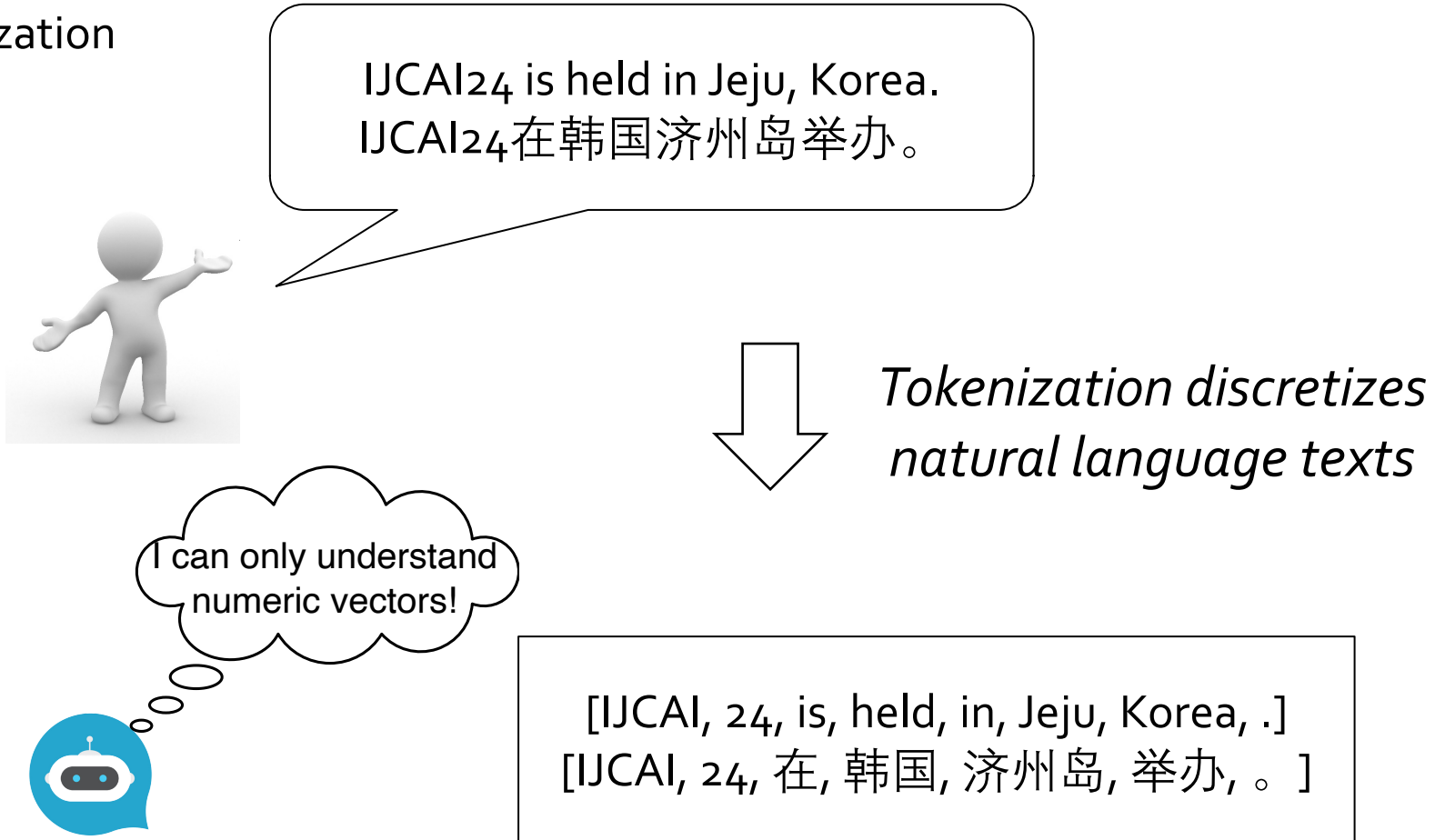


# Pre-training in NLP



## □ Why pre-training?

### ❖ Recall: Tokenization



# Pre-training in NLP



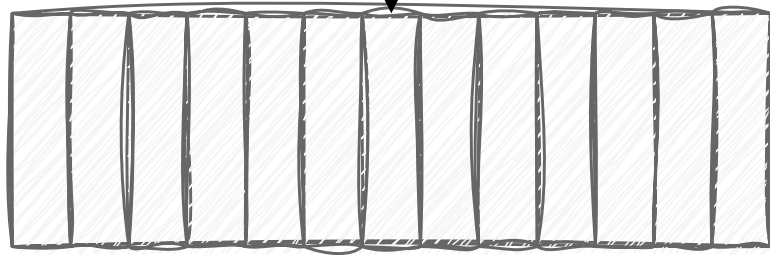
## □ Why pre-training?

❖ Recall: Tokenization

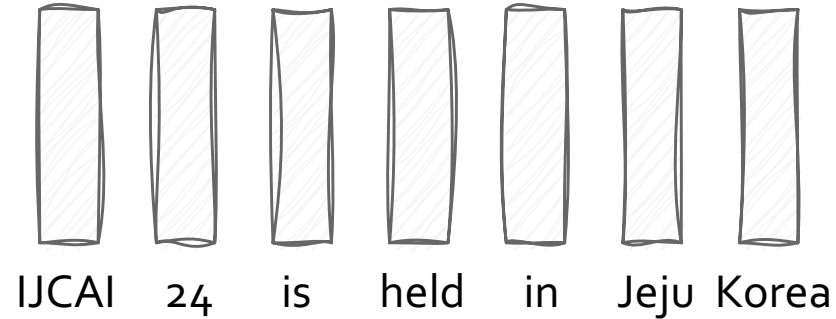
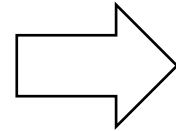
*Tokenized natural language texts are mapped to embedding vectors*

[IJCAI, 24, is, held, in, Jeju, Korea, .]  
[IJCAI, 24, 在, 韩国, 济州岛, 举办, 。]

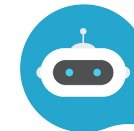
Look up



Embedding Matrix



How to initialize?

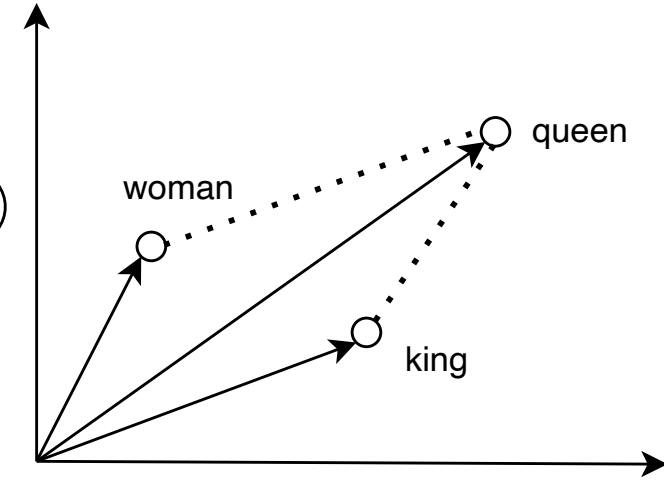


# Pre-training in NLP



## Word embeddings?

- ❖ king:  $[-0.5, -0.9, 1.4, \dots]$
- ❖ queen:  $[-0.6, -0.8, -0.2, \dots]$
- ❖ woman:  $[-0.1, -0.1, -1.6, \dots]$



Static word embeddings (word2vec, Glove) are pre-trained on text corpus from co-occurrence statistics.

- ❖ He is the king of the country



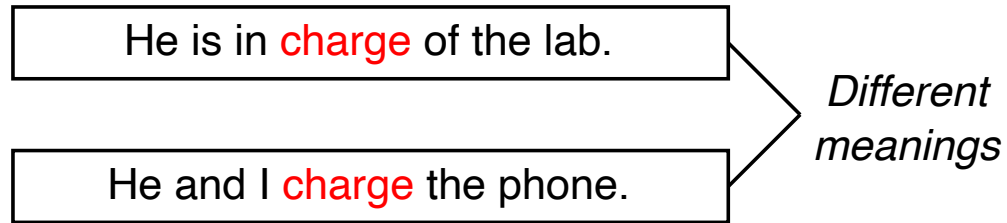
- ❖ She is the queen of the country



# Pre-training in NLP



- ❑ Problem of static word embedding – Context-Free



- ❑ How to solve it? – Contextual representations

- ❖ He is in **charge** of the lab
  - charge: [0.2, 0.8, 1.4, ...]
- ❖ He and I **charge** the phone
  - charge: [-0.3, -0.4, 0.7, ...]



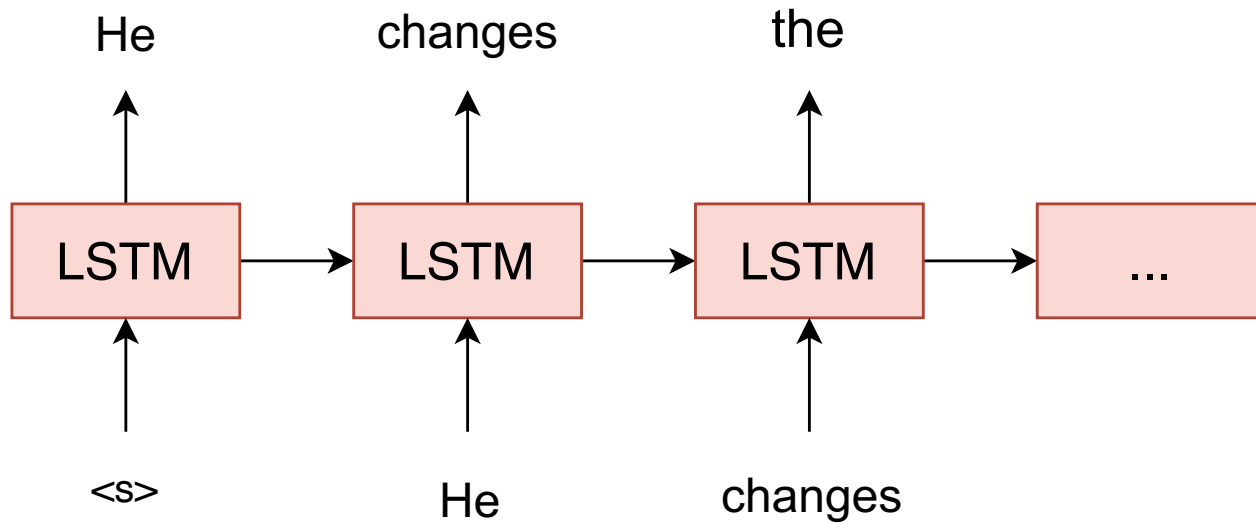


# Pre-training in NLP

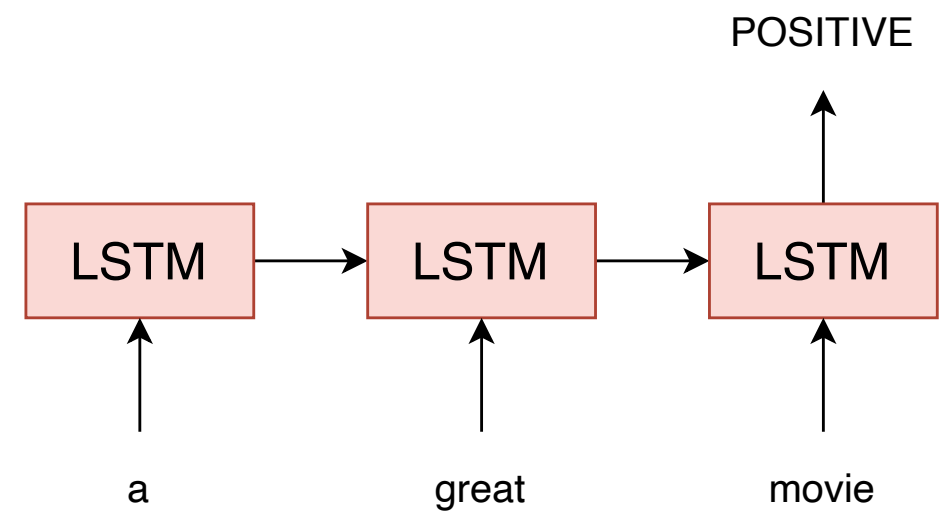


## □ Semi-Supervised Sequence Learning

**Training LSTM as Language Model**



**Fine-tuning on Sentiment Classification**

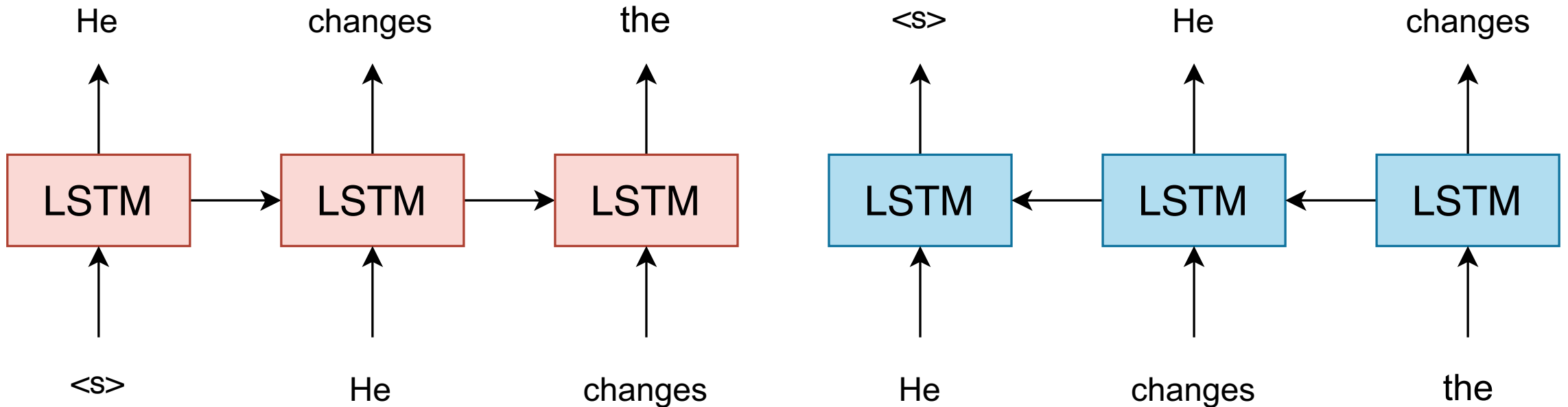


# Pre-training in NLP



- ELMo: Deep Contextual Word Embeddings

## Training Separate Left-to-Right and Right-to-Left Language Models



# Pre-training in NLP



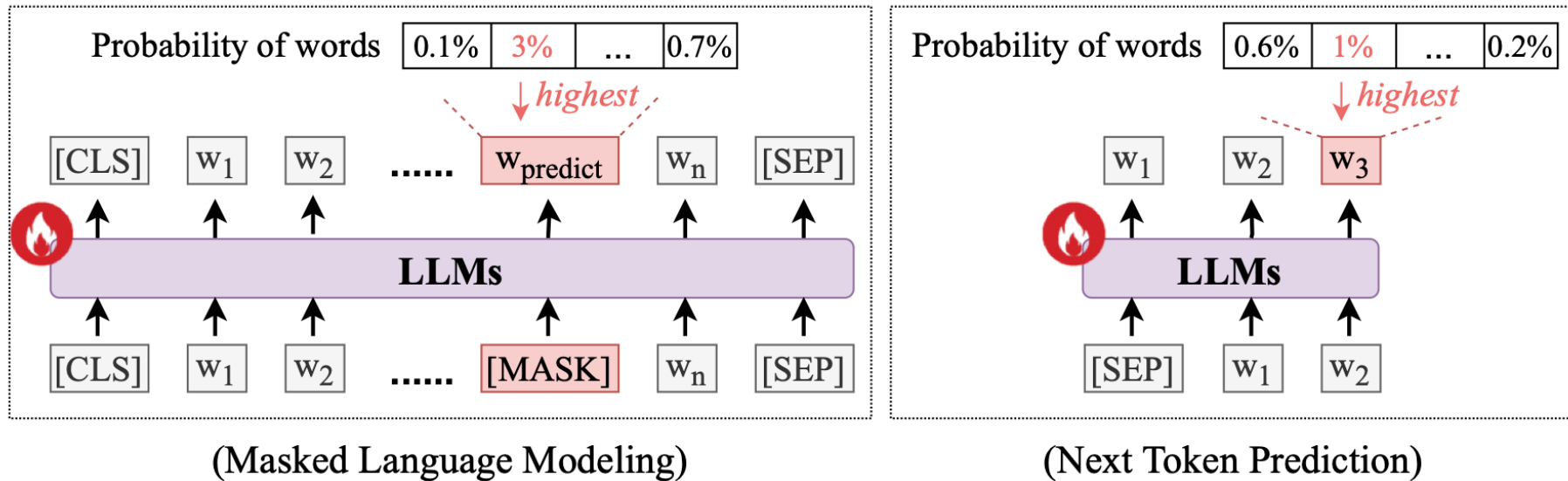
## Most Favored Pre-training Tasks in NLP

- ❖ Design specific pre-training tasks that could introduce knowledge
  - **Masked Language Modelling** (For Encoder-Decoder and Encoder-only Structures)
  - **Next Token Prediction** (For Decoder-only Structures)

### Pre-training



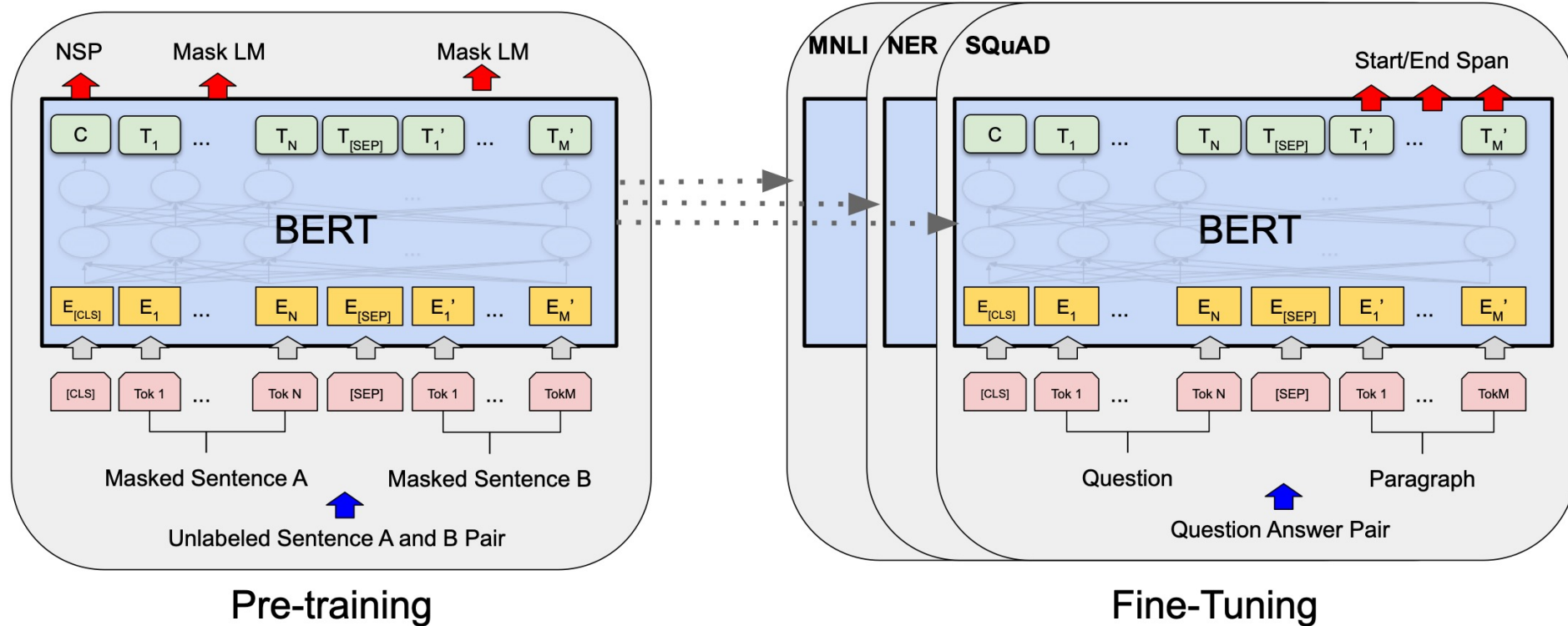
Large corpus  
unlabeled data



# Pre-training in NLP



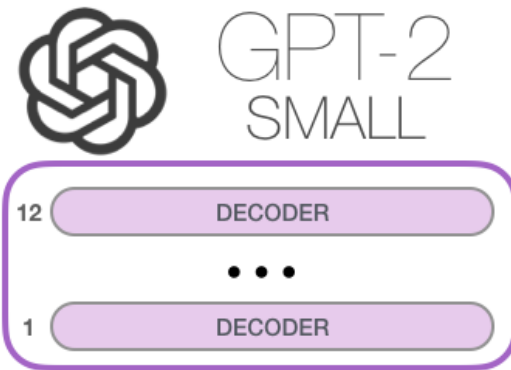
## □ BERT: Bidirectional Encoder Representations from Transformers



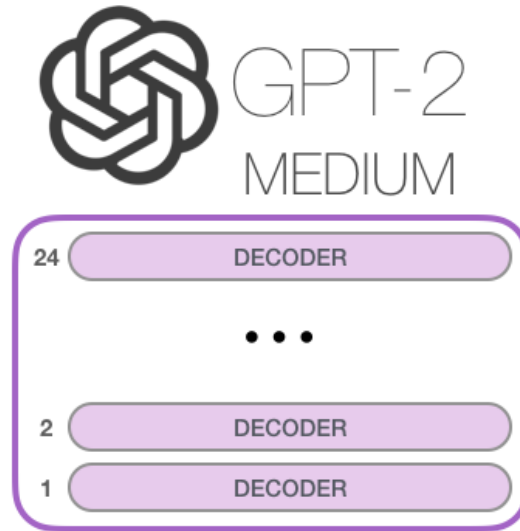
# Pre-training in NLP



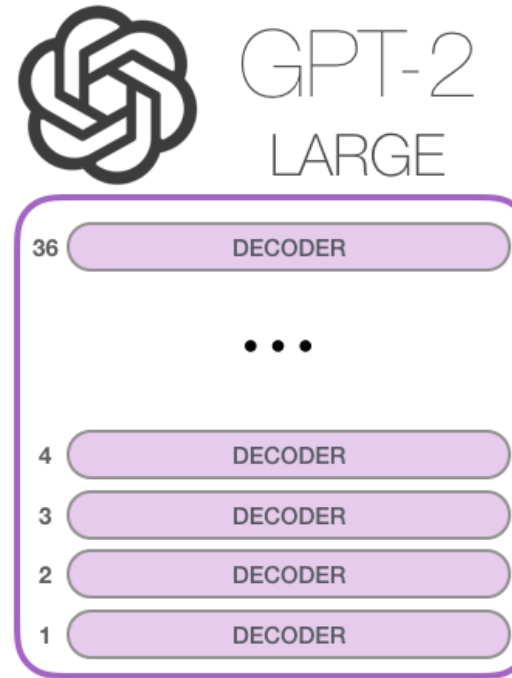
## □ GPT



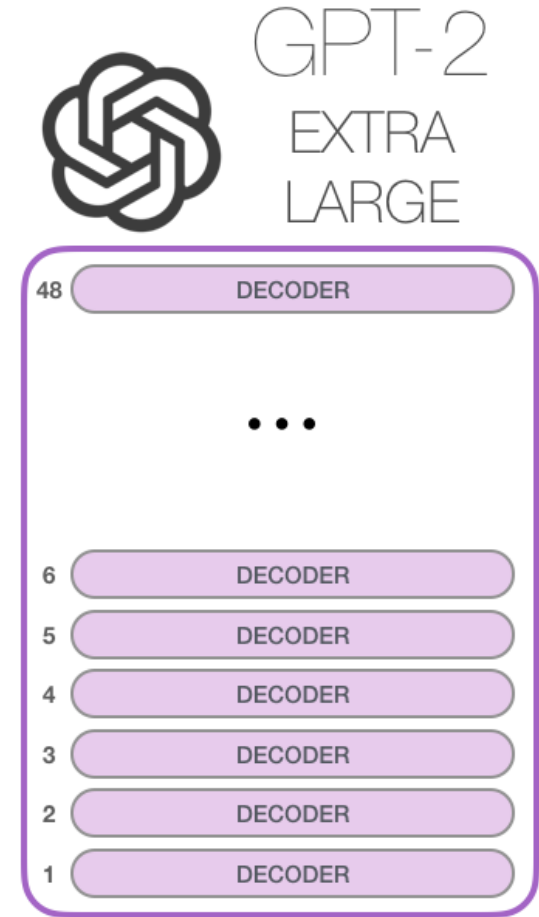
Model Dimensionality: 768



Model Dimensionality: 1024



Model Dimensionality: 1280



Model Dimensionality: 1600



# PART 3: RecSys Pre-training



Website of this tutorial

- ⦿ **Pre-training in NLP**
  - ⦿ What is pre-training?
  - ⦿ Why is pre-training needed?
  - ⦿ NLP pre-training methods
- **Pre-training LLM-based RecSys**
  - What and why?
  - RecSys pretraining methods

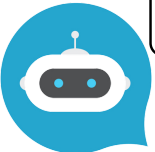


# Pre-training LLM-based RecSys




## □ What is Pre-training in LLM-based RecSys and Why is it Necessary?

- ❖ General pre-training vs. domain-specific pre-training
- ❖ Domain knowledge is essential for relieving the knowledge gap

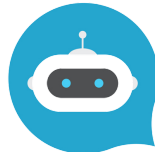


😊 I now know how to generate natural language like a human!

General Pre-training



User\_4782 has bought 3472, 7653, 0192, 4271.  
What will he buy next?

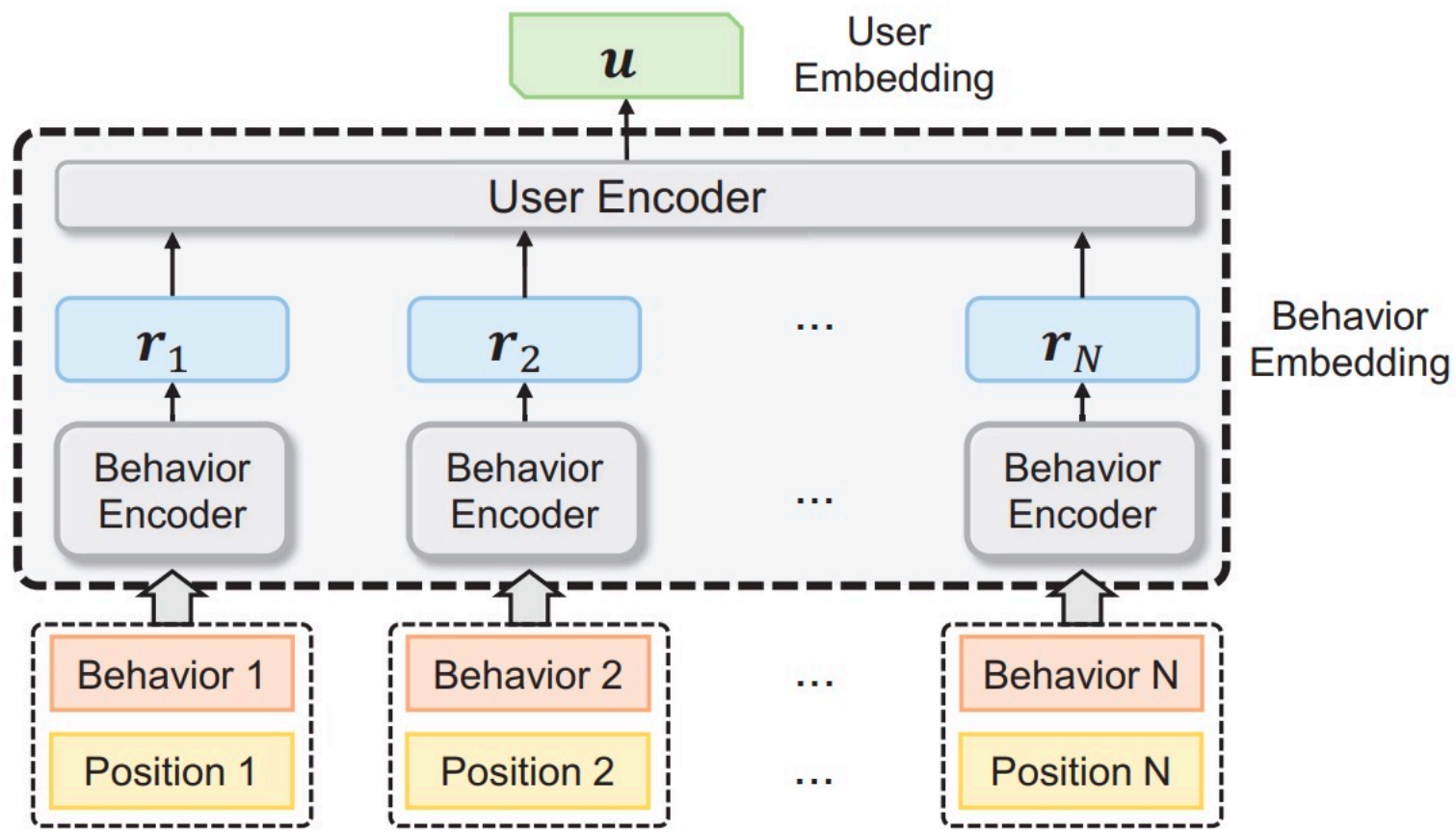


😓 I need more knowledge to make recommendations!

Recommendation-corpus is required



## Pre-training User Model from Unlabeled User Behaviors via Self-supervision

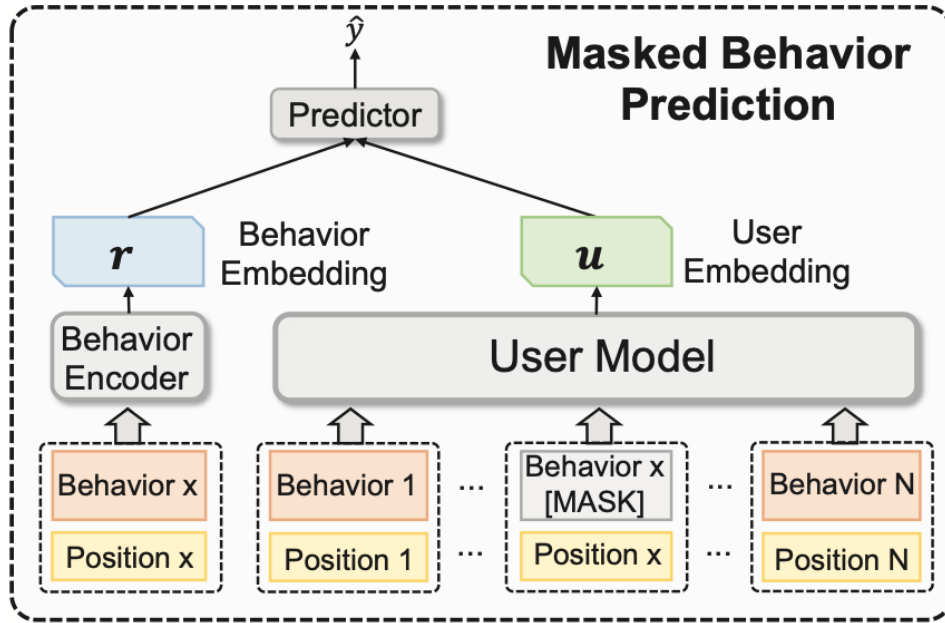




# PTUM pre-training tasks

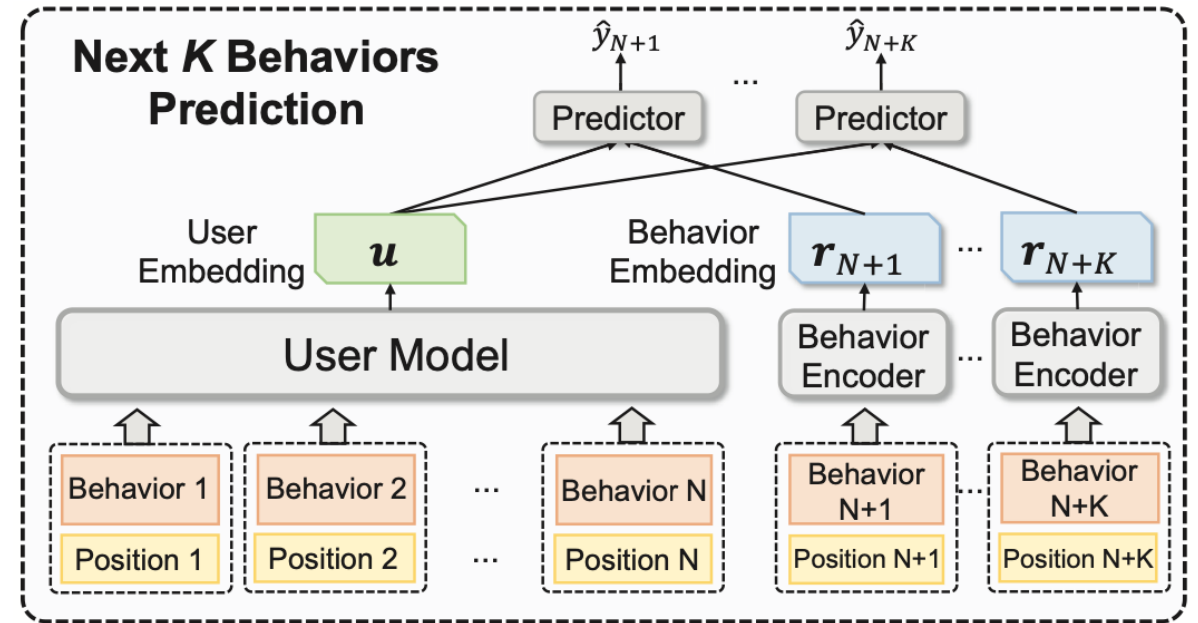


- ❑ Masked Behavior Prediction (MBP)
- ❑ Next K Behaviors Prediction (NBP)



(a) Masked Behavior Prediction (MBP) task.

$$\mathcal{L}_{MBP} = - \sum_{y \in \mathcal{S}_1} \sum_{i=1}^{P+1} y_i \log(\hat{y}_i)$$



(b) Next  $K$  Behaviors Prediction (NBP) task.

$$\mathcal{L}_{NBP} = - \frac{1}{K} \sum_{y \in \mathcal{S}_2} \sum_{k=1}^K \sum_{i=1}^{P+1} y_{i,k} \log(\hat{y}_{i,k})$$

$$\mathcal{L} = \mathcal{L}_{MBP} + \lambda \mathcal{L}_{NBP}$$



# PTUM application tasks



## □ Dataset

Demo			
# users	20,000	avg. # behaviors per user	224.7
# behaviors	4,494,771	avg. # words per webpage title	9.28
CTR			
# users	374,584	avg. # words per webpage title	10.23
# ads	4,159	avg. # words per ad title	11.95
# impressions	400,000	avg. # words per ad description	15.80
# clicked samples	364,281	# non-clicked samples	568,716
# users for pre-training	500,000	# behaviors for pre-training	63,178,293

## □ Ads CTR prediction

Methods	20%		50%		100%	
	AUC	AP	AUC	AP	AUC	AP
GRU4Rec	71.45	73.20	71.78	73.85	72.20	74.40
GRU4Rec+PTUM (no finetune)	71.76	73.66	71.95	74.15	72.33	74.77
GRU4Rec+PTUM (finetune)	72.33	74.55	72.42	74.72	72.79	75.40
NativeCTR	71.64	73.47	71.96	74.03	72.35	74.56
NativeCTR+PTUM (no finetune)	71.99	73.95	72.14	74.33	72.50	74.94
NativeCTR+PTUM (finetune)	72.52	74.79	72.59	74.91	72.91	75.57
BERT4Rec	71.82	73.97	72.39	74.89	72.99	75.45
BERT4Rec+PTUM (no finetune)	72.16	74.46	72.58	75.21	73.15	75.83
BERT4Rec+PTUM (finetune)	<b>72.74</b>	<b>75.34</b>	<b>73.03</b>	<b>75.81</b>	<b>73.59</b>	<b>76.48</b>

➤ Helpful across all datasets and setting scenarios

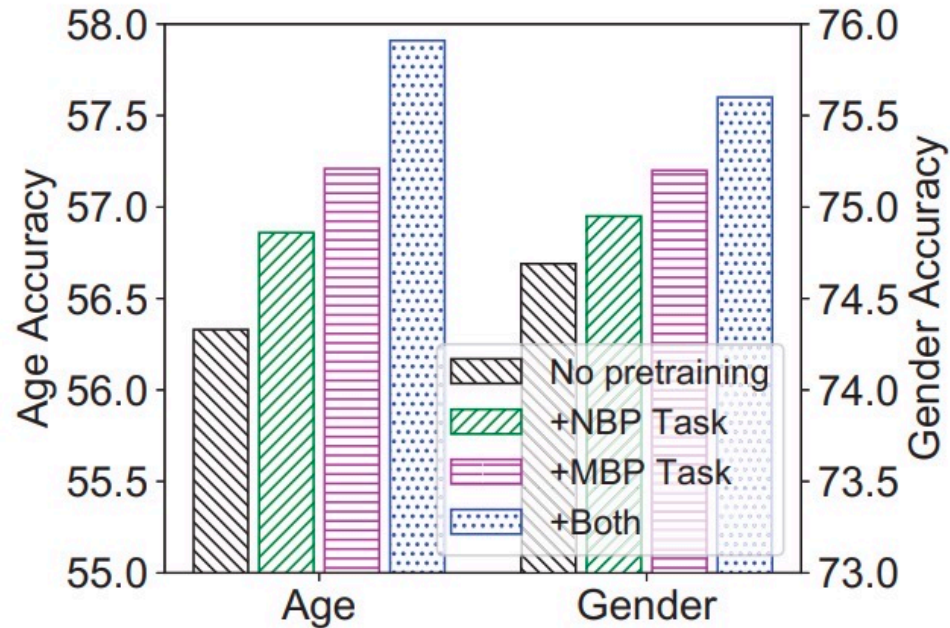


# PTUM effect of pre-training tasks

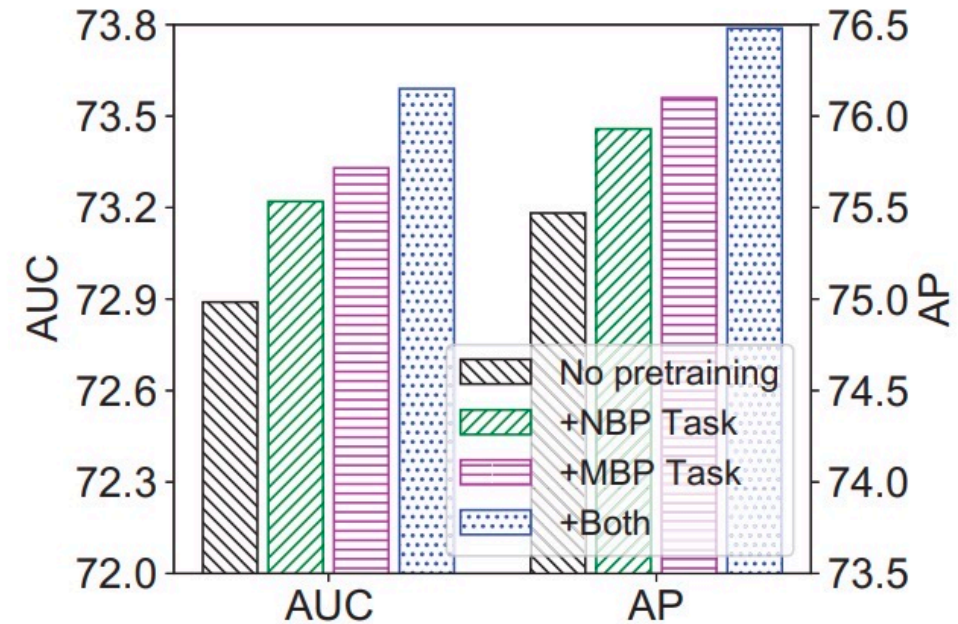


$$\mathcal{L}_{MBP} = - \sum_{y \in \mathcal{S}_1} \sum_{i=1}^{P+1} y_i \log(\hat{y}_i), \quad \mathcal{L}_{NBP} = - \frac{1}{K} \sum_{y \in \mathcal{S}_2} \sum_{k=1}^K \sum_{i=1}^{P+1} y_{i,k} \log(\hat{y}_{i,k})$$

$$\mathcal{L} = \mathcal{L}_{MBP} + \lambda \mathcal{L}_{NBP}$$



(a) *Demo* Dataset.



(b) *CTR* Dataset.





- ❑ Foundation recommendation model: one model to facilitate diverse domains and a myriad of tasks
- ❑ Challenges
  - ❖ the potentially unlimited set of downstream domains and tasks
  - ❖ the real-world systems' emphasis on computational efficiency
- ❑ Backbone: **M6**
  - ❖ is a series of visual-linguistic pretrained models
  - ❖ supports both Chinese and English
  - ❖ is a multi-modal model which aligns well with our plan to incorporate multi-modal features in the future, has achieved widespread success in Alibaba Group's ecosystem when deployed into real-world businesses



# M6-Rec: task convert



## □ Behavior Modeling as Language Modeling

- ❖ **Scoring tasks** (estimate the probability of a user clicking or purchasing an item)

[BOS'] December. Beijing, China. Cold weather. A male user in early twenties, searched "winter stuff" 23 minutes ago, clicked a product of category "jacket" named "men's lightweight warm winter hooded jacket" 19 minutes ago, clicked a product of category "sweatshirt" named "men's plus size sweatshirt stretchy pullover hoodies" 13 minutes ago, clicked . . . [EOS']

[BOS] The user is now recommended a product of category "boots" named "waterproof hiking shoes mens outdoor". The product has a high population-level CTR in the past 14 days, among the top 5%. The user clicked the category 4 times in the last 2 years. [EOS]



# M6-Rec: task convert



## ❑ Behavior Modeling as Language Modeling

- **Generation tasks** (personalized product design, explainable recommendation, personalized search query generation and conversational recommendation)

[BOS'] December. Beijing, China. Cold weather. A male user in early twenties, searched "winter stuff" 23 minutes ago, clicked a product of category "jacket" named "men's lightweight warm winter hooded jacket" 19 minutes ago, . . . [EOS'] [BOS] [EOS]

- **Zero-shot scoring tasks**

[BOS'] A user clicks hiking shoes [EOS'] [BOS] also clicks trekking poles [EOS] [BOS'] A user clicks hiking shoes [EOS'] [BOS] also clicks yoga knee pads [EOS]

- **Retrieval tasks**

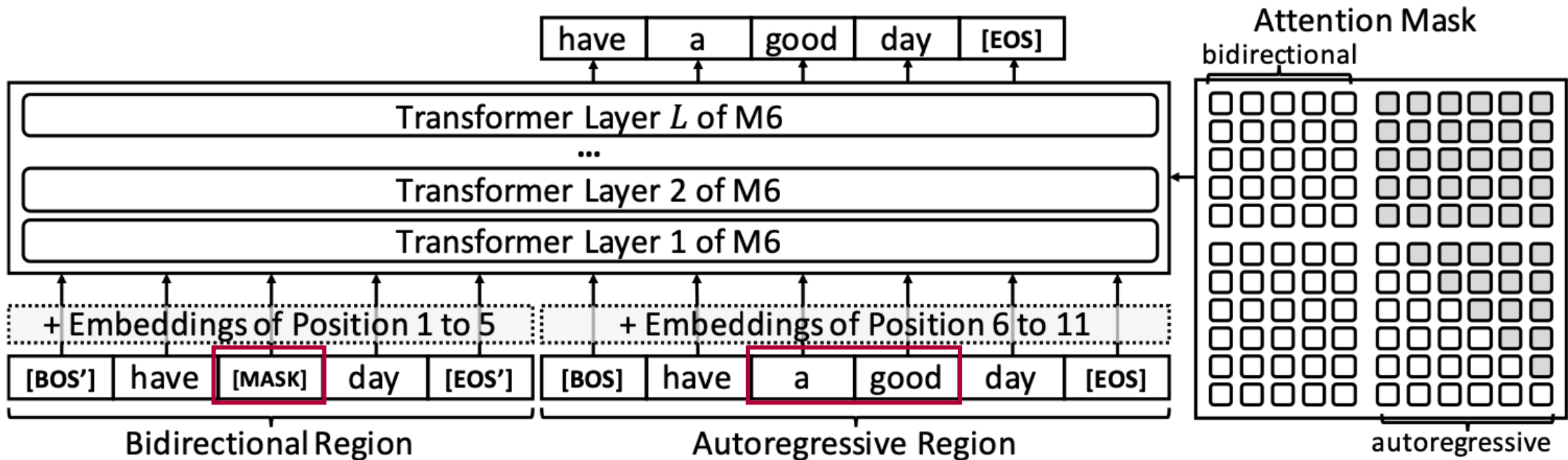
[BOS'] . . . [EOS'] [BOS] The user now purchases a product of category ". . ." named ". . ." . Product details: . . . The user likes it because . . . [EOS]



# M6-Rec: pre-training task



- ❑ Text infilling: masking small spans in a sentence
- ❑ Autoregressive language generation: masking the whole sentence



# M6-Rec: evaluation task results



## ❑ Click-through rate (CTR)

Method	Method Type	Datasets	
		AlipayQuery↑	TaoProduct↑
DIN	ID embeddings	0.7332	0.7611
M6-Rec	Text semantics	<b>0.7508</b>	<b>0.7995</b>

## ❑ kNN retrieval

Method	Method Type	Test Sets	
		All Items↑	Unseen Items↑
YouTubeDNN	ID embeddings	54.4%	<i>fail</i>
TwinBERT	Text semantics	69.6%	49.6%
M6-Rec	Text semantics	<b>74.1%</b>	<b>57.0%</b>

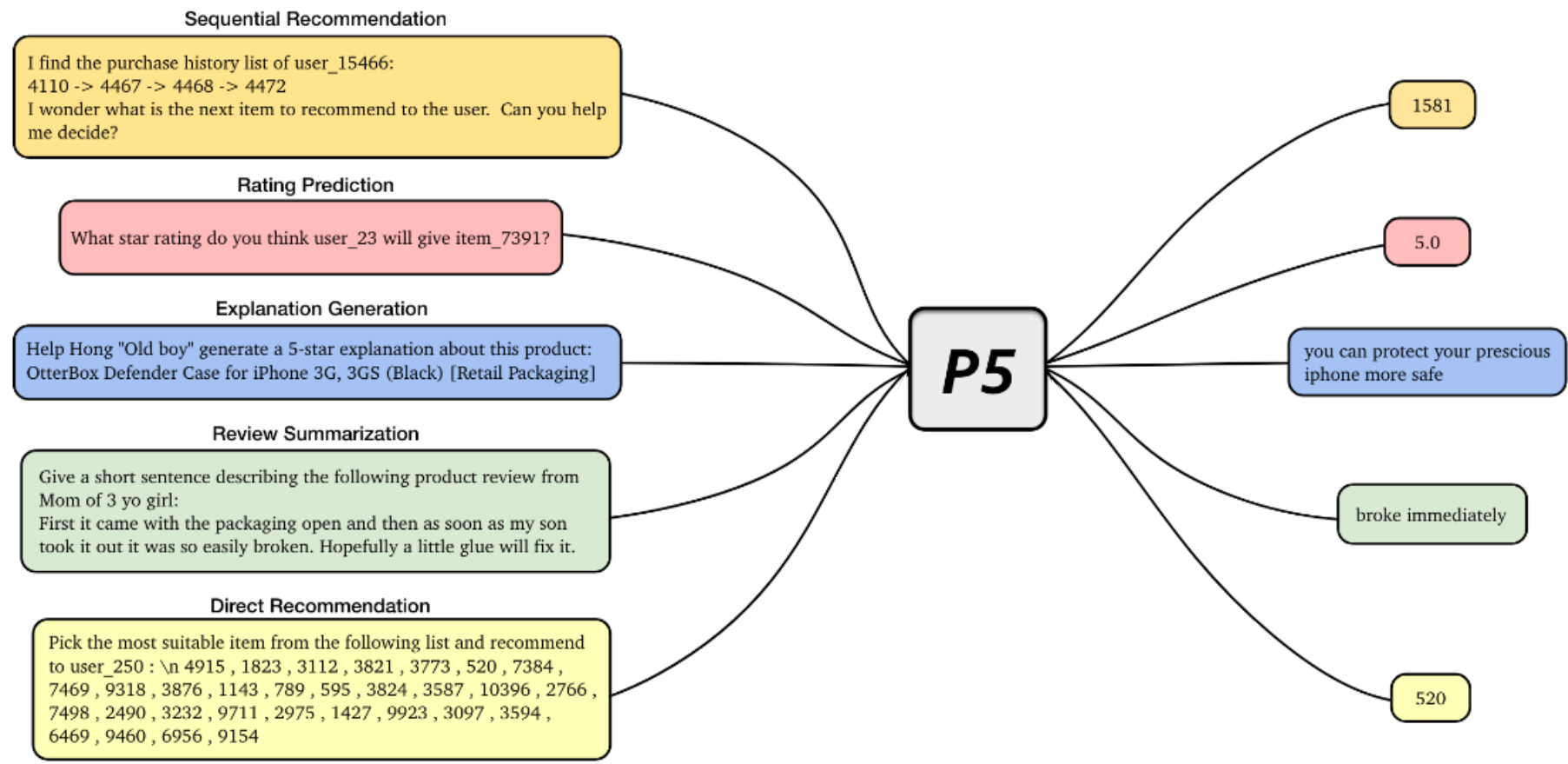
## ❑ Conversational recommendation

Method	Metrics				
	PPL↓	BLEU-2↑	BLEU-3↑	Dist-3↑	Dist-4↑
Transformer	20.44	0.026	0.014	0.27	0.39
KBRD [3]	17.90	0.060	0.024	0.30	0.45
KGSF [66]	10.73	0.033	0.022	0.40	0.46
M6-Rec	<b>10.25</b>	<b>0.122</b>	0.021	<b>0.46</b>	<b>0.64</b>





- ❑ Pretrain, Personalized Prompt, and Predict Paradigm
- ❑ Multi-task Pretraining with Personalized Prompt Collection



Multi-task Pretraining with Personalized Prompt Collection



# P5: task convert



## Rating / Review / Explanation raw data for *Beauty*

```
user_id: 7641      user_name: stephanie
item_id: 2051
item_title: SHANY Nail Art Set (24 Famouse Colors
Nail Art Polish, Nail Art Decoration)
review: Absolutely great product. I bought this for my fourteen year
old niece for Christmas and of course I had to try it out, then I
tried another one, and another one and another one. So much fun!
I even contemplated keeping a few for myself!
star_rating: 5
summary: Perfect!
explanation: Absolutely great product      feature_word: product
```

(a)

Which star rating will user\_{{user\_id}} give item\_{{item\_id}}?  
(1 being lowest and 5 being highest)

→ {{star\_rating}}

Based on the feature word {{feature\_word}}, generate an  
explanation for user\_{{user\_id}} about this product:  
{{item\_title}}

→ {{explanation}}

Give a short sentence describing the following product review  
from {{user\_name}}: {{review}}

→ {{summary}}



# P5: task convert



## Sequential Recommendation raw data for *Beauty*

```
user_id: 7641      user_name: Victor
purchase_history: 652 -> 460 -> 447 -> 653 -> 654 -> 655 -> 656 -> 8
-> 657
next_item: 552
candidate_items: 4885 , 4280 , 4886 , 1907 , 870 , 4281 , 4222 ,
4887 , 2892 , 4888 , 2879 , 3147 , 2195 , 3148 , 3179 , 1951 ,
..... , 1982 , 552 , 2754 , 2481 , 1916 , 2822 , 1325
```

(b)

Here is the purchase history of user\_{{user\_id}}:  
{{purchase\_history}}  
What to recommend next for the user?

→ {{next\_item}}

## Direct Recommendation raw data for *Beauty*

```
user_id: 250      user_name: moriah rose
target_item: 520
random_negative_item: 9711
candidate_items: 4915 , 1823 , 3112 , 3821 , 3773 , 520 , 7384 ,
7469 , 9318 , 3876 , 1143 , 789 , 595 , 3824 , 3587 , 10396 ,
..... , 2766 , 7498 , 2490 , 3232 , 9711 , 2975 , 1405 , 8051
```

(c)

Choose the best item from the candidates to recommend for  
{{user\_name}}? \n {{candidate\_items}}

→ {{target\_item}}

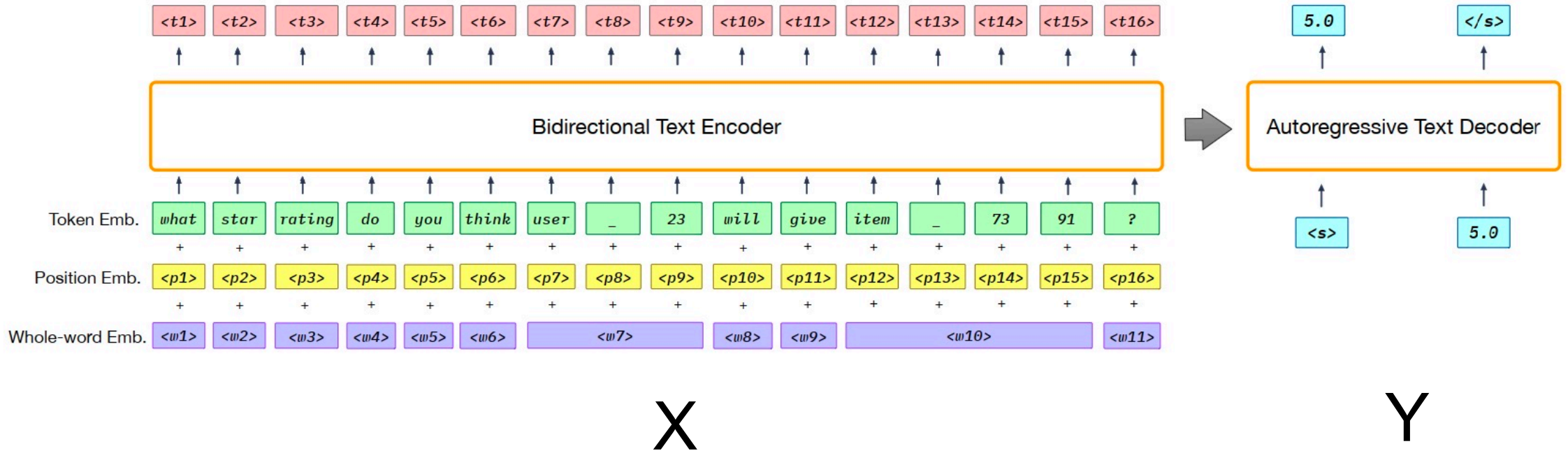


# P5: pre-training task



- Label token prediction

$$\mathcal{L}_\theta^{P5} = - \sum_{j=1}^{|y|} \log P_\theta (y_j | y_{<j}, \mathbf{x})$$

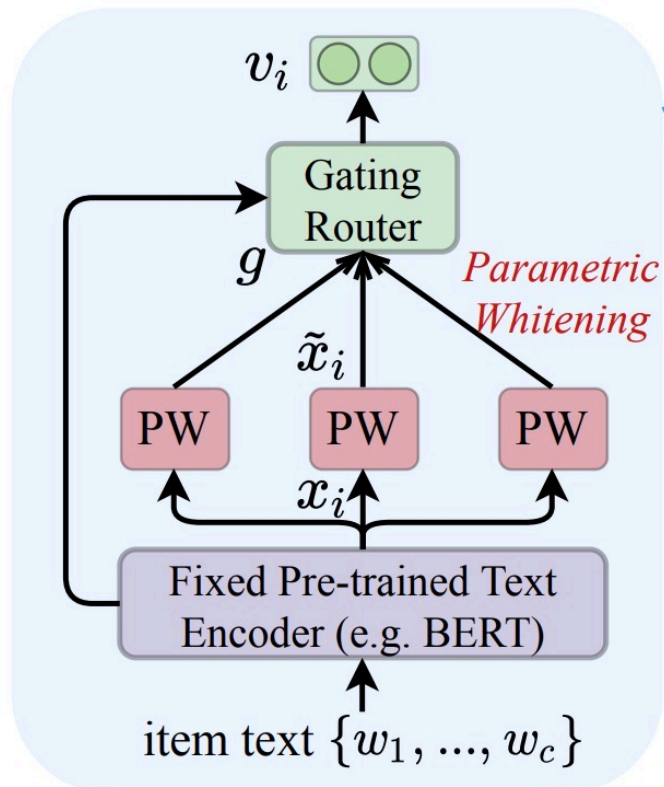


# UniSRec: universal sequence representation learning



- Utilizing the **associated description text** of items to learn **transferable** representations across different recommendation scenarios.

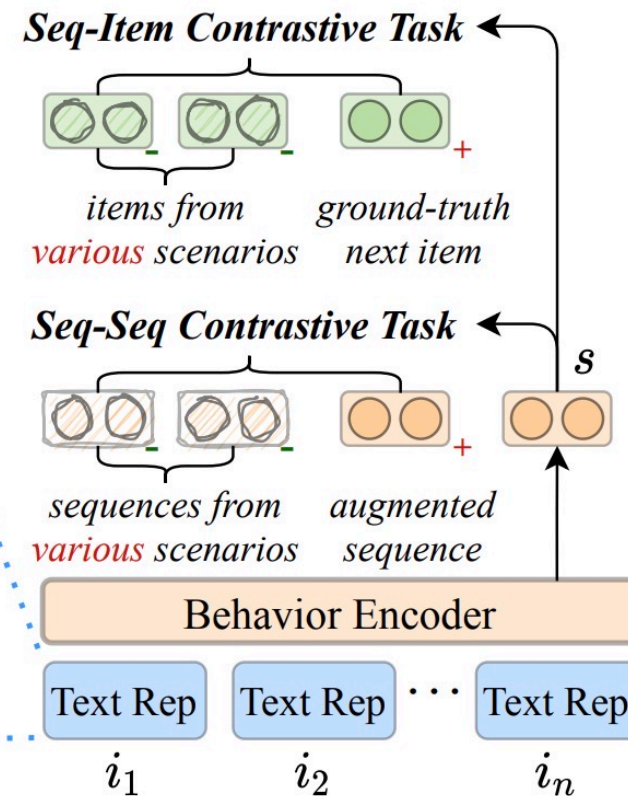
## Universal Item Representation



*MoE-enhanced Adaptor*

Mixture-of-Expert

## Universal Sequence Representation Pre-training



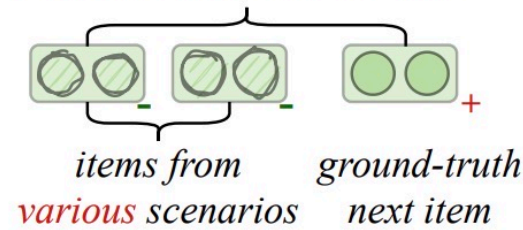
# UniSRec: Pre-training task



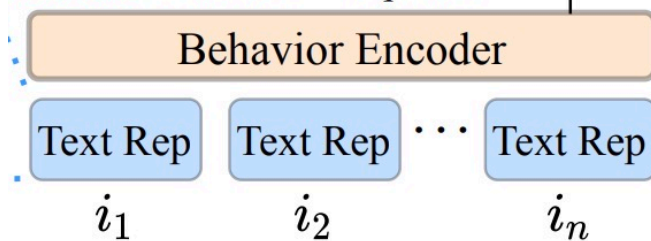
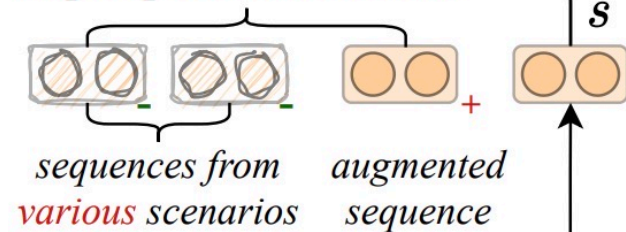
- ❑ Sequence-item contrastive learning
- ❑ Sequence-sequence contrastive learning

## Universal Sequence Representation Pre-training

### Seq-Item Contrastive Task



### Seq-Seq Contrastive Task



$$\ell_{S-I} = - \sum_{j=1}^B \log \frac{\exp(s_j \cdot v_j / \tau)}{\sum_{j'=1}^B \exp(s_j \cdot v_{j'} / \tau)}$$

$$\ell_{S-S} = - \sum_{j=1}^B \log \frac{\exp(s_j \cdot \tilde{s}_j / \tau)}{\sum_{j'=1}^B \exp(s_j \cdot s_{j'} / \tau)}$$

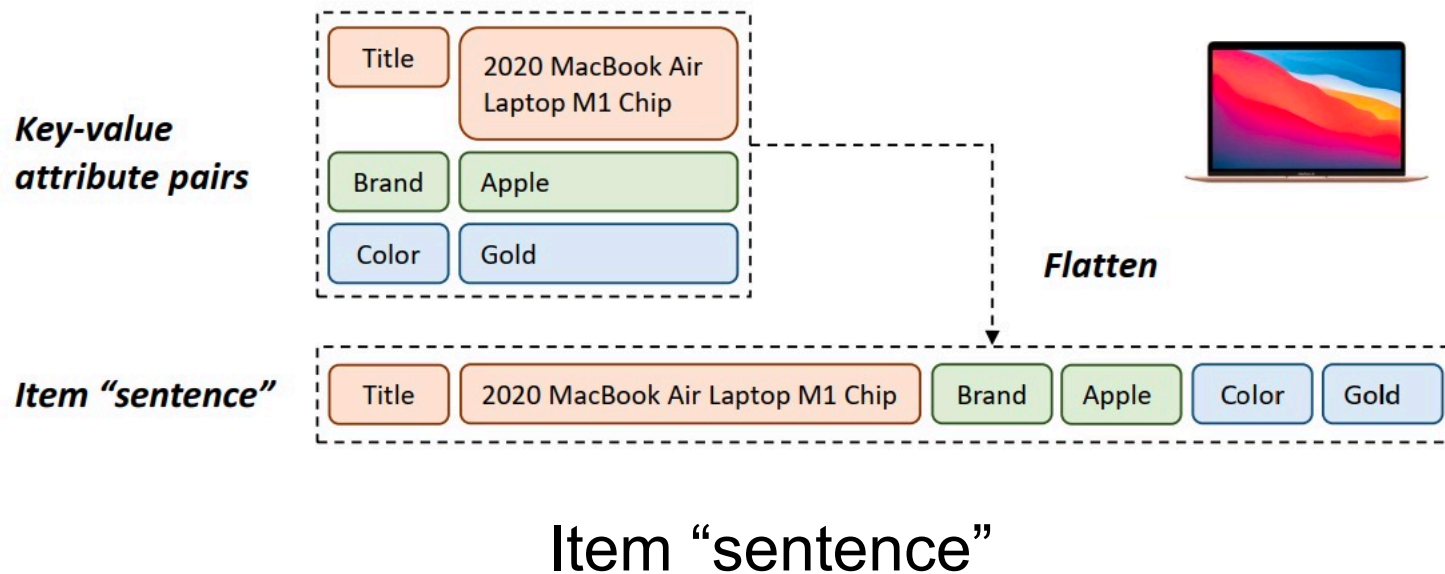
$$\mathcal{L}_{PT} = \ell_{S-I} + \lambda \cdot \ell_{S-S}$$



# Recformer



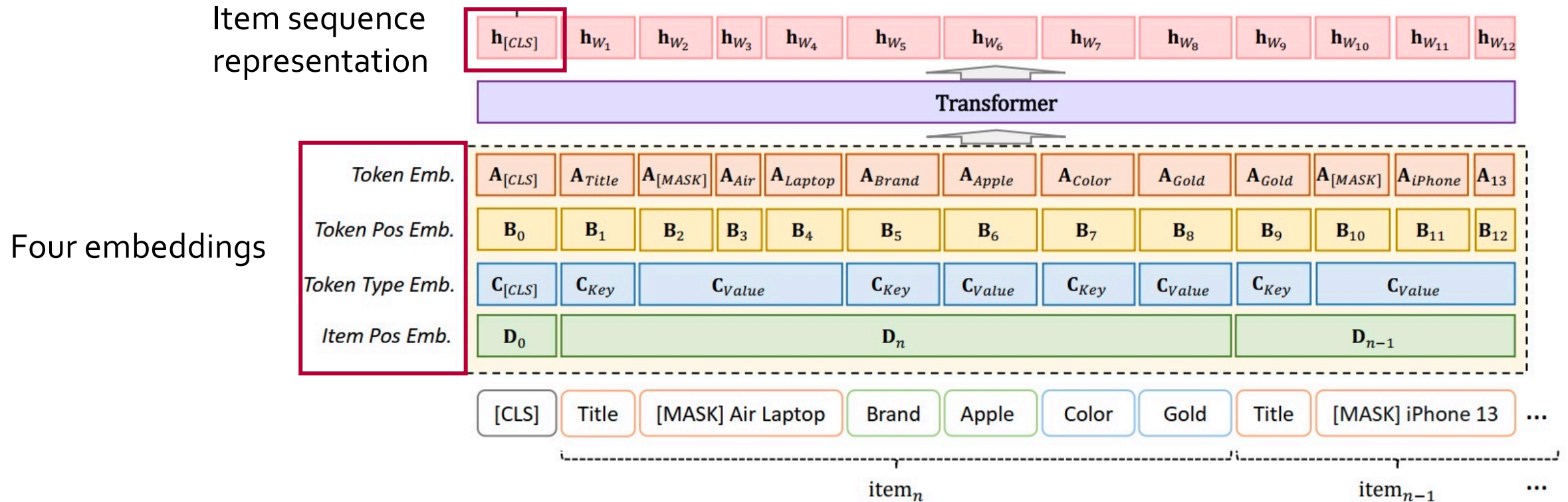
- Item → “sentence” (word sequence): **flattening item key-value attributes** described by text so that an item sequence for a user becomes a sequence of sentences.



# Recformer: model structure



- A similar structure as Longformer: a **multi-layer bidirectional Transformer** with an attention mechanism that scales linearly with sequence length.
- Considering computational efficiency, but also open to other bidirectional Transformer structures such as BERT.

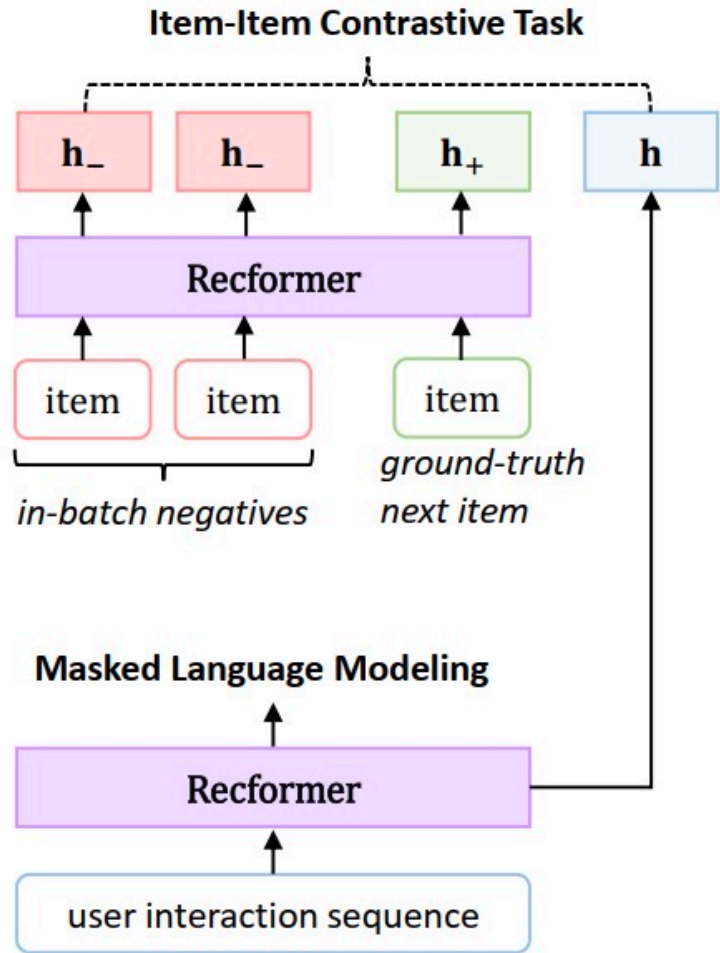




# Recformer: pre-training task



- ❑ Masked Language Modeling (MLM)
- ❑ Item-item contrastive task (IIC)



$$\mathcal{L}_{\text{IIC}} = -\log \frac{e^{\text{sim}(\mathbf{h}_s, \mathbf{h}_i^+) / \tau}}{\sum_{i \in \mathcal{B}} e^{\text{sim}(\mathbf{h}_s, \mathbf{h}_i) / \tau}}$$

$$\mathcal{L}_{\text{MLM}} = -\sum_{i=0}^{|\mathcal{V}|} y_i \log(p_i)$$

$$\mathcal{L}_{\text{PT}} = \mathcal{L}_{\text{IIC}} + \lambda \cdot \mathcal{L}_{\text{MLM}}$$



# Recformer: performance



- Different categories of Amazon review datasets

Dataset	Metric	ID-Only Methods				ID-Text Methods		Text-Only Methods			Improv.
		GRU4Rec	SASRec	BERT4Rec	RecGURU	FDSA	S <sup>3</sup> -Rec	ZESRec	UniSRec	RECFORMER	
Scientific	NDCG@10	0.0826	0.0797	0.0790	0.0575	0.0716	0.0451	0.0843	<u>0.0862</u>	<b>0.1027</b>	19.14%
	Recall@10	0.1055	<u>0.1305</u>	0.1061	0.0781	0.0967	0.0804	0.1260	0.1255	<b>0.1448</b>	10.96%
	MRR	0.0702	0.0696	0.0759	0.0566	0.0692	0.0392	0.0745	<u>0.0786</u>	<b>0.0951</b>	20.99%
Instruments	NDCG@10	0.0633	0.0634	0.0707	0.0468	0.0731	<u>0.0797</u>	0.0694	0.0785	<b>0.0830</b>	4.14%
	Recall@10	0.0969	0.0995	0.0972	0.0617	0.1006	<u>0.1110</u>	0.1078	<b>0.1119</b>	0.1052	-
	MRR	0.0707	0.0577	0.0677	0.0460	0.0748	<u>0.0755</u>	0.0633	0.0740	<b>0.0807</b>	6.89%
Arts	NDCG@10	<u>0.1075</u>	0.0848	0.0942	0.0525	0.0994	0.1026	0.0970	0.0894	<b>0.1252</b>	16.47%
	Recall@10	0.1317	0.1342	0.1236	0.0742	0.1209	<u>0.1399</u>	0.1349	0.1333	<b>0.1614</b>	15.37%
	MRR	0.1041	0.0742	0.0899	0.0488	0.0941	<u>0.1057</u>	0.0870	0.0798	<b>0.1189</b>	12.49%
Office	NDCG@10	0.0761	0.0832	<u>0.0972</u>	0.0500	0.0922	0.0911	0.0865	0.0919	<b>0.1141</b>	17.39%
	Recall@10	0.1053	0.1196	0.1205	0.0647	<u>0.1285</u>	0.1186	0.1199	0.1262	<b>0.1403</b>	9.18%
	MRR	0.0731	0.0751	0.0932	0.0483	<u>0.0972</u>	0.0957	0.0797	0.0848	<b>0.1089</b>	12.04%
Games	NDCG@10	0.0586	0.0547	<u>0.0628</u>	0.0386	0.0600	0.0532	0.0530	0.0580	<b>0.0684</b>	8.92%
	Recall@10	0.0988	0.0953	<u>0.1029</u>	0.0479	0.0931	0.0879	0.0844	0.0923	<b>0.1039</b>	0.97%
	MRR	0.0539	0.0505	<u>0.0585</u>	0.0396	0.0546	0.0500	0.0505	0.0552	<b>0.0650</b>	11.11%
Pet	NDCG@10	0.0648	0.0569	0.0602	0.0366	0.0673	0.0742	<u>0.0754</u>	0.0702	<b>0.0972</b>	28.91%
	Recall@10	0.0781	0.0881	0.0765	0.0415	0.0949	<u>0.1039</u>	0.1018	0.0933	<b>0.1162</b>	11.84%
	MRR	0.0632	0.0507	0.0585	0.0371	0.0650	<u>0.0710</u>	0.0706	0.0650	<b>0.0940</b>	32.39%

