

# Aligning LLMs to Recommendation: Progresses and Future Prospects

Fuli Feng

University of Science and Technology of China

# Outline

- **Background**
- **Alignment**
  - **Align to Recommendation Task :**
    - **Fast Alignment : TALLRec [Keqin Bao et al. RecSys 23]**
    - **Align to Generative Recommendation [Keqin Bao et al. arXiv 23]**
    - **Align Language to Recommendation Items [Xinyu Lin et al. arXiv 23]**
  - **Align to Recommendation Modality:**
    - **Empowering LLM Recommendation with Modality Alignment [Yang Zhang et al. arXiv 23]**
    - **Align to Understand Recommendation Modality [Zhengyi Yang et al. arXiv 23]**
- **Future Work**

# Outline

- **Background**
- **Alignment**
  - **Align to Recommendation Task :**
    - Fast Alignment : TALLRec [Keqin Bao et al. RecSys '23]
    - Align to Generative Recommendation [Keqin Bao et al. arXiv '23]
    - Align Language to Recommendation Items [Xinyu Lin et al. arXiv 23]
  - **Align to Recommendation Modality:**
    - Empowering LLM Recommendation with Modality Alignment [Zhengyi Yang et al. arXiv '23]
    - Align to Understand Recommendation Modality [Zhengyi Yang et al. arXiv '23]
- **Future Work**

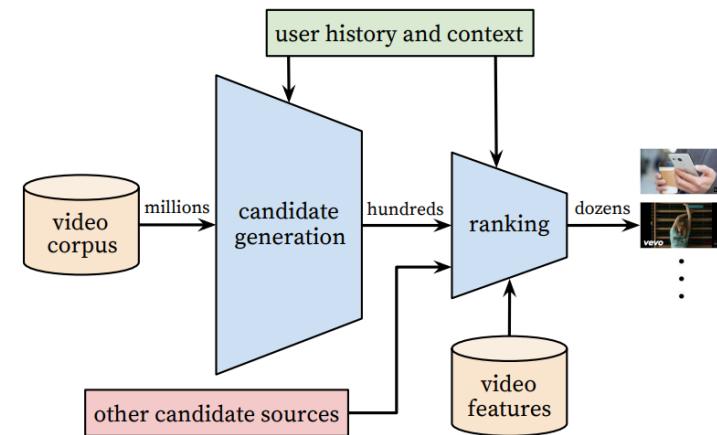
# Background

## □ Information explosion era

- E-commerce: **12 million items** in Amazon.
- Social networks: **2.8 billion users** in Facebook.
- Content sharing platforms: **720,000 hours videos** uploaded to Youtube per day.



## □ Recommender system

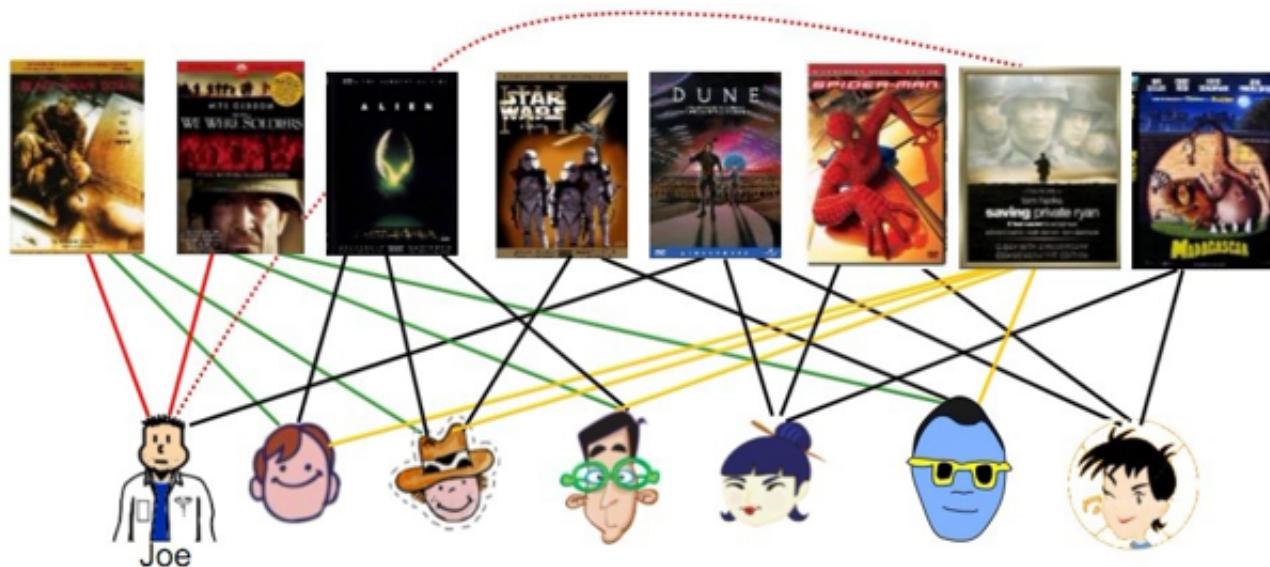


# Background

## □ Core idea of personalized recommendation

- **Collaborative filtering (CF):**

Making automatic predictions (filtering) about the interests of a user by collecting preferences from many users (collaborating).



		item			
		1	2	3	4
user	1	5	?	?	?
	2	3	4	?	?
3	?	1	2	4	...
...	...	...	...	...	...

Interaction Matrix

Memory-based CF

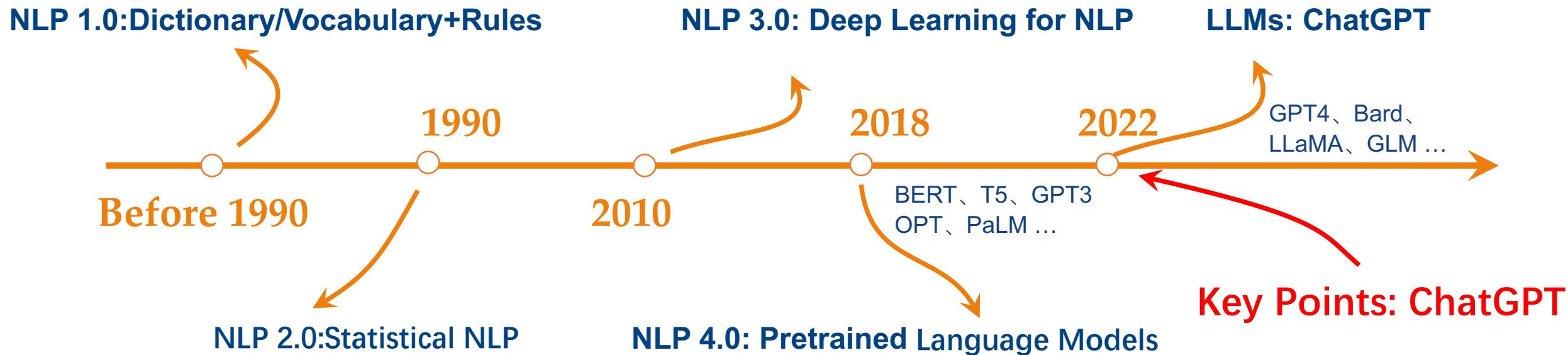
- User CF
- Item CF

Model-based CF

- MF
- FISM
- ...

# Background

## □ Development of NLP techniques



# Background

## □ How recommender systems benefit from LLMs

- **Representation:**

Textual item representation,  
Knowledge representation

- **Generalization:**

Few/zero-shot inference,  
Cross-domain,  
World knowledge

- **Generation:**

Personalized content generation,  
Generative recommendation

**However, LLMs are **not** naturally suitable for recommendation**

- **Reasons:**

Lack of recommendation-related training data in the LLM pre-training stage.  
Hard to capture the collaborative information.

# Outline

- **Background**
- **Alignment**
  - **Align to Recommendation Task:**
    - **Fast Alignment : TALLRec [Keqin Bao et al. RecSys 23]**
    - **Align to Generative Recommendation: BIGRec [Keqin Bao et al. arXiv 23]**
    - **Align Language to Recommendation Items: TransRec [Xinyu Lin et al. arXiv 23]**
  - **Align to Recommendation Modality:**
    - Empowering LLM Recommendation with Modality Alignment [Zhengyi Yang et al. arXiv 23]
    - Align to Understand Recommendation Modality [Zhengyi Yang et al. arXiv 23]
- **Future Work**

# Fast Alignment

- **In-context learning can recommend ?**
- **Example:** Given the user' s interaction history, directly ask the LLM to predict whether he/she will like the new target item via answering "Yes" or "No" .

---

## Instruction Input

---

Task Instruction: Given the user's historical interactions, please determine whether the user will enjoy the target new movie by answering "Yes" or "No".

---

Task Input: User's liked items: GodFather.  
User's disliked items: Star Wars.  
Target new movie: Iron Man

---

## Instruction Output

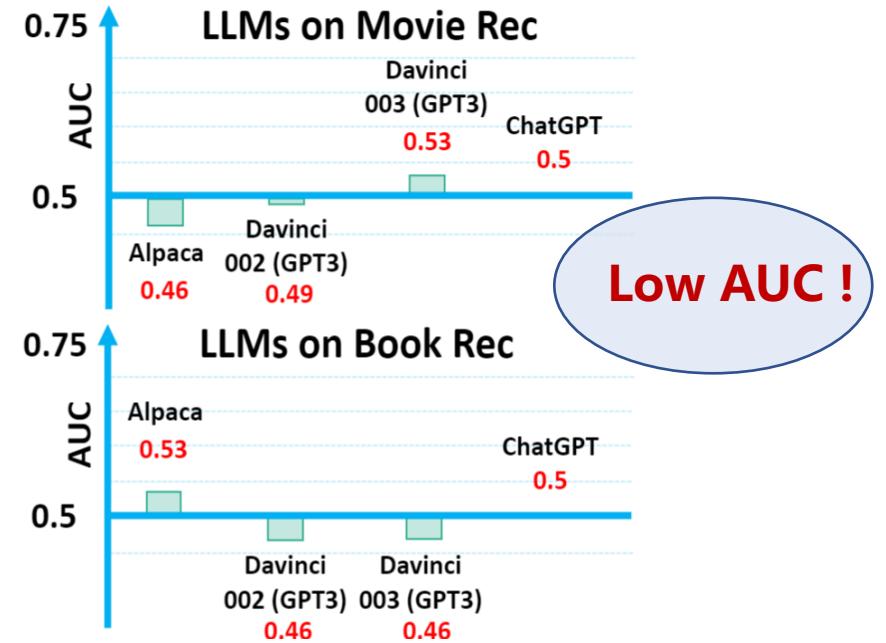
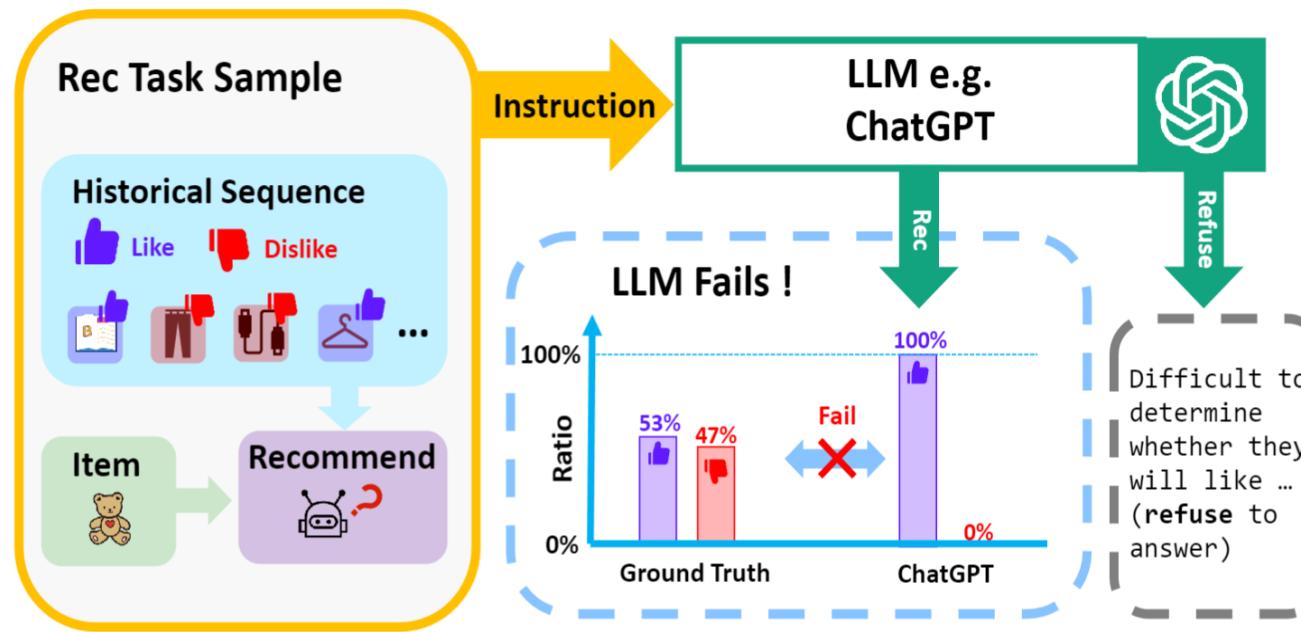
---

Task Output: No.

---

# Fast Alignment

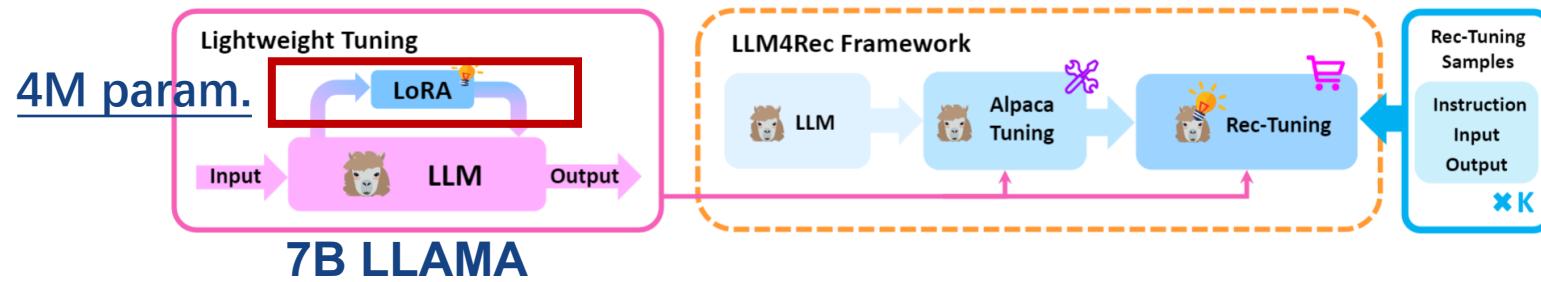
- In-context learning is not enough.
- In complex scenarios, ChatGPT usually gives positive ratings or refuse to answer.



Need to **align** LLM with recommendation!

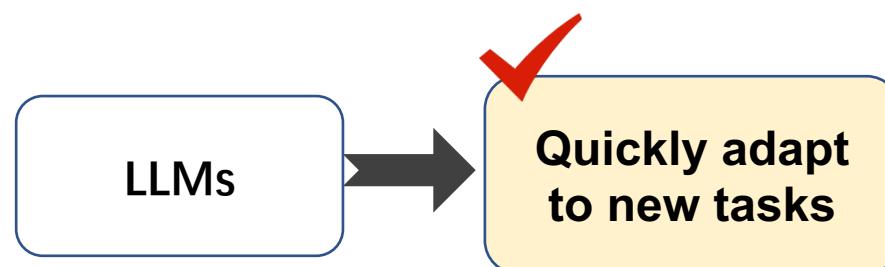
# Fast Alignment

## □ Instruction-tuning



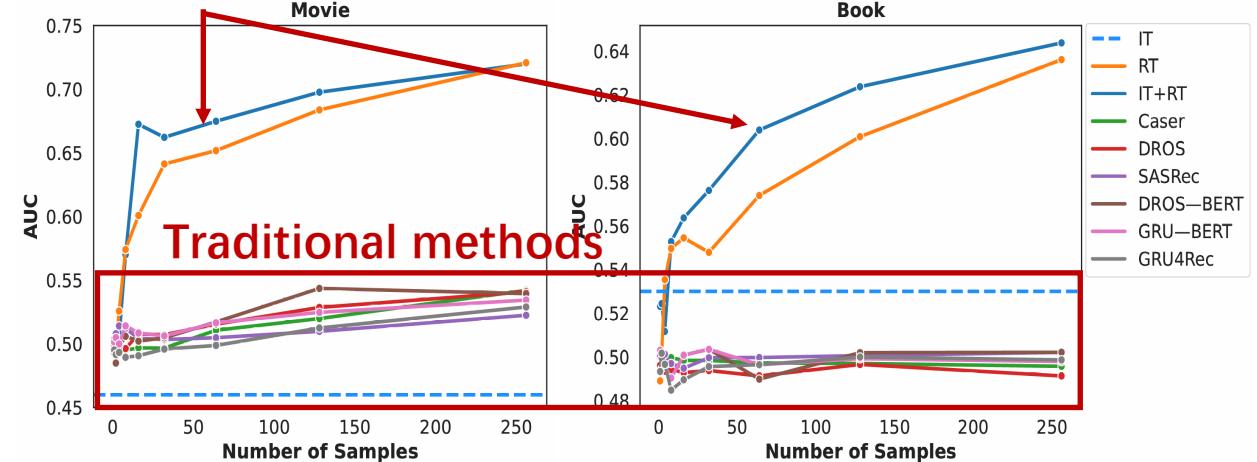
$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi+\Theta}(y_t|x, y_{<t})),$$

Fine-tune 4M parameters by few-shot samples via generative loss



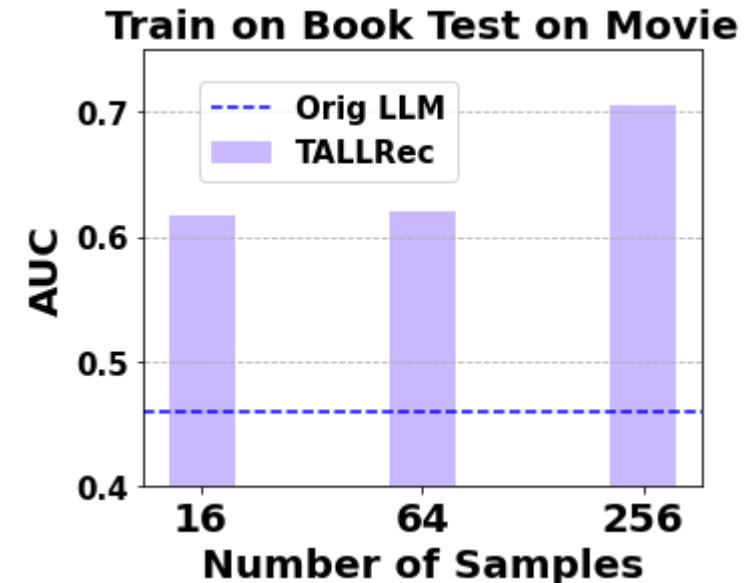
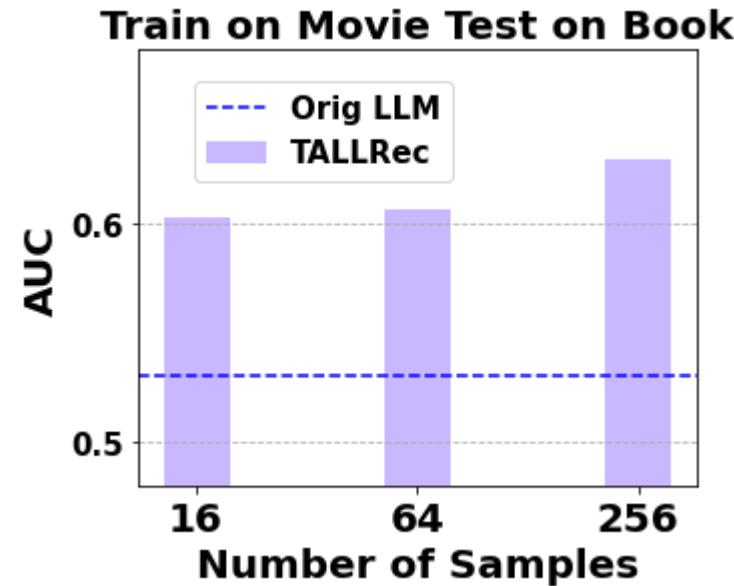
- Use item titles as the input
- Better for cold-start recommendation

Performance significantly improves by fine-tuning few-shot samples.



# Fast Alignment

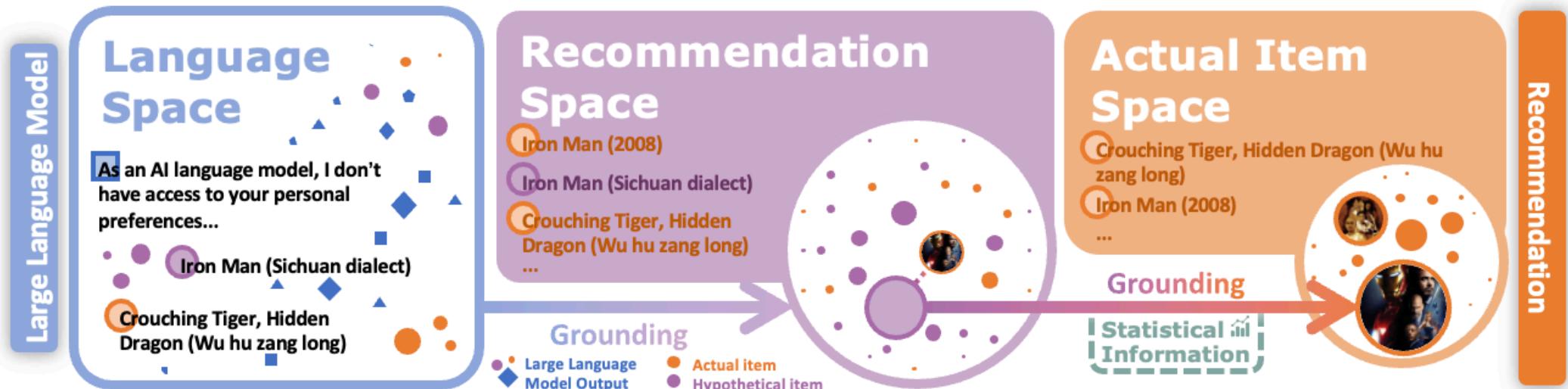
- ❑ Cross-domain generalization
- ❑ Learning from movie scenario can directly recommend on books, and vice versa
- ❑ LLM can leverage domain knowledge to accomplish recommendation tasks after acquiring the ability to recommend.



# Align with Grounding

## □ Generation + Grounding

- Generation ability is the important feature of the LLM, and it almost can generate **all conceivable language sequences**.
- However, LLMs don't know which kind of **sequences describe a item** in the recommendation scenario.
- The item described by the LLM may not in **the actual world**.



## Grounding Paradigm

### Language Space

↓ Step1: instruction tuning

### Recommendation Space

↓ Step2: L2 distance between representations

### Actual Item Space

# Align with Grounding

## □ Generation + Grounding

### □ Few-shot training

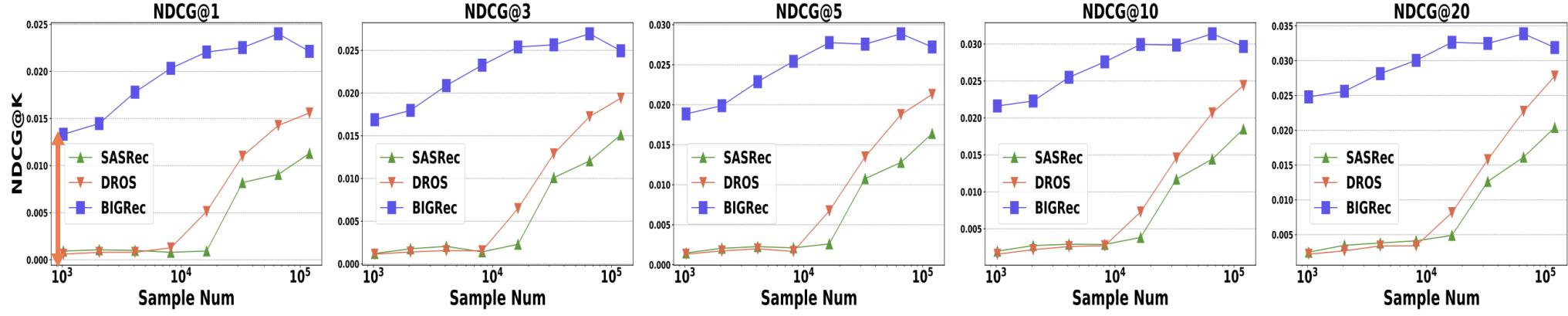
Dataset	Model	NG@1	NG@3	NG@5	NG@10	NG@20	HR@1	HR@3	HR@5	HR@10	HR@20
Movie	GRU4Rec	0.0015	0.0034	0.0047	0.0070	0.0104	0.0015	0.0047	0.0079	0.0147	0.0281
	Caser	0.0020	0.0035	0.0052	0.0078	0.0109	0.0020	0.0046	0.0088	0.0171	0.0293
	SASRec	0.0023	0.0051	0.0062	0.0082	0.0117	0.0023	0.0070	0.0097	0.0161	0.0301
	P5	0.0014	0.0026	0.0036	0.0051	0.0069	0.0014	0.0035	0.0059	0.0107	0.0176
	DROS	0.0022	0.0040	0.0052	0.0081	0.0112	0.0022	0.0051	0.0081	0.0173	0.0297
	GPT4Rec-LLaMA	0.0016	0.0022	0.0024	0.0028	0.0035	0.0016	0.0026	0.0030	0.0044	0.0074
	<b>BIGRec (1024)</b>	<b>0.0176</b>	<b>0.0214</b>	<b>0.0230</b>	<b>0.0257</b>	<b>0.0283</b>	<b>0.0176</b>	<b>0.0241</b>	<b>0.0281</b>	<b>0.0366</b>	<b>0.0471</b>
Game	<b>Improve</b>	<b>654.29%</b>	<b>323.31%</b>	<b>273.70%</b>	<b>213.71%</b>	<b>142.55%</b>	<b>654.29%</b>	<b>244.71%</b>	<b>188.39%</b>	<b>111.97%</b>	<b>56.55%</b>
	GRU4Rec	0.0013	0.0016	0.0018	0.0024	0.0030	0.0013	0.0018	0.0024	0.0041	0.0069
	Caser	0.0007	0.0012	0.0019	0.0024	0.0035	0.0007	0.0016	0.0032	0.0048	0.0092
	SASRec	0.0009	0.0012	0.0015	0.0020	0.0025	0.0009	0.0015	0.0021	0.0037	0.0057
	P5	0.0002	0.0005	0.0007	0.0010	0.0017	0.0002	0.0007	0.0012	0.0023	0.0049
	DROS	0.0006	0.0011	0.0013	0.0016	0.0022	0.0006	0.0015	0.0019	0.0027	0.0052
	GPT4Rec-LLaMA	0.0000	0.0000	0.0000	0.0001	0.0001	0.0000	0.0000	0.0000	0.0002	0.0002
<b>BIGRec (1024)</b>	<b>0.0133</b>	<b>0.0169</b>	<b>0.0189</b>	<b>0.0216</b>	<b>0.0248</b>	<b>0.0133</b>	<b>0.0195</b>	<b>0.0243</b>	<b>0.0329</b>	<b>0.0457</b>	
	<b>Improve</b>	<b>952.63%</b>	<b>976.26%</b>	<b>888.19%</b>	<b>799.64%</b>	<b>613.76%</b>	<b>952.63%</b>	<b>985.19%</b>	<b>660.42%</b>	<b>586.11%</b>	<b>397.10%</b>

- Baselines exhibit significantly worse performance than BIGRec.
- Improvement of BIGRec is significantly higher on Game compared to on Movie.
  - possibly due to the varying properties of popularity bias between the two datasets.

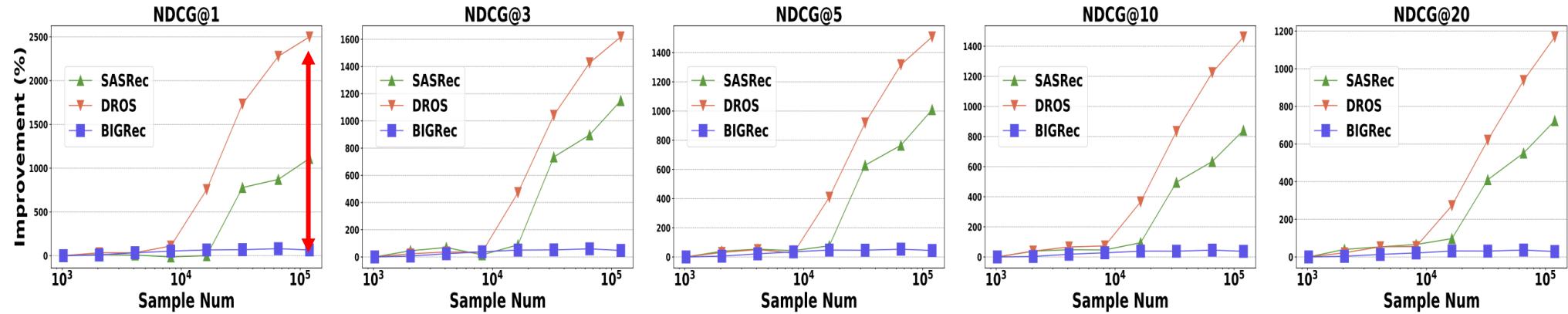
# Align with Grounding

## □ Generation + Grounding

Quickly Adapt to Recommendation



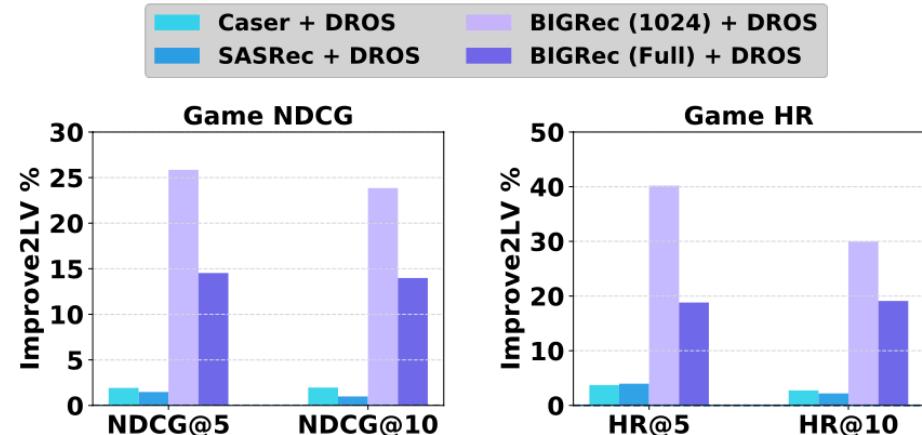
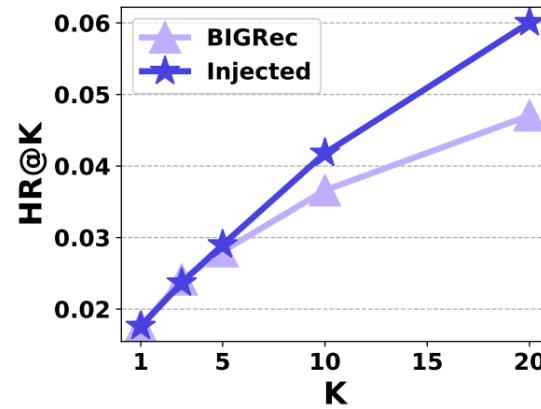
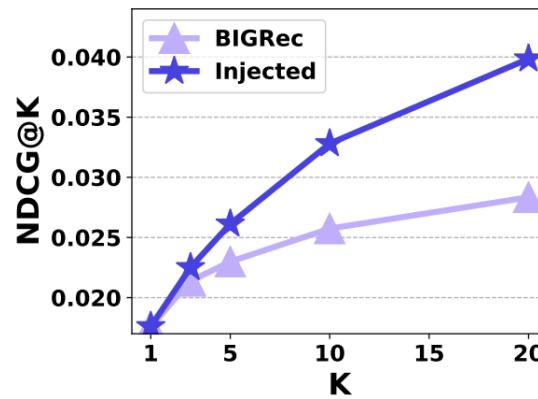
Not proficient in utilizing CF info.



# Align with Grounding

## □ Generation + Grounding

- In-depth analysis
- Injecting statistical information into BIGRec at step2



- By incorporating popularity, BIGRec achieves significant improvements w.r.t.  $\text{NDCG}@K$  and  $\text{HR}@K$  particularly for a larger  $K$ .
- Incorporating collaborative information into BIGRec yields more significant enhancements than conventional models.

# Align with Grounding

## □ Align Language and Items: Two Key Steps

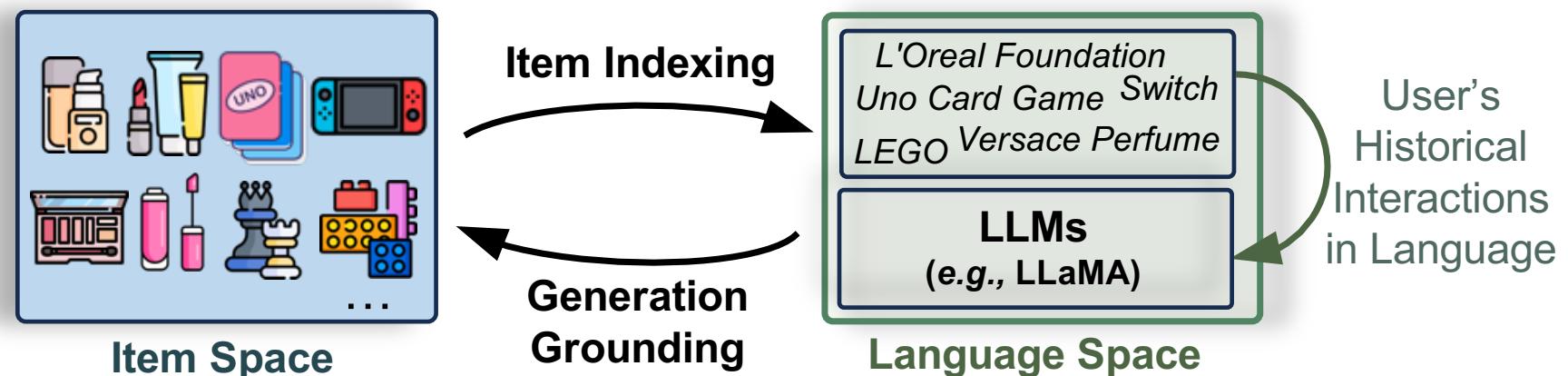
### □ Item indexing

- How to represent an item in natural language?

### □ Generation grounding

- How to generate and ground the natural language to the actual items?

- ✓ Two criteria for identifiers
  - **distinctiveness**
  - **semantics**
- ✓ Consideration for generation
  - **constrained generation**
  - **position-free generation**
  - ⚠ rely heavily on first token

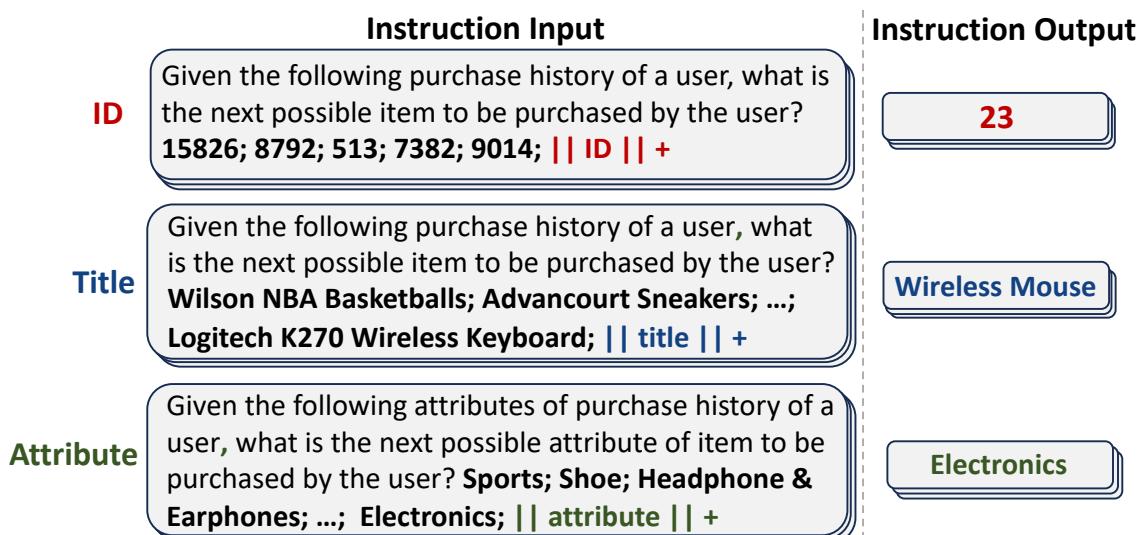


# Align with Grounding

- Item indexing: multi-facet identifier



- Instruction data reconstruction

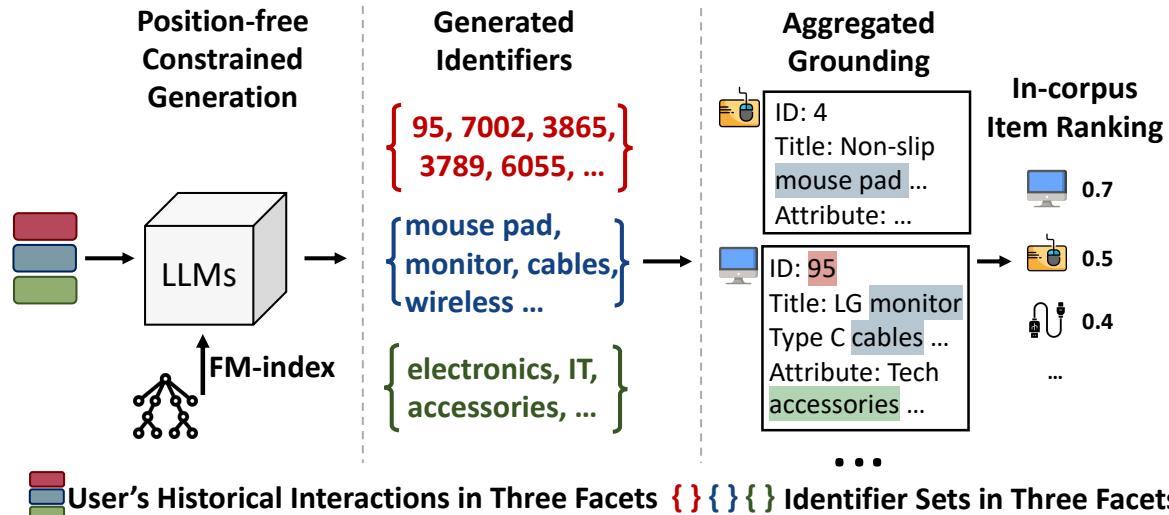


- Generation grounding:

- position-free constrained generation
- aggregated grounding



FM-index: special prefix tree that supports search from any position of the identifier corpus.



# Align with Grounding

## □ Enhanced recommendation under full training

Table 2: Overall performance comparison between the baselines and TransRec instantiated on BART on three datasets. The best results are highlighted in bold and the second-best results are underlined. \* implies the improvements over the second-best results are statistically significant ( $p$ -value < 0.01) under one-sample t-tests.

Model	Beauty				Toys				Yelp			
	R@5	R@10	N@5	N@10	R@5	R@10	N@5	N@10	R@5	R@10	N@5	N@10
MF	0.0294	0.0474	0.0145	0.0191	0.0236	0.0355	0.0153	0.0192	0.0220	0.0381	0.0138	0.0190
LightGCN	0.0305	0.0511	0.0194	0.0260	0.0322	0.0508	0.0215	0.0275	0.0255	0.0427	0.0163	0.0218
SASRec	0.0380	0.0588	0.0246	0.0313	0.0470	0.0659	0.0312	0.0373	0.0183	0.0296	0.0116	0.0152
ACVAE	0.0503	0.0710	0.0356	0.0422	0.0488	0.0679	0.0350	0.0411	0.0211	0.0356	0.0127	0.0174
SID	0.0430	0.0602	0.0288	0.0368	0.0164	0.0218	0.0120	0.0139	0.0346	0.0486	0.0242	0.0287
SemID+IID	0.0501	0.0724	0.0344	0.0411	0.0145	0.0260	0.0069	0.0123	0.0229	0.0382	0.0150	0.0199
CID+IID	<u>0.0512</u>	<u>0.0732</u>	0.0356	0.0427	0.0169	0.0276	0.0104	0.0154	0.0287	0.0468	0.0195	0.0254
P5	0.0508	0.0664	<u>0.0379</u>	<u>0.0429</u>	0.0608	0.0688	<u>0.0507</u>	0.0534	0.0506	<u>0.0648</u>	<u>0.0343</u>	0.0389
TransRec-B	<b>0.0550*</b>	<b>0.0766*</b>	<b>0.0395*</b>	<b>0.0464*</b>	<b>0.0679*</b>	<b>0.0881*</b>	<b>0.0512*</b>	<b>0.0567*</b>	<b>0.0549*</b>	<b>0.0679*</b>	<b>0.0408*</b>	<b>0.0450*</b>

- Superior performance compared to both traditional models and LLM-based models.
- TransRec significantly outperform P5, which utilizes additional training data from user reviews. In contrast, TransRec uses only user-item interactions.

# Align with Grounding

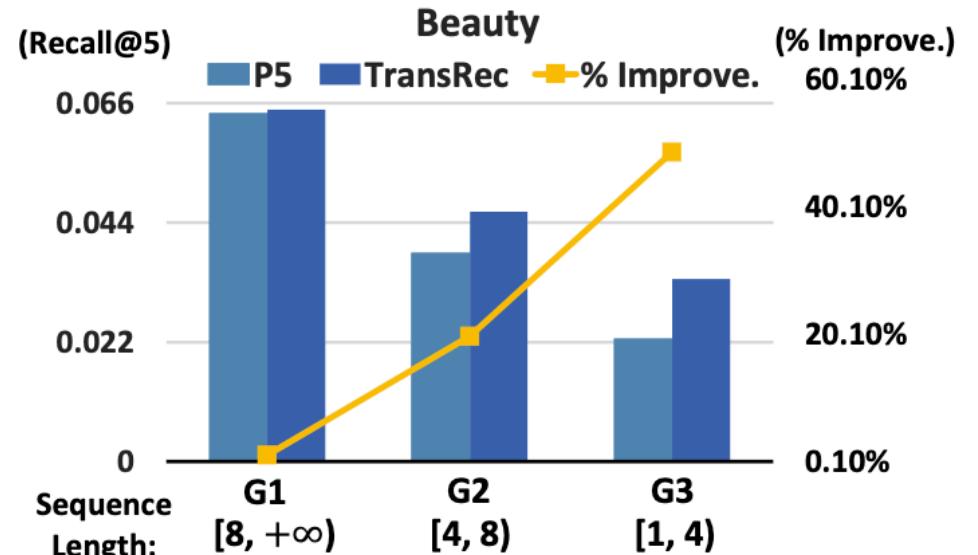
## □ Strong generalization ability

- Few-shot training
  - warm- and cold-start testing
- User group analysis
  - from dense users to sparse users

N-shot	Model	Warm		Cold	
		R@5	N@5	R@5	N@5
1024	LightGCN	0.0205	0.0125	0.0005	0.0003
	ACVAE	0.0098	0.0057	0.0047	0.0026
	P5	0.0040	0.0016	0.0025	0.0015
	TransRec-B	0.0039	0.0024	0.0025	0.0016
	TransRec-L	<b>0.0141</b>	<b>0.0070</b>	<b>0.0159</b>	<b>0.0097</b>
2048	LightGCN	0.0186	0.0117	0.0005	0.0004
	ACVAE	0.0229	0.0136	0.0074	0.0044
	P5	0.0047	0.0030	0.0036	0.0012
	TransRec-B	0.0052	0.0027	0.0039	0.0017
	TransRec-L	<b>0.0194</b>	<b>0.0126</b>	<b>0.0206</b>	<b>0.0126</b>

\* The bold results highlight the superior performance compared to the best LLM-based recommender baseline.

- Remarkable generalization ability of LLMs with vase knowledge base, especially on cold-start recommendation under limited data.
- On user side, TransRec significantly improves the performance of sparse users with fewer interactions.



# Outline

- **Background**
- **Alignment**
  - **Align to Recommendation Task :**
    - Fast Alignment : TALLRec [Keqin Bao et al. RecSys 23]
    - Align to Generative Recommendation [Keqin Bao et al. arXiv 23]
    - Align Language to Recommendation Items: TransRec [Xinyu Lin et al. arXiv 23]
  - **Align to Recommendation Modality :**
    - Empowering LLM Recommendation with Modality Alignment [Yang Zhang et al. arXiv 23]
    - Align to Understand Recommendation Modality [Zhengyi Yang et al. arXiv 23]
- **Future Work**

# Empowering LLM Rec with Modality Alignment

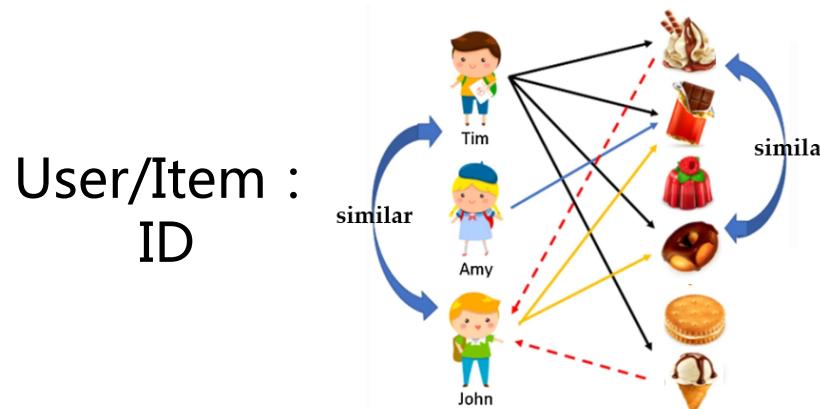
LLM Rec is not good at modeling collaborative information as traditional models

## LLM Rec vs Traditional CF Model :

- Excellent at old-start scenarios
- Poor at warm-start scenarios



## Traditional CF Model

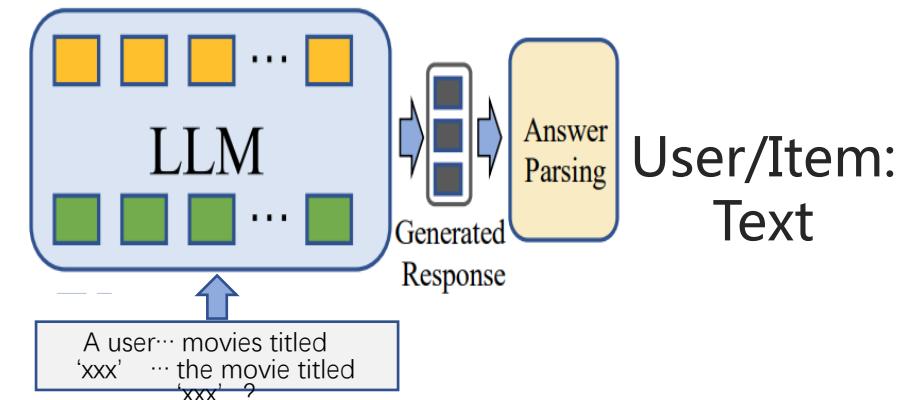


Rely on collab. Info. --- co-occurrence similarities among interactions (Good for Warm)

**Lack** of modeling  
collab. Info.

Textually similar item  
may have distinct  
collab. info.

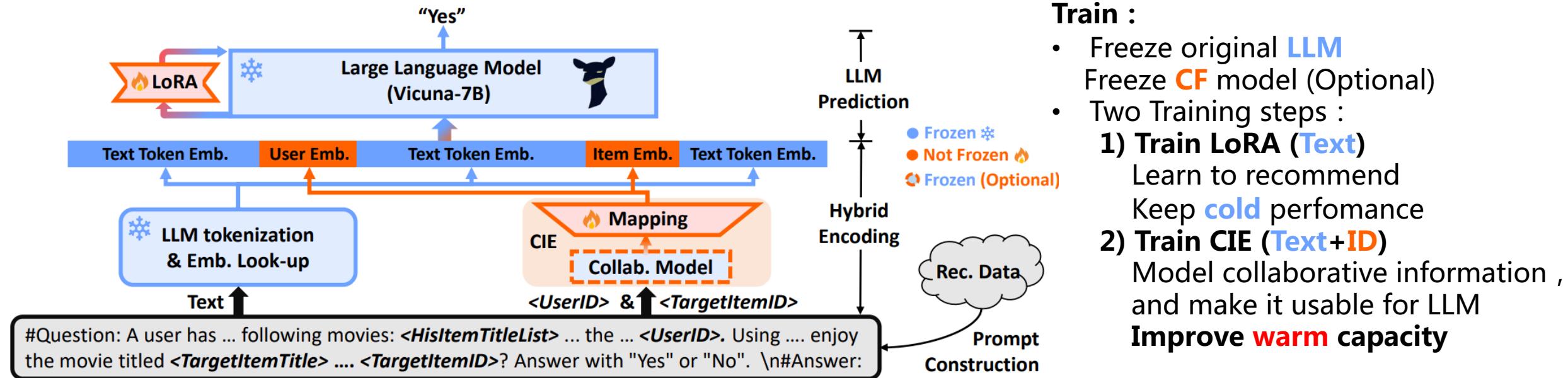
## LLM Recommendation



Relying on text semantics  
(Good for Cold)

# Empowering LLM Rec with Modality Alignment

## CoLLM : Integrating Collaborative Embedding into LLM Rec – Align with Rec Modality



### Train :

- Freeze original **LLM**
- Freeze **CF** model (Optional)
- Two Training steps :

#### 1) Train LoRA (**Text**)

Learn to recommend  
Keep **cold** performance

#### 2) Train CIE (**Text+ID**)

Model collaborative information ,  
and make it usable for LLM

**Improve warm capacity**

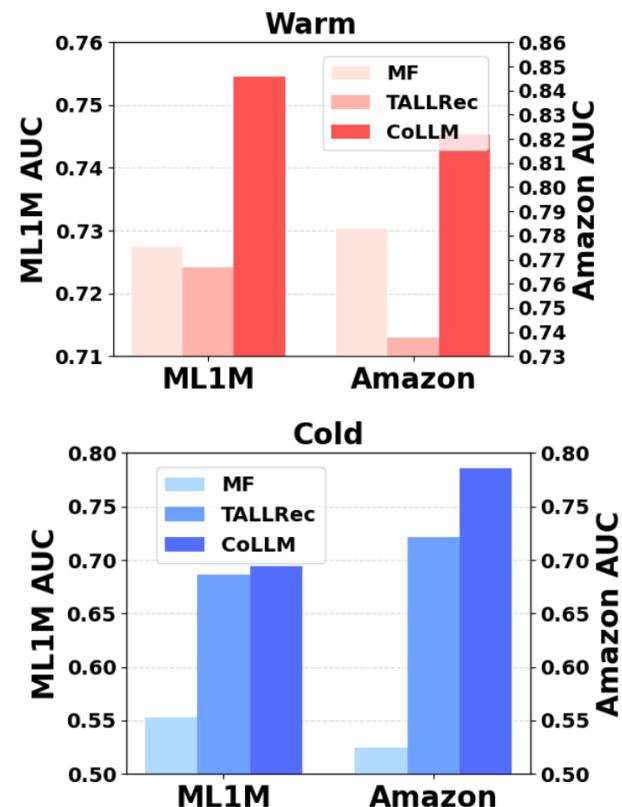
- **Prompt construction:** add <UserID> and <TargetID> for placing the Collab. Info.
- **Hybrid Encoding:**
  - text: tokenization & LLM emb Lookup;
  - user/item ID: CIE --- extract info with collab. model (**low rank**), then map it to the token embedding space
- **LLM prediction:** add a LoRA module for recommendation task learning

# Empowering LLM Rec with Modality Alignment

## The effect of alignment with rec modality

### Overall Performance

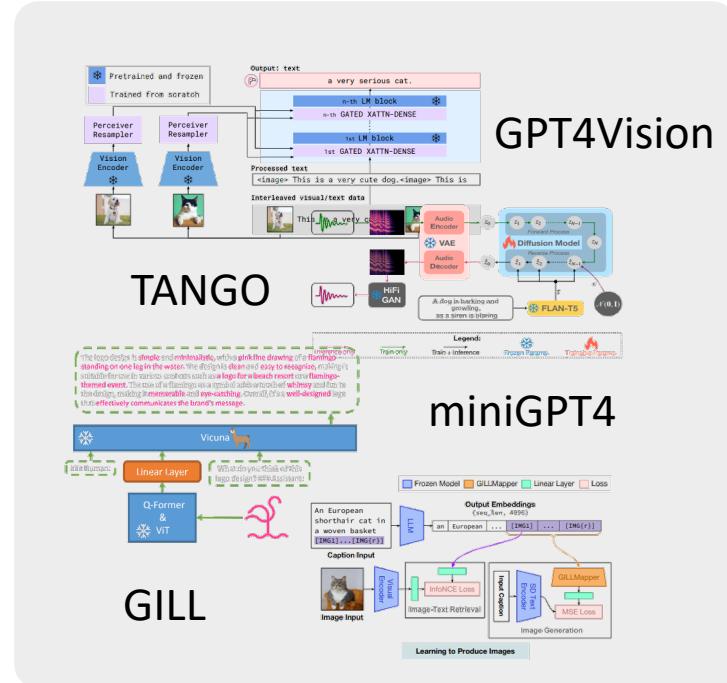
Dataset		ML-1M			Amazon-Book		
Methods		AUC	UAUC	Rel. Imp.	AUC	UAUC	Rel. Imp.
Collab.	MF	0.6482	0.6361	10.3%	0.7134	0.5565	12.8%
	LightGCN	0.5959	0.6499	13.2%	0.7103	0.5639	10.7%
	SASRec	0.7078	0.6884	1.9%	0.6887	0.5714	8.4%
LLMRec	ICL	0.5320	0.5268	33.8%	0.4820	0.4856	48.2%
	Soft-Prompt	0.7071	0.6739	2.7%	0.7224	0.5881	10.4%
	TALLRec	0.7097	0.6818	1.8%	0.7375	0.5983	8.2%
Ours	CoLLM-MF	0.7295	0.6875	-	0.8109	0.6225	-
	CoLLM-LightGCN	0.7100	0.6967	-	0.7978	0.6149	-
	CoLLM-SASRec	0.7235	0.6990	-	0.7746	0.5962	-



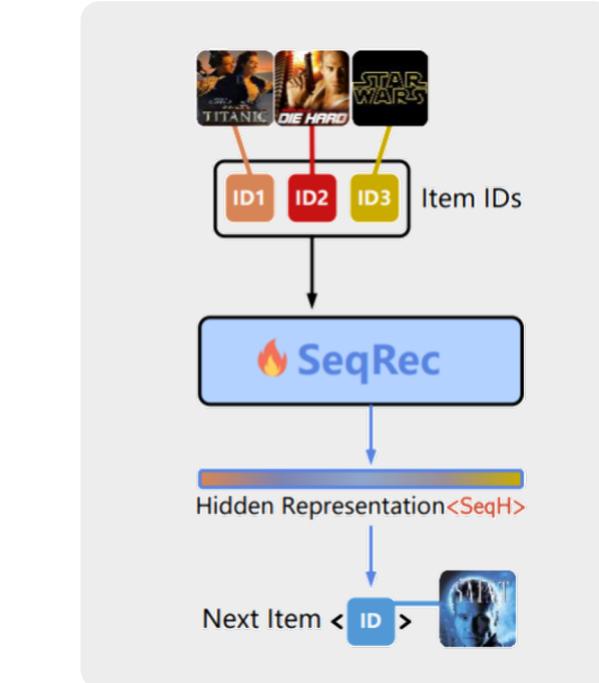
- CoLLM brings performance improvements over traditional collaboration models and current LLM Rec in most cases
- CoLLM can significantly improve the warm performance of LLM Rec (TALLRec), while ensuring cold scene performance

# Align to Understand Recommendation Modality

Trend of universal LLM: align models of various modalities with LLM.  
 Basic research question: can we align recommenders with LLM?



Models of various modalities can be unified by LLM.



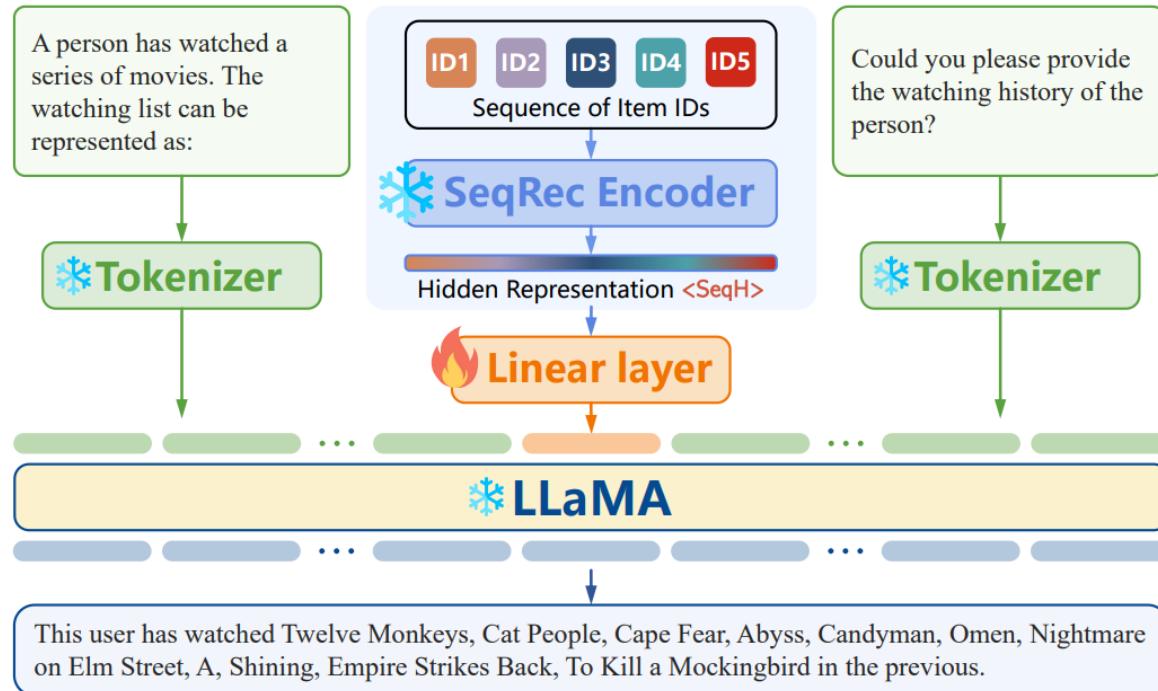
Conventional recommenders:  
 ID-based, interact at hidden space

**Basic research question:**

- Can conventional recommenders be unified to LLM w.r.t., the hidden representation?

# Align to Understand Recommendation Modality

Straightforward idea: recover items encoded in the hidden representation.



The sequential recommender and LLM are both frozen, tuning a linear projection layer.

Sequence-Recovery Prompt	
Input Prompt	A person has watched a series of movies. The watching list can be represented as: <SeqH>. Describe this watching history of the person in detail.
Target Response	This user has watched <b>Twelve Monkeys, Cat People, Cape Fear, Abyss, Candyman, Omen, Nightmare on Elm Street, Shining, Empire Strikes Back, To Kill a Mockingbird</b> in the previous.

Learning objective:

$$p(\mathbf{X}_{Target} | \mathbf{X}_{Input}) = \prod_{i=1}^N p(\mathbf{X}_{Target}^i | \mathbf{X}_{Input}, \mathbf{X}_{Target}^{[1:i-1]}),$$

Inference phase:

- Only provide the input prompt containing the hidden representation <SeqH>.

# Align to Understand Recommendation Modality

**Table 1: Statistics of datasets.**

Dataset	MovieLens	Steam
#sequences	943	11,938
#items	1,682	3,581
#interactions	100,000	274,726

MovieLens, SASRec, Max\_SeqLength=10

recover 4 items	recover 5 items	recover 6 items
25.26%	16.84%	9.47%

Steam, DreamRec, Max\_SeqLength=10

recover 4 items	recover 5 items	recover 6 items
30.54%	24.58%	5.12%

## MoviesLens



A person has watched a series of movies. The watching list can be represented as: <SeqH>. Describe this watching history of the person in detail.



This user has watched [Air Force One](#), [Saint, Conspiracy Theory](#), [Murder at 1600](#), [Cop Land](#), [Liar Liar](#), [Peacemaker](#), [Kiss the Girls](#), [Scream](#), [Devil's Advocate](#).

## Steam



A person has played a series of video games on Steam. The playing list can be represented as: <SeqH>. What games do you think the person has played before?

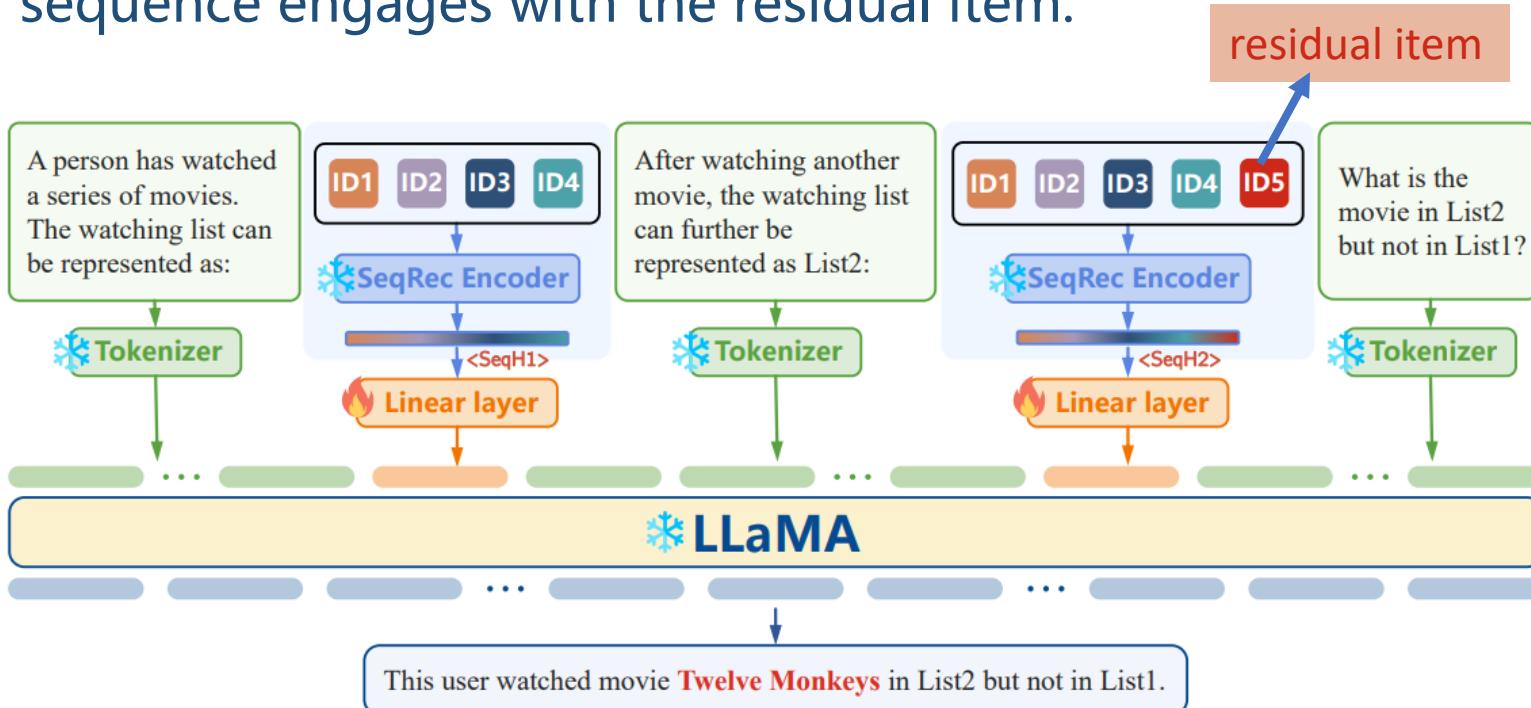


This user has played [Mark of the Ninja](#), [Brothers - A Tale of Two Sons](#), [The Walking Dead: Season 2](#), [The Witcher 2: Assassins of Kings Enhanced Edition](#), [The Evil Within](#), [The Last of Us](#), [Far Cry 3](#), [The Darkness II](#), [Hotline Miami](#).

The **blue text** indicates the correctly recovered items

# Align to Understand Recommendation Modality

More delicate design: identify the **residual item** from hidden representations before and after the sequence engages with the residual item.



Sequence-Residual Prompt	
Input Prompt	A person has watched a series of movies. The watching list can be represented as List1: <SeqH1>. After watching another movie, the watching list can further be represented as List2: <SeqH2>. What is the movie in List2 but not in List1?
Target Response	This user watched movie <b>Twelve Monkeys</b> in List2 but not in List1.

Accuracy of the correctly identifying residual item:

Dataset	GRU4Rec	Caser	SASRec	DreamRec
MovieLens	52.63%	78.95%	93.68%	97.89%
Steam	17.11%	55.03%	52.60%	86.33%

LLM can correctly identify the residual item, especially with more advanced recommenders

# Outline

- **Background**
- **Alignment**
- **Future Work**
  - **Deployment Challenge**
  - **Fairness & Bias [Jizhi Zhang et al. RecSys 23]**
  - **Generative Recommendation [Wenjie Wang et al. arXiv 23]**
  - **Context enhanced LLM Recommendation**

# Deployment Challenge

- **Inference Efficiency**
  - The recommended scenario requires low latency
  - In some scenarios, there are tens of thousands of historical interaction sequences
  - LLMs involve tens of billions or even hundreds of billions, which has extremely high requirements for GPU resources
- **Marginal Profit**
  - Inference cost paid for recommendation/expected gains after recommendation

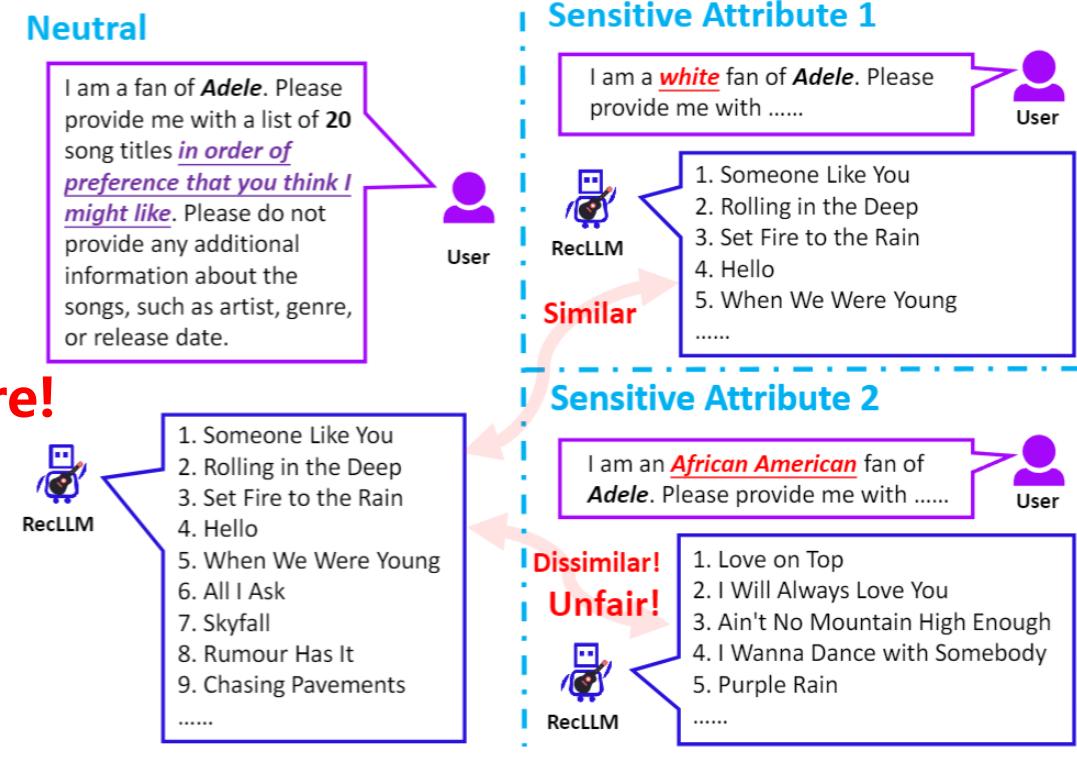
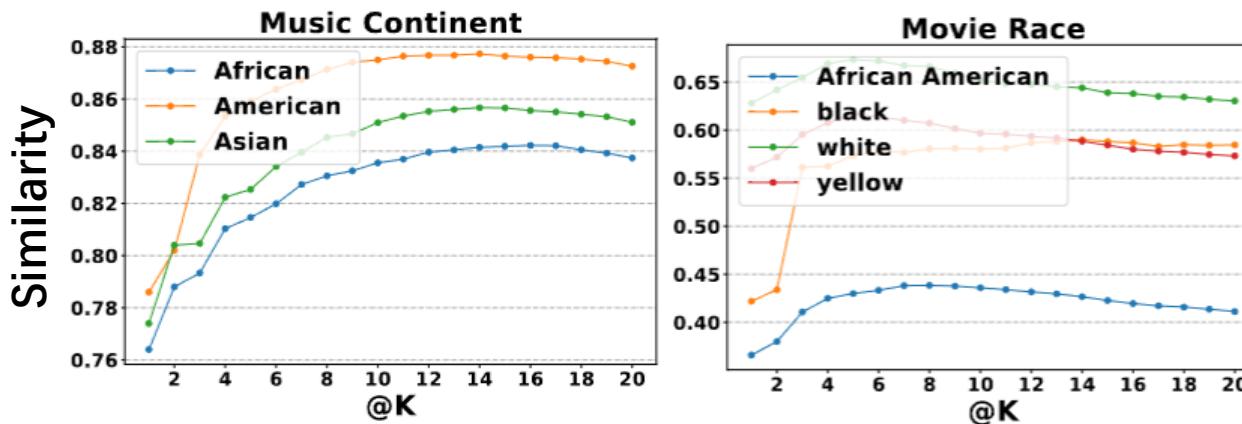
# Fairness and Bias

## □ Fairness in LLM4Rec

- **Discrimination** against certain groups
- LLM4Rec **inherits existing social biases!**
- It is essential to protect vulnerable groups.

**LLM4Rec favors certain populations!**

**LLM4Rec needs to improve fairness in the future!**

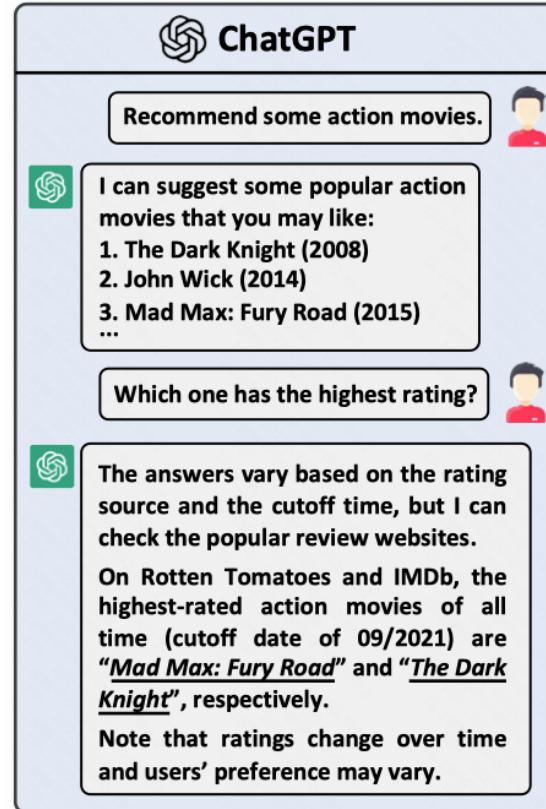


# Generative Recommendation

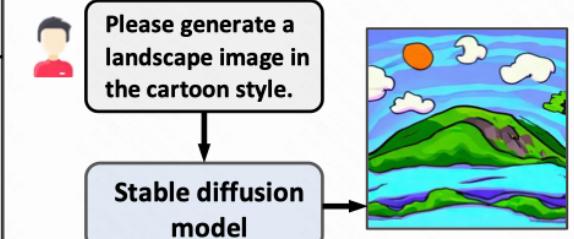
**Generate Content for user' s personalized information needs**

## Application:

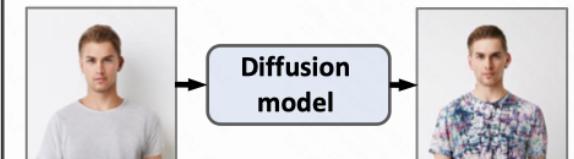
- 1. Image Generative Recommendation
- 2. Music Generative Recommendation
- 3. Fashion Generative Recommendation
- 4. Video Generative Recommendation
- 5. Conversational Recommendation
- ○ ○



(a) A conversation between a user and ChatGPT.



(b) An example of conditional image generation via stable diffusion.



(c) An example of changing image attributes (color change in clothes).

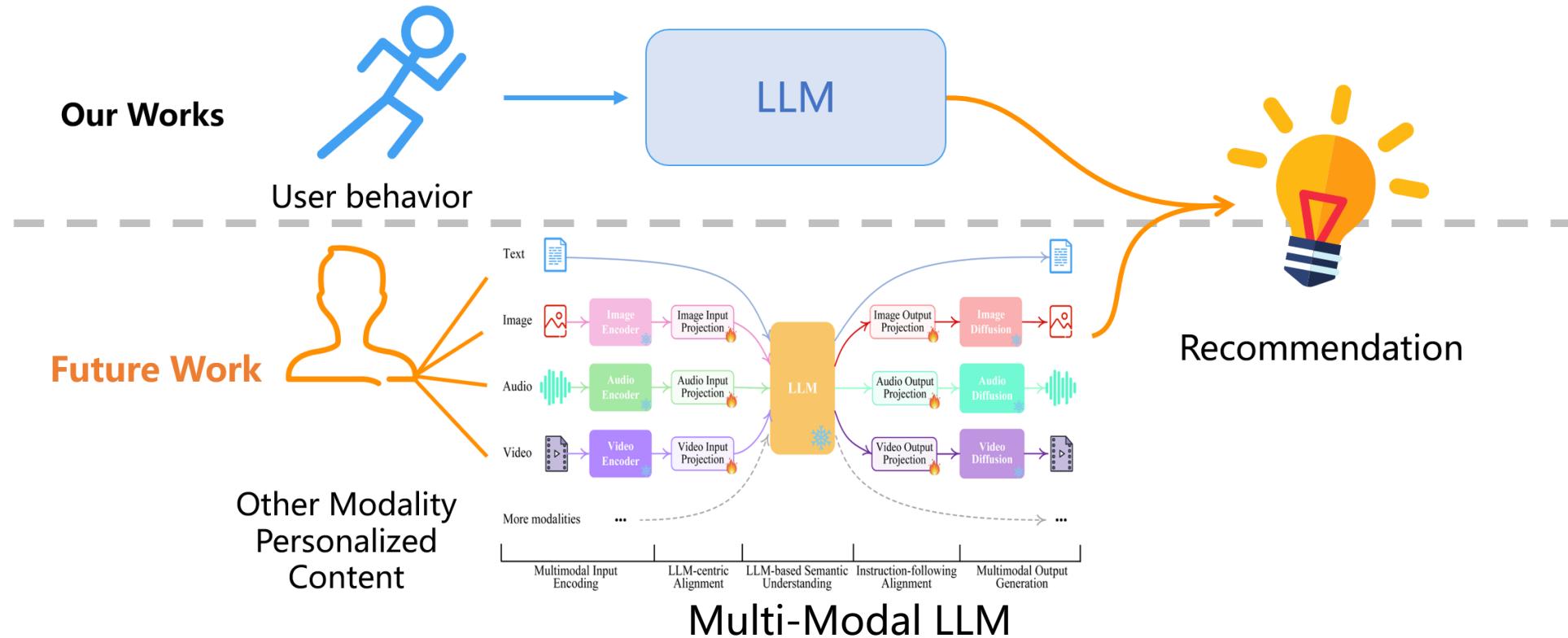


(d) An example of image style transfer (to a cartoon style).

# Context enhanced LLM Recommendation

Our previous work focused primarily on user behavior data. In addition, LLM may also make use of multimodal **context** to enhance recommendations.

LLM can now **understand** and **generate** multimodal context! This ability may be used in the future recommendation.



# Advertisement

## ➤ <@iData> We are hiring: engineer and research scientist

- 合肥总部 + 北京分部, <http://idata.ah.cn/> [hr@idata.ah.cn](mailto:hr@idata.ah.cn)
- 以数据重构网络空间为核心理念, 聚焦大数据、人工智能和网络安全, 汇聚全球顶尖科技人才, 开展前沿技术研究和应用落地:
  - ✓ 顶天: 面向国家战略需求, 承担国家重大工程任务
  - ✓ 立地: 面向市场产业需求, 孵化若干“小而美”的科技公司



岗位	学历	薪资	专业
高级工程师/副研究员/研究员	博士	50W-120W	计算机科学、网络安全、大数据、人工智能、软件工程、统计学、数学等

## ➤ <@USTC> We are hiring: intern, master, PhD, postdoc, and faculty (tenure-track).

- ✓ Details: <https://fulifeng.github.io>
- ✓ Contact: [fengfl@ustc.edu.cn](mailto:fengfl@ustc.edu.cn)



**中国科学技术大学**  
University of Science and Technology of China

## ➤ <@SIGIR-AP 2023> Tutorial On LLM for ecommendation

# THANKS !