# 05

# Summary

Open Challenges and Beyond

Human-readable Data 👨 👩













Machine-readable Data





Human-readable Data 👨 👩













Machine-readable Data





i3487, i124, i19240, i772

One token per action?



### Human-readable Data 👨 👩













# Machine-readable Data





Premium Men's Short Sleeve Athletic Training T-Shirt Made of Lightweight Breathable Fabric, Ideal for Running, Gym Workouts, and Casual Sportswear in All Seasons; High-Performance Breathable Cotton Crew Socks for Men with Arch Support, Cushioned Heel and Toe, and Moisture Control, Perfect for Sports, Walking, and Everyday Comfort; Men's Loose-Fit Basketball Shorts with Elastic Drawstring Waistband, Quick-Dry Mesh Fabric, and Printed Number 11 for Professional and Recreational Play; Official Size 7 Composite Leather Basketball Designed for Indoor and Outdoor Use, Deep Channel Design for Enhanced Grip and Ball Control, Ideal for Training and Competitive Matches;

# Text description of each action?



# **Semantic IDs**

(also called: SemID or SID)

A few tokens that jointly index one item.

t3, t321, t643, t1011



### Can we find a sweet spot between

### One token per action

### i3487, i124, i19240, i772

### **Semantic IDs**

†34, †392, †600, †891, †21, †502, †592, †1002, †123, †403, †611, †821, †21, †502, †711, †1022

### **Text description of each action**

Premium Men's Short Sleeve Athletic
Training T-Shirt Made of Lightweight
Breathable Fabric, Ideal for Running,
Gym Workouts, and Casual Sportswear in
All Seasons; High-Performance
Breathable Cotton Crew Socks for Men
with Arch Support, Cushioned Heel and
Toe, and Moisture Control, Perfect for
Sports, Walking, and Everyday Comfort;
Men's Loose-Fit Basketball Shorts with
Elastic Drawstring Waistband, Quick-Dry
Mesh Fabric, and Printed Number 11 for
Professional and Recreational Play;
Official Size 7 Composite Leather ...

# **Open Challenges**

### Part 1: What becomes harder?

Comparing to traditional RecSys, what challenges may large generative models face?

# **Open Challenges**

### Part 1: What becomes harder?

Comparing to traditional RecSys, what challenges may large generative models face?

### Part 2: What becomes possible?

What new opportunities may large generative models unlock for recommender systems?

# Part 1: What Becomes Harder?

Comparing to traditional RecSys, what challenges may large generative models face?

# **Cold-Start Recommendation**

Do semantic ID-based models really good at cold-start recommendation?

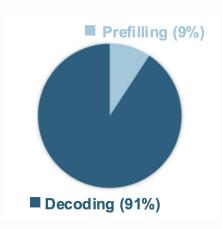
### **Cold-Start Recommendation**

Do semantic ID-based models really good at cold-start recommendation?

Model	#Params. (M)	Video Games						Cell Phones and Accessories					
		Overall		<b>In-Sample</b> (39.7%)		<b>Unseen</b> (60.3%)		Overall		<b>In-Sample</b> (31.8%)		<b>Unseen</b> (68.2%)	
		R@50	N@50	R@50	N@50	R@50	N@50	R@50	N@50	R@50	N@50	R@50	N@50
UniSRec	2.90	0.0621	0.0200	0.1386	0.0461	0.0118	0.0029	0.0233	0.0077	0.0604	0.0211	0.0060	0.0014
Recformer	233.73	0.0740	0.0218	0.1082	0.0333	0.0514	0.0142	0.0236	0.0070	0.0340	0.0103	0.0188	0.0055
TIGER	13.26	0.0584	0.0193	0.1472	0.0486	-	-	0.0232	0.0078	0.0730	0.0245	-	-
$TIGER_C$	13.26	0.0611	0.0198	0.1447	0.0482	0.0061	0.0011	0.0233	0.0078	0.0691	0.0238	0.0019	0.0003
SpecGR <sub>Aux</sub>	16.16	0.0726	0.0220	0.1399	0.0436	0.0283	0.0078	0.0269	0.0084	0.0722	0.0230	0.0058	0.0015
SpecGR++	13.28	0.0717	0.0225	0.1323	0.0439	0.0318	$\underline{0.0084}$	0.0275	0.0090	0.0730	0.0246	0.0063	0.0017

Retrieval Models: K Nearest Neighbor Search

Generative Models (e.g., AR models): Beam Search



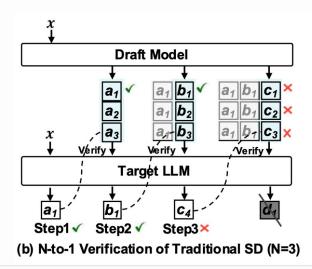
# How to accelerate LLMs? Speculative Decoding

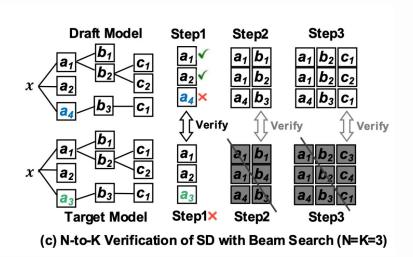
- Use a "cheap" model to generate candidates
- "Expensive" model can accept or reject (and perform inference if necessary)

```
[START] japan 's benchmark bend n
[START] japan 's benchmark nikkei 22 75
[START] japan 's benchmark nikkei 225 index rose 22 76
[START] japan 's benchmark nikkei 225 index rose 226 69 7 points
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 9 1
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 9859
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 7 in
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in tokye late
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in tokye late
[START] japan 's benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in late morning trading [END]
```

Speculative decoding for generative rec? X

### N-to-K verification



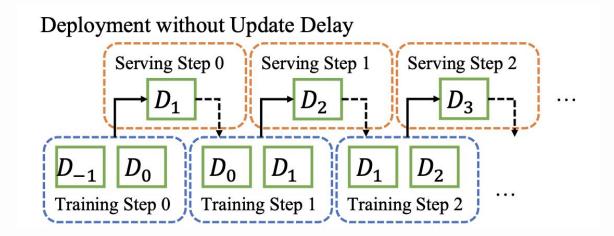


In addition to single-model acceleration methods, what about "serving throughout"?

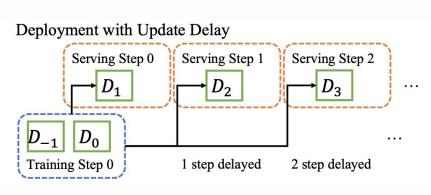
Example: VLLM offers solutions for high-throughput and memory-efficient inference and serving

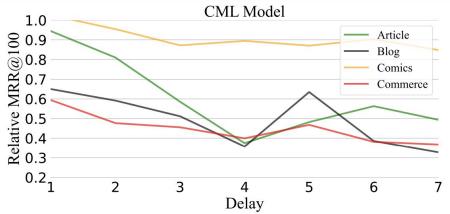
What's unique for generative rec?

Recommendation models favor timely updates



Delayed updates lead to performance degradation

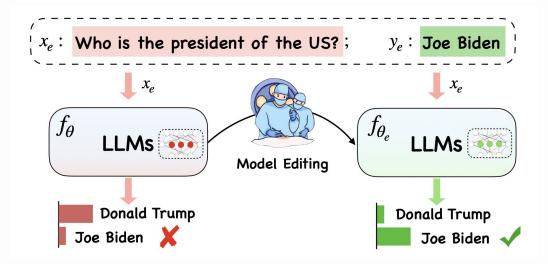




How to update large generative rec models timely?

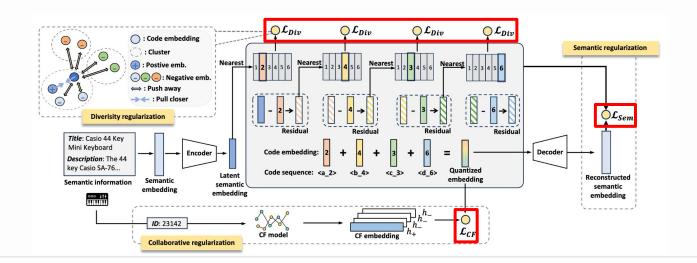
(Frequently retraining large generative models may be resource consuming)

How to update large generative rec models timely?



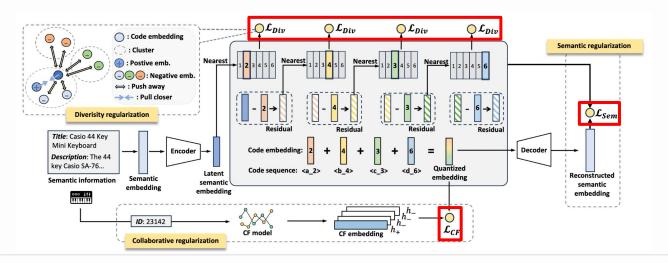
Knowledge editing?

Multiple objectives for optimizing item tokenization ...



Multiple objectives for optimizing item tokenization ...

But none of them is directly related to rec performance



### reconstruction loss ≠ downstream performance

How to connect tokenization objective with recommendation performance?

Zipf's distribution? Entropy? Linguistic metrics?

Language Tokenization

2014~2015: Word / Char

Context-independent ⇒ Context-aware

### Language Tokenization

2014~2015: Word / Char 2016~present: BPE / WordPiece

Context-independent ⇒ Context-aware

### Language Tokenization

2014~2015: Word / Char

2016~present: BPE / WordPiece

Context-independent ⇒ Context-aware

### **SemID Construction**

2023~2024:

RQ / PQ / Clustering /

LM-based Generator

# Language Tokenization

2014~2015: Word / Char 2016~present: BPE / WordPiece

Context-independent ⇒ Context-aware

### **SemID Construction**

2023~2024:

RQ / PQ / Clustering /

LM-based Generator

2025:

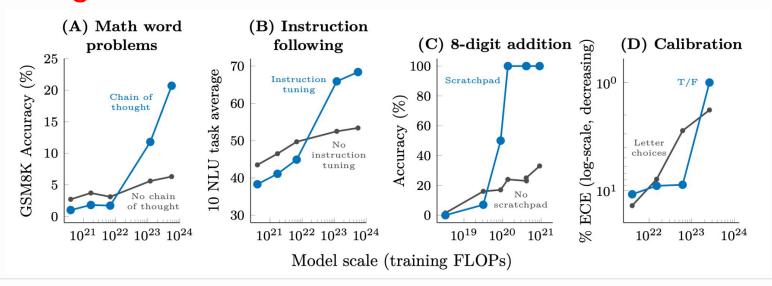
ActionPiece / Pctx / ?

# Part 2: What Becomes Possible?

What new opportunities may large generative models unlock for recommender systems?

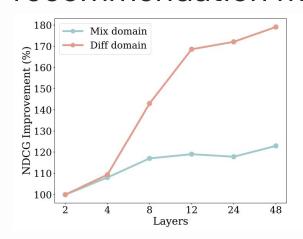
# **Emergent Ability**

Abilities not present in smaller models but is present in larger models

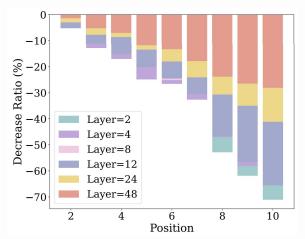


# **Emergent Ability**

Do we have emergent abilities in large generative recommendation models?

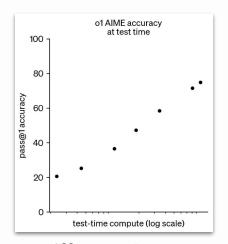


Cross-domain



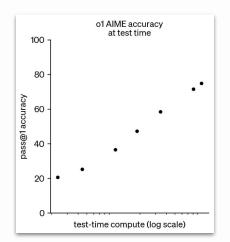
**Trajectory Prediction** 

There have been explorations on model / data scaling of recommendation models



Test-time scaling is still under exploration

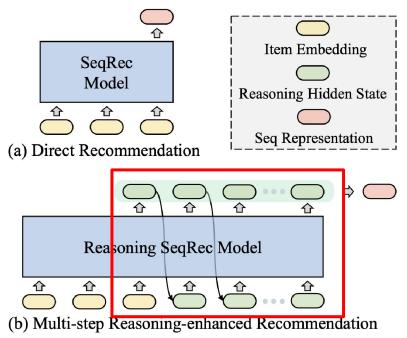
There have been explorations on model / data scaling of recommendation models



Test-time scaling is still under being actively exploration

Reasoning over latent hidden states to scale up test-time computation.





Reasoning over text tokens

### Recommendation in Conversation

我现在心情不好,请给我推一些轻松的内容 I'm not in a good mood right now, please recommend some relaxing content for me



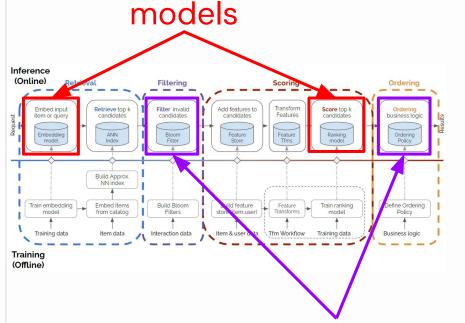


根据观看历史,用户可能对动作冒险类游戏或运动感兴趣。然而,由于用户明确表示希望看到温暖治愈的内容,因此需要避免推荐过于激烈或紧张的视频。相反,可以考虑推荐一些轻松愉快的游戏或动画,或者提供一些能够缓解压力的视频内容。最终,用户更倾向于选择能带来积极情绪体验的内容,而不是那些可能引发焦虑或不适的内容。我将优先考虑以下几类视频: 1. 动画电影中的温馨故事,如《冰雪奇缘》或《疯狂动物城》; 2. 宠物日常vlog,尤其是猫咪和狗狗的可爱瞬间; 3. 自然风光纪录片,比如北极熊的生存挑战或珊瑚礁生态。此外,我会排除任何涉及暴力、恐怖元龄内容,并增加正能量励志短片的比例。通过这样的调整,我相信可以更好地满足用户的情感需求。所以,用户可能点击视频《litem\_begin]~sitem\_a\_234~sitem\_b\_1262~sitem\_c\_941>< litem\_empl>

Based on viewing history, users may be interested in action-adventure games or sports content.

However, since users have explicitly expressed a desire to view warm and healing content, it is therefore necessary to avoid recommending overly intense or tense videos. Conversely, one might consider recommending some light-hearted and enjoyable games or animations, or providing video content that can alleviate stress. Ultimately, users are more inclined to select content that provides positive emotional experiences rather than content that may trigger anxiety or discomfort. I will prioritize the following categories of videos: 1. Heartwarming stories from animated films, such as *Frozen* or *Zootopia*; 2. Pet daily vlogs, especially adorable moments featuring cats and dogs; 3. Nature documentary films, such as the survival challenges of polar bears or coral reef ecosystems. Additionally, I will exclude any content involving violence or horror elements and increase the proportion of positive and inspirational short videos. Through such adjustments, I believe users' emotional needs can be better satisfied. Therefore, users may click on video

<|item\_begin|><item\_a\_234><item\_b\_1262><item\_c\_941><|item\_end|>. .....

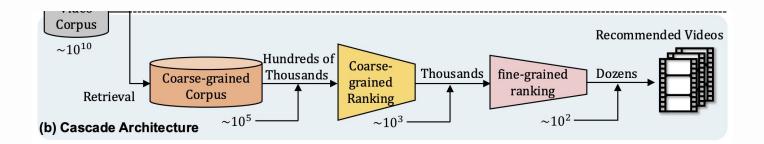


Complicated Architecture

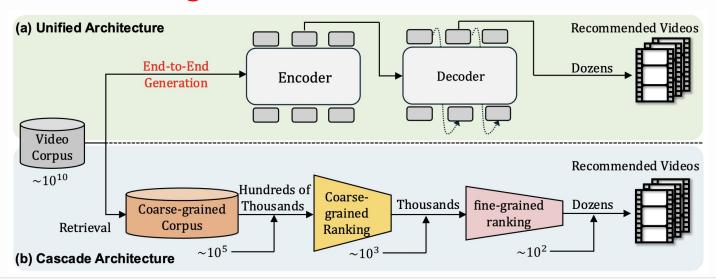
- Difficult to be optimized in an end-to-end way
- Latency between / within different modules

rules, strategies, heuristics

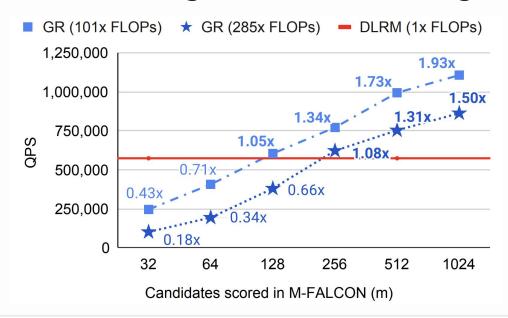
Is it possible to replace traditional cascade architecture

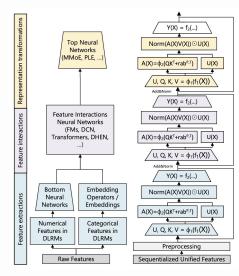


Is it possible to replace traditional cascade architecture with a **unified generative model**?



Better throughout when ranking more candidates









# Q&A

Thank you for coming!

Please kindly refer to

large-genrec.github.io

for slides, paper list, .....