

# MISSRec: Pre-training and Transferring Multi-modal Interest-aware Sequence Representation for Recommendation

Jinpeng Wang, Ziyun Zeng, Yunxiao Wang, Yuting Wang, Xingyu Lu, Tianxiang Li, Jun Yuan, Rui Zhang, Hai-Tao Zheng, and Shu-Tao Xia.  
Tsinghua Shenzhen International Graduate School, Tsinghua University    Huawei Noah's Ark Lab    Peng Cheng Laboratory  
Data & Code: <https://github.com/gimpong/MM23-MISSRec>    Contact: [wjp20@mails.tsinghua.edu.cn](mailto:wjp20@mails.tsinghua.edu.cn)



## Highlights

- ★ A multi-modal pre-training and transfer learning framework for sequential recommendation.
- ★ A clustering-based multi-modal interest discovery module.
- ★ An interest-aware transformer-based encoder-decoder model for multi-modal sequence representation.

## Introduction

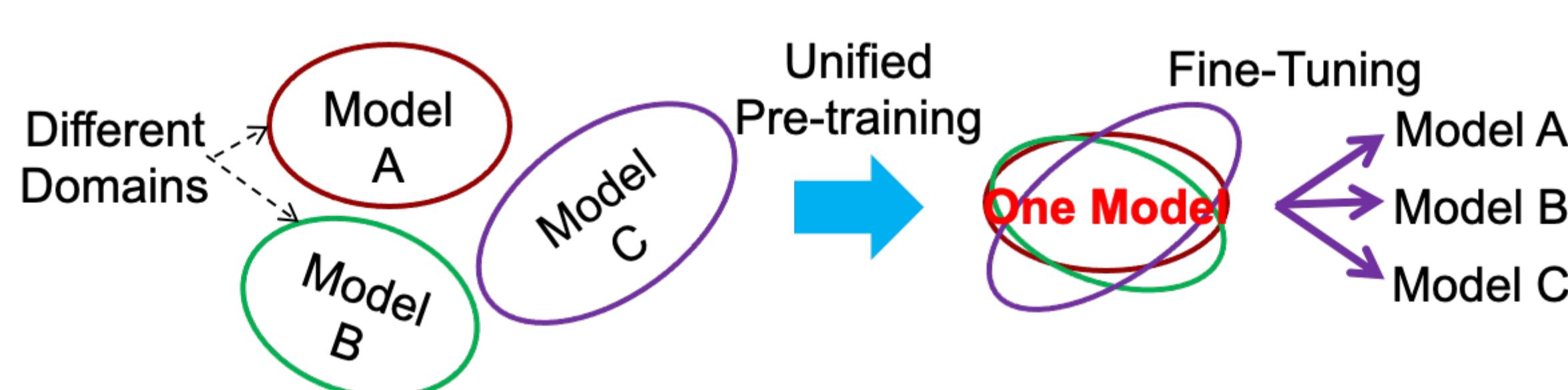
- **Sequential recommendation:** Predicting a user's potential interested items based on her / his historical interactions.



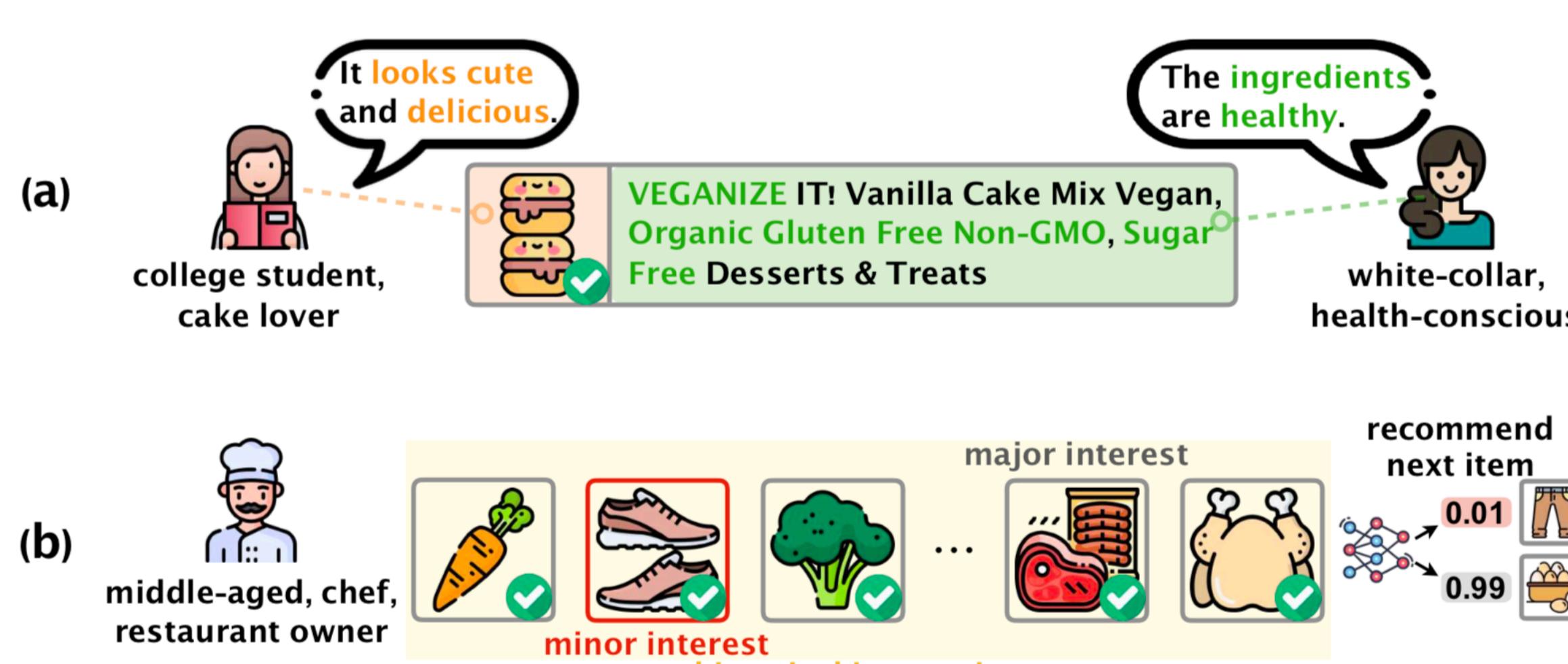
- Most existing SR approaches are developed based on ID features, despite the widespread application, suffers from
  - Under-fitting sparse IDs;
  - Cold-start problem;
  - Limited transferability across different domains.
- Multi-modal information (e.g., texts and images) is ubiquitous and plays an important role in attracting user attention.



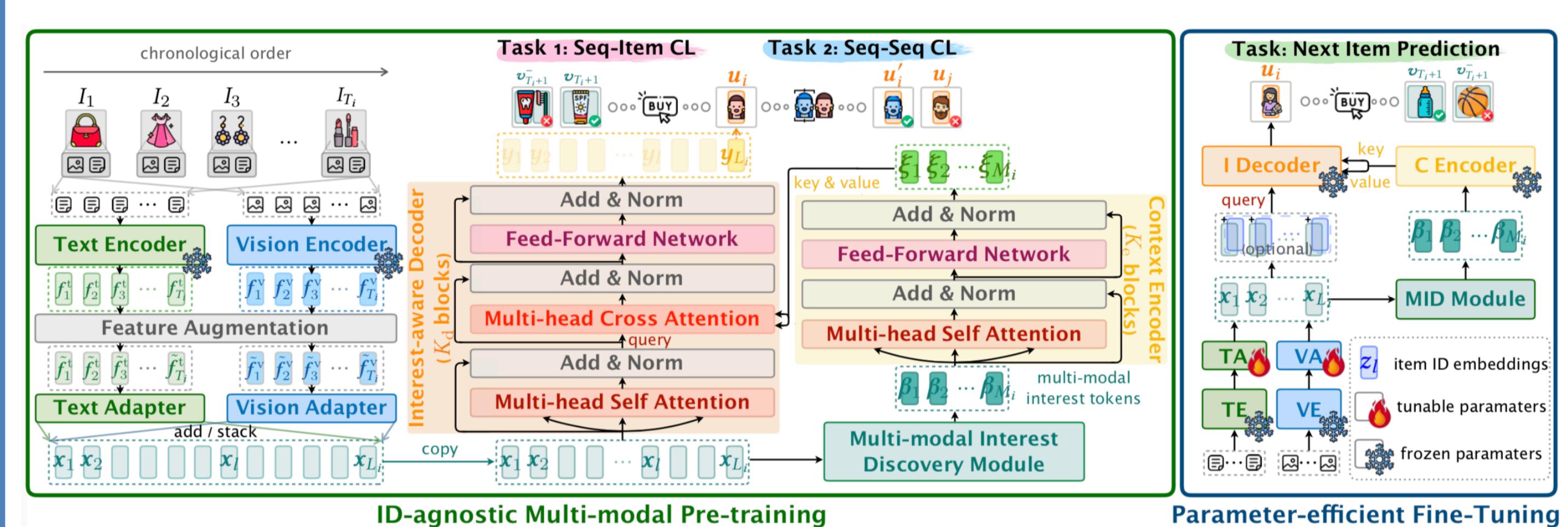
- Our solution
  - ◆ Use **multi-modal** features for universal representation.
  - ◆ **Pre-train** a unified model and **fine-tune** it for transferable recommendation.



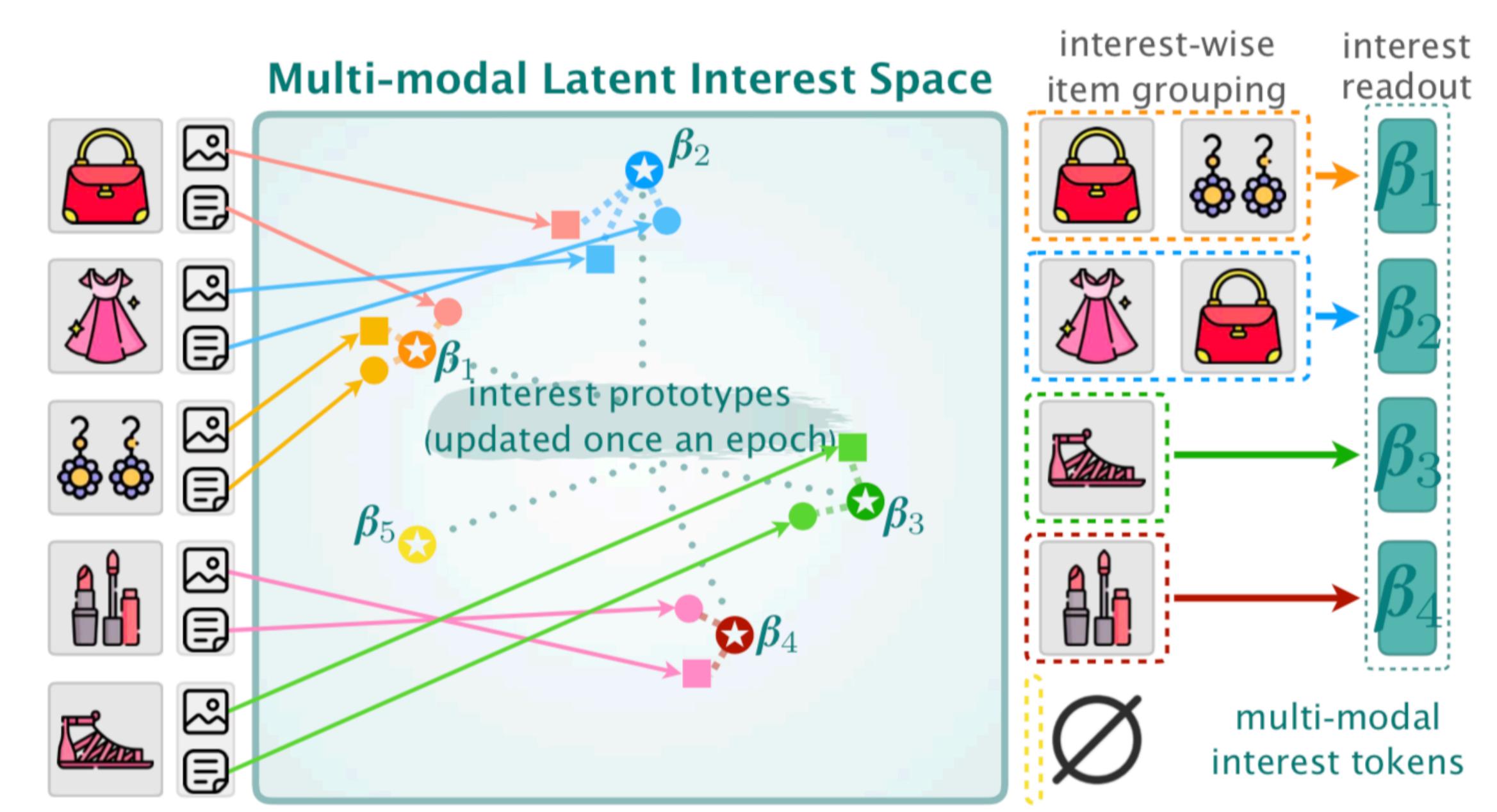
- Challenges
  - \* The multi-modal synergy is **dynamic** and **user-dependent** in the user-item interaction process.
  - \* **Long-tailed** proportion of user interests in the historical interaction.



## Our Approach (MISSRec)



- Two stages
  1. ID-agnostic multi-modal pre-training;
  2. **Parameter-efficient** fine-tuning in target domains.
- Multi-modal Interest Discovery (MID)
  - ◆ Clustering for index construction: ***k*-nearest neighbor-based density peaks clustering (DPC-KNN)** algorithm;
  - ◆ Interest discovery via indexing.



## Experiments

- Comparison with state-of-the-arts

Dataset	Input Type & Model →	ID					T+ID			T+V+ID			Improv. w/ ID	T			T+V		Improv. w/o ID				
		SASRec	BERT4Rec	FDSA	S <sup>3</sup> -Rec	UniSRec	MISSRec	SASRec	ZESRec	UniSRec	MISSRec	SASRec	ZESRec	UniSRec	MISSRec	SASRec	ZESRec						
Scientific	R@10	0.1080	0.0488	0.0899	0.0525	<b>0.1235</b>	<b>0.1360</b>	10.12%	0.0994	0.0851	0.1188	<b>0.1278</b>	7.58%										
	N@10	0.0553	0.0243	0.0580	0.0275	0.0634	<b>0.0753</b>	18.77%	0.0561	0.0475	0.0641	<b>0.0658</b>	2.65%										
	R@50	0.2042	0.1185	0.1732	0.1418	<b>0.2473</b>	<b>0.2431</b>		0.2162	0.1746	<b>0.2394</b>	<b>0.2375</b>											
	N@50	0.0760	0.0393	0.0759	0.0468	<b>0.0904</b>	<b>0.0983</b>		0.0815	0.0670	<b>0.0903</b>	<b>0.0893</b>											
Pantry	R@10	0.0501	0.0308	0.0395	0.0444	<b>0.0693</b>	<b>0.0779</b>	12.41%	0.0585	0.0454	0.0636	<b>0.0771</b>	21.23%										
	N@10	0.0218	0.0152	0.0209	0.0214	<b>0.0311</b>	<b>0.0365</b>	17.36%	0.0285	0.0230	0.0306	<b>0.0345</b>	12.75%										
	R@50	0.1322	0.1030	0.1151	0.1315	<b>0.1827</b>	<b>0.1875</b>	2.63%	0.1647	0.1141	0.1658	<b>0.1833</b>	10.55%										
	N@50	0.0394	0.0305	0.0370	0.0400	<b>0.0556</b>	<b>0.0598</b>	7.55%	0.0523	0.0378	<b>0.0527</b>	<b>0.0571</b>	8.35%										
Instruments	R@10	0.1118	0.0813	0.1070	0.1056	<b>0.1267</b>	<b>0.1300</b>	2.60%	0.1127	0.0783	0.1189	<b>0.1201</b>	1.01%										
	N@10	0.0612	0.0620	0.0796	0.0713	<b>0.0748</b>	<b>0.0843</b>	5.90%	0.0661	0.0497	0.0680	<b>0.0771</b>	13.38%										
	R@50	0.2106	0.1454	0.1890	0.1927	<b>0.2387</b>	<b>0.2370</b>		0.2104	0.1387	<b>0.2255</b>	<b>0.2218</b>											
	N@50	0.0826	0.0756	0.0972	0.0901	<b>0.0991</b>	<b>0.1071</b>		0.0873	0.0627	0.0912	<b>0.0988</b>	8.33%										
Arts	R@10	0.1108	0.0722	0.1002	0.1003	<b>0.1239</b>	<b>0.1314</b>	6.05%	0.0977	0.0664	0.1066	<b>0.1119</b>	4.97%										
	N@10	0.0587	0.0479	0.0714	0.0601	<b>0.0712</b>	<b>0.0767</b>	7.42%	0.0562	0.0375	0.0586	<b>0.0625</b>	6.66%										
	R@50	0.2030	0.1367	0.1779	0.1888	<b>0.2347</b>	<b>0.2410</b>	2.68%	0.1916	0.1323	0.2049	<b>0.2100</b>	2.49%										
	N@50	0.0788	0.0619	0.0883	0.0793	<b>0.0955</b>	<b>0.1002</b>	4.92%	0.0766	0.0518	0.0799	<b>0.0836</b>	4.63%										
Office	R@10	0.1056	0.0825	0.1118	0.1030	<b>0.1280</b>	<b>0.1275</b>		-	0.0929	0.0641	0.1013	<b>0.1038</b>	2.47%									
	N@10	0.0710	0.0634	<b>0.0868</b>	0.0653	0.0831	<b>0.0856</b>		-	0.0582	0.0391	0.0619	<b>0.0666</b>	7.59%									
	R@50	0.1627	0.1227	0.1665	0.1613	<b>0.2016</b>	<b>0.2005</b>		-	0.1580	0.1113	<b>0.1702</b>	<b>0.1701</b>	-									
	N@50	0.0835	0.0721	0.0987	0.0780	<b>0.0991</b>	<b>0.1012</b>		2.12%	0.0723	0.0493	0.0769	<b>0.0808</b>	5.07%									

- Influence of multi-modal information

Input Type	Variant	Scientific				Office				T	T+V	Improv.
		R@10	R@50	N@10	N@50	R@10	R@50	N@10	N@50			
ID+T	UniSRec	0.1676	0.3200	0.0872	0.1208	0.1376	0.2152	0.0937	0.1104			