

# Large-scale Generative and Multimodal Recommendation Systems: An Overview

Junwei Pan

**Tencent** 腾讯



# Outline

## Generative Recommendation

- Token Organization
- Action Handling

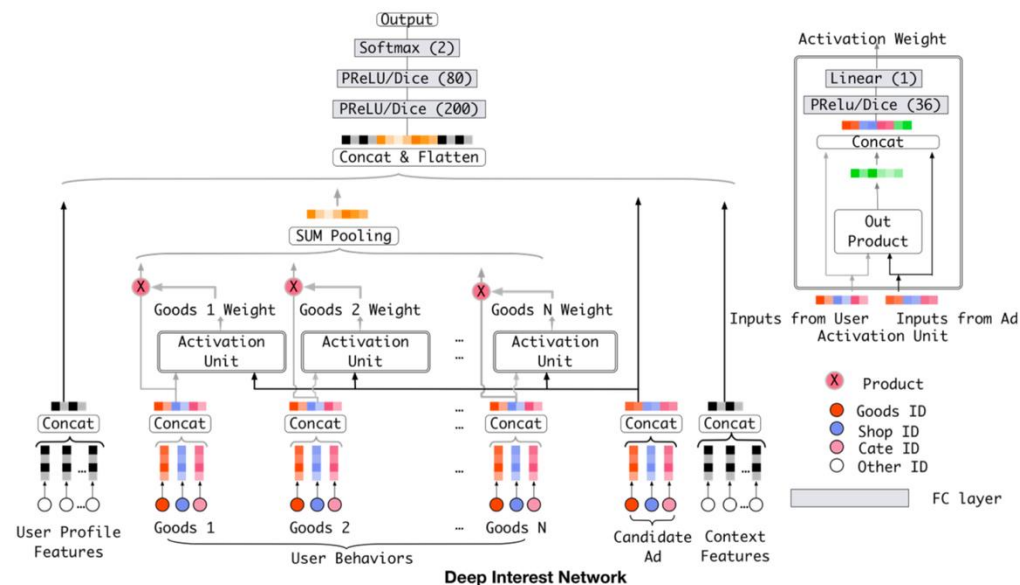
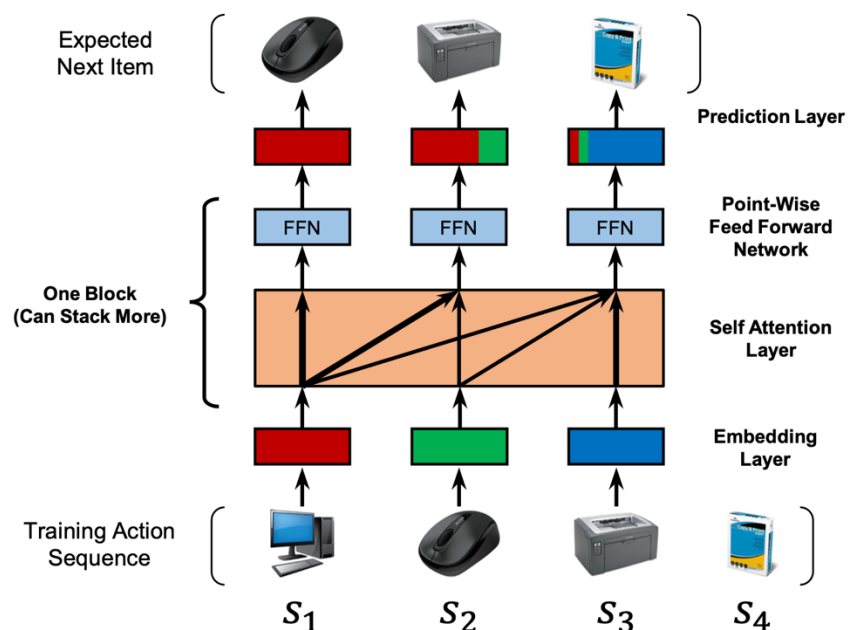
## Multimodal Recommendation

- Alignment
- Distance Transfer
- Semantic IDs



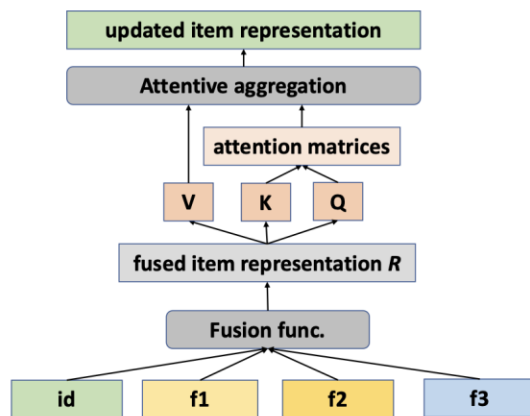
# Item-Oriented: SASRec, DIN

- Next-item prediction paradigm
- Fuse all side infos

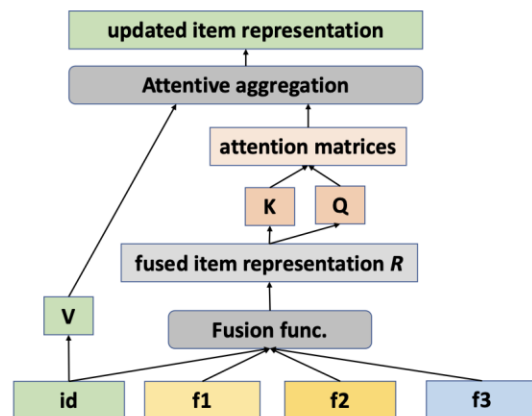


# Side Infos

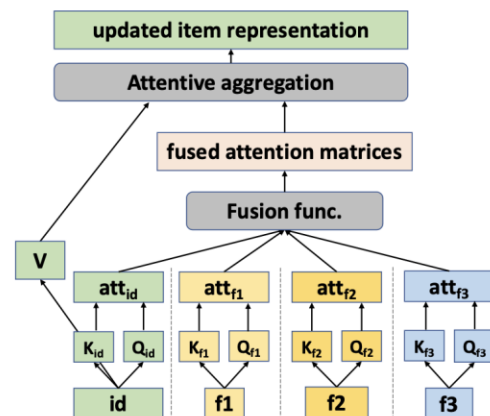
- SASRec fuse all side infos
- NOVA employs only the item ID in the **Value**
- DIF-SR further decouple each feature in the **Attention (Q, K)**
- DSI-TIN decouples feature groups in the **Attention**



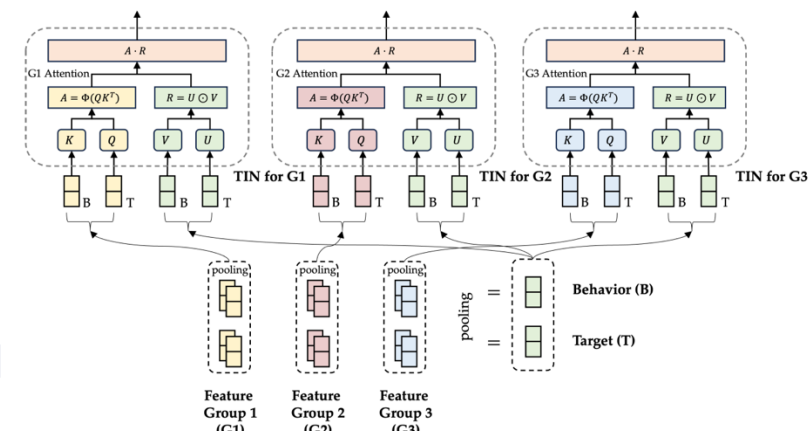
(a) SASRec.



(b) NOVA-SR.

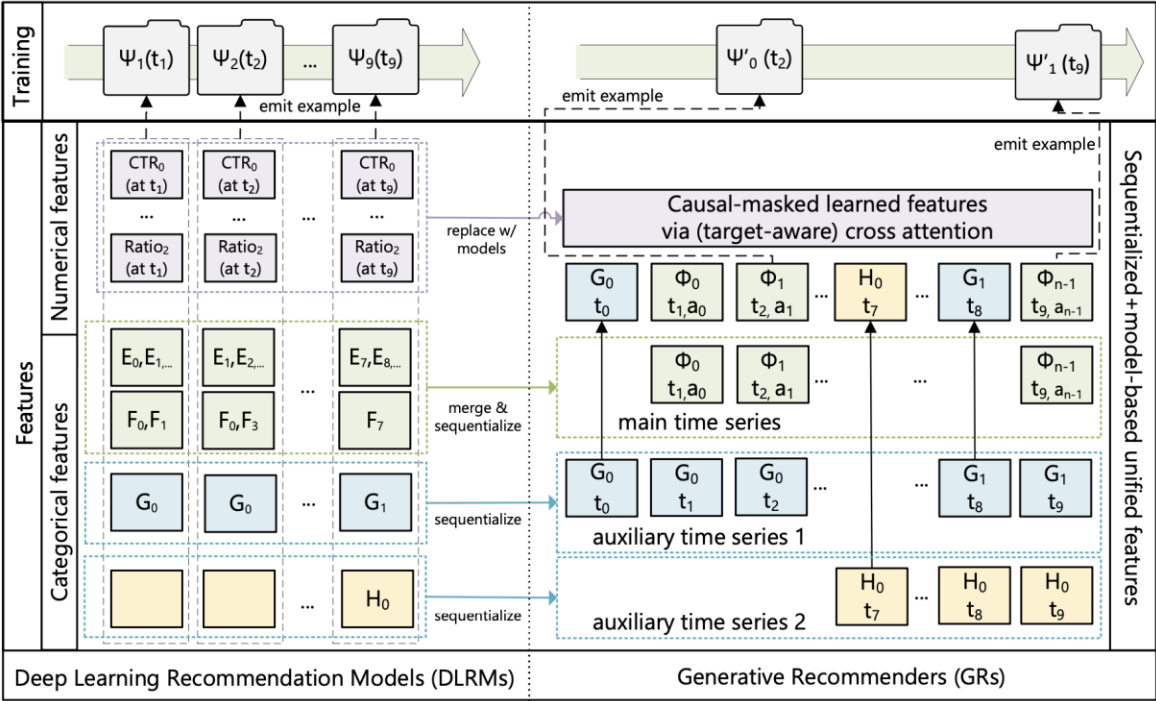


(c) DIF-SR.



# HSTU

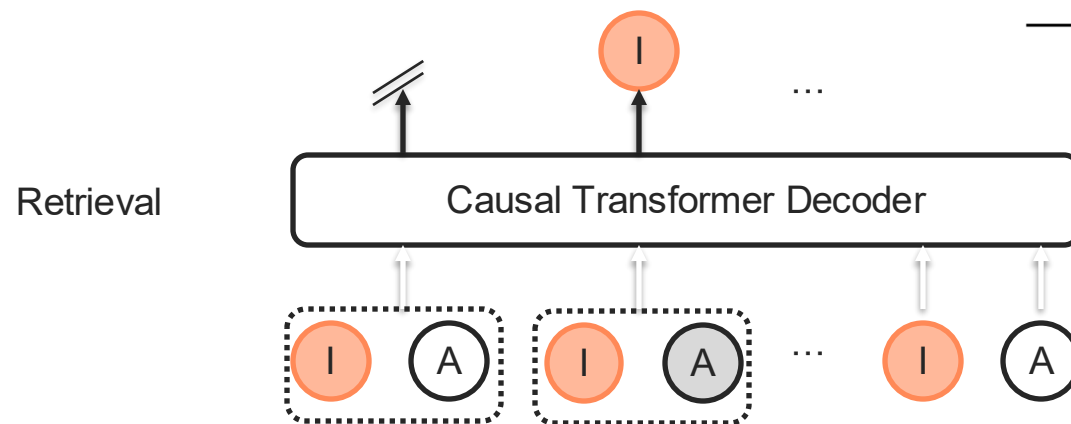
- HSTU first selects the longest time series, typically by merging the features that represent items user engaged with, i.e., item ID, as the **main time series**
- It then *compress* the remaining time series by keeping the earliest entry per consecutive segment and then merge the results into the main time series
- **Remove numerical features**, relying on **the target attention** mechanism to handle them



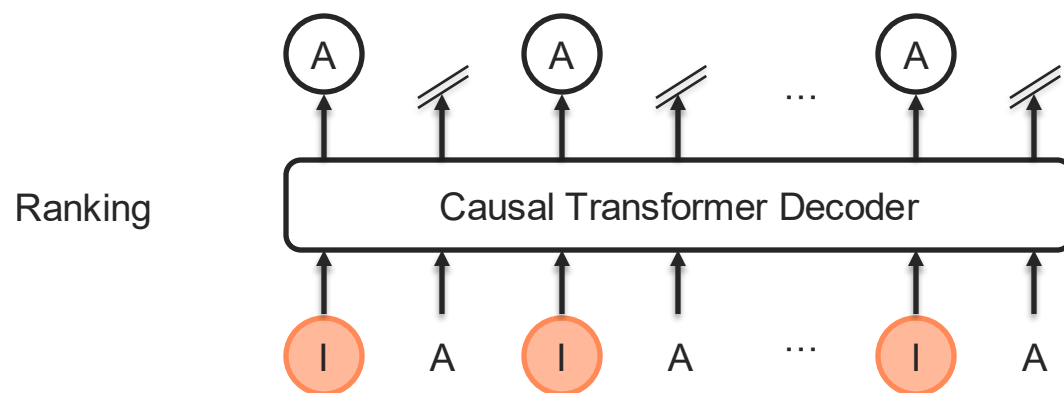
Task		Specification (Inputs / Outputs)
Ranking	$x_i$ s	$\Phi_0, a_0, \Phi_1, a_1, \dots, \Phi_{n_c-1}, a_{n_c-1}$
	$y_i$ s	$a_0, \emptyset, a_1, \emptyset, \dots, a_{n_c-1}, \emptyset$
Retrieval	$x_i$ s	$(\Phi_0, a_0), (\Phi_1, a_1), \dots, (\Phi_{n_c-1}, a_{n_c-1})$
	$y_i$ s	$\Phi'_1, \Phi'_2, \dots, \Phi'_{n_c-1}, \emptyset$ ( $\Phi'_i = \Phi_i$ if $a_i$ is positive, otherwise $\emptyset$ )

# HSTU

Task		Specification (Inputs / Outputs)
Ranking	$x_i$ s	$\Phi_0, a_0, \Phi_1, a_1, \dots, \Phi_{n_c-1}, a_{n_c-1}$
	$y_i$ s	$a_0, \emptyset, a_1, \emptyset, \dots, a_{n_c-1}, \emptyset$
Retrieval	$x_i$ s	$(\Phi_0, a_0), (\Phi_1, a_1), \dots, (\Phi_{n_c-1}, a_{n_c-1})$
	$y_i$ s	$\Phi'_1, \Phi'_2, \dots, \Phi'_{n_c-1}, \emptyset$ ( $\Phi'_i = \Phi_i$ if $a_i$ is positive, otherwise $\emptyset$ )



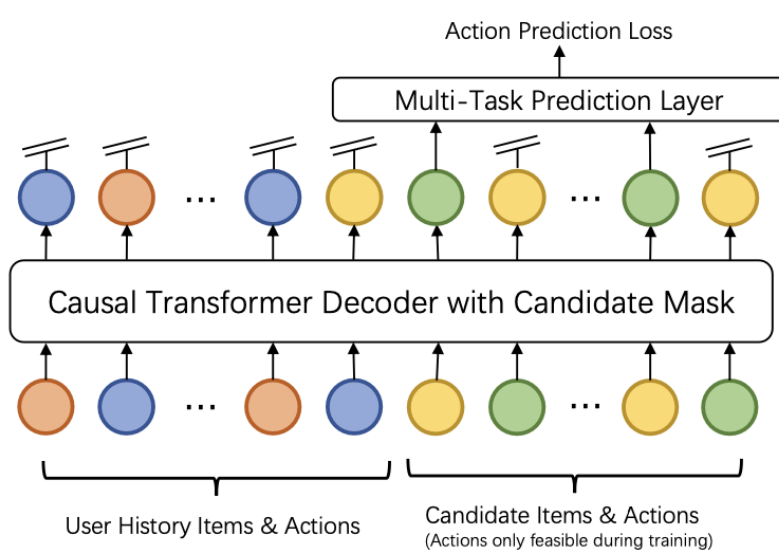
- If there is **positive** feedback on the next item, then predict this item; otherwise, predict **empty**.



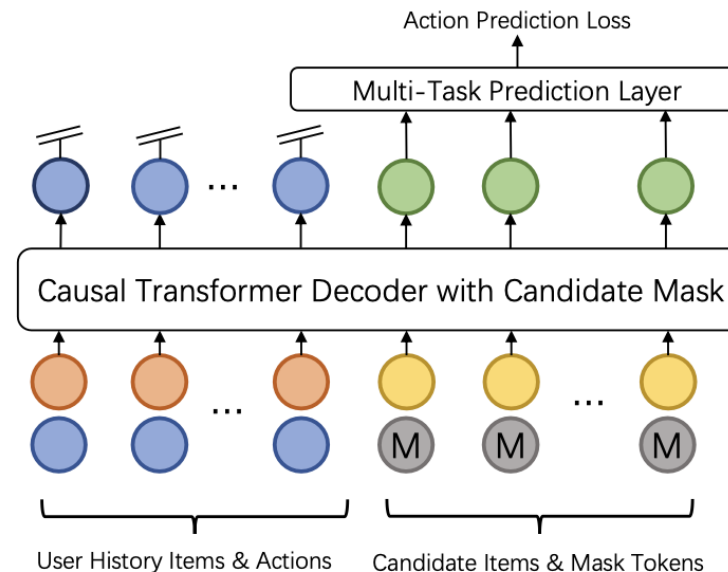
- Interleave** items and actions in the sequence. Predict the action for each item, and predict **empty** for the action.

# GenRank

- **DON'T interleave** items and actions in the sequence, but treat actions as side info, as done in traditional methods like SASRec or DIN for the history behaviors
- For the target, predict the action of each item
- **Left side (history): SASRec; right side (target): HSTU**



(a) Existing Approach with Item-Oriented Organization



(b) Our Approach with Action-Oriented Organization

- Legend:
- History Item Embedding (orange circle)
  - History Action Embedding (blue circle)
  - Candidate Item Embedding (yellow circle)
  - Candidate Action Embedding (green circle)
  - Mask Embedding (grey circle with 'M')
  - Stop Gradient Operation (double slashes)

# Overview

Generative Recommendation

- Token Organization
- Action Handling

Multimodal Recommendation

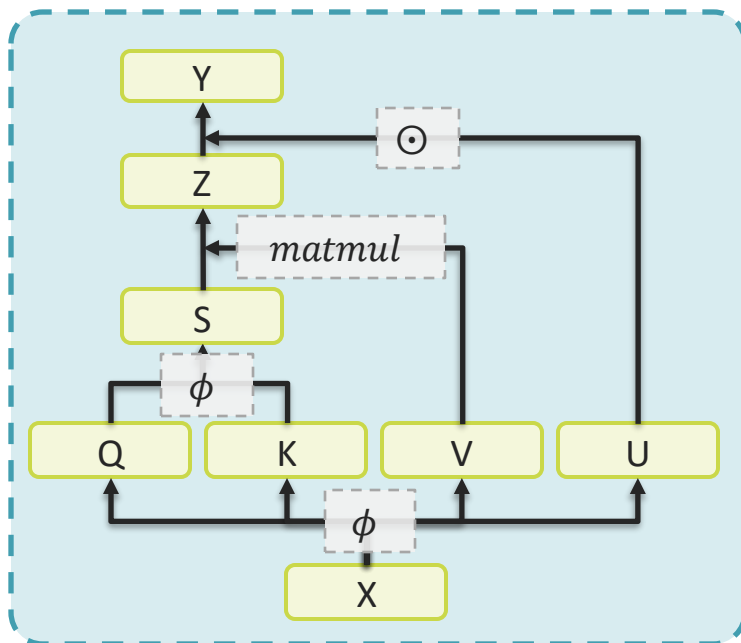
- Alignment
- Distance Transfer
- Semantic IDs





# HSTU

- DNN-based methods **employs MLPs**  $f_{MLP}(b_i, t)$  to learn the interaction between behaviors and the target
- Employs a SiLU or SwiGLU activation function, due to “**the difficulty of approximating dot products with learned MLPs**”



Task		Specification (Inputs / Outputs)
Ranking	$x_i s$	$\Phi_0, a_0, \Phi_1, a_1, \dots, \Phi_{n_c-1}, a_{n_c-1}$
	$y_i s$	$a_0, \emptyset, a_1, \emptyset, \dots, a_{n_c-1}, \emptyset$
Retrieval	$x_i s$	$(\Phi_0, a_0), (\Phi_1, a_1), \dots, (\Phi_{n_c-1}, a_{n_c-1})$
	$y_i s$	$\Phi'_1, \Phi'_2, \dots, \Phi'_{n_c-1}, \emptyset$ $(\Phi'_i = \Phi_i \text{ if } a_i \text{ is positive, otherwise } \emptyset)$

$$U(X), V(X), Q(X), K(X) = \text{Split}(\phi_1(f_1(X)))$$

$$A(X)V(X) = \phi_2 \left( Q(X)K(X)^T + \text{rab}^{p,t} \right) V(X)$$

$$Y(X) = f_2 (\text{Norm} (A(X)V(X)) \odot U(X))$$

# Temporal Interest Network (TIN)

- Target-aware Temporal Encoding
- Target-aware Attention
- Target-aware Representation

$$u_{TIN} = \alpha(Q, K) \odot (U \odot V)$$

$$= \alpha(\tilde{v}_t W_Q, \tilde{e}_i Q_K) \cdot (\tilde{v}_t W_U \odot \tilde{e}_i W_V)$$

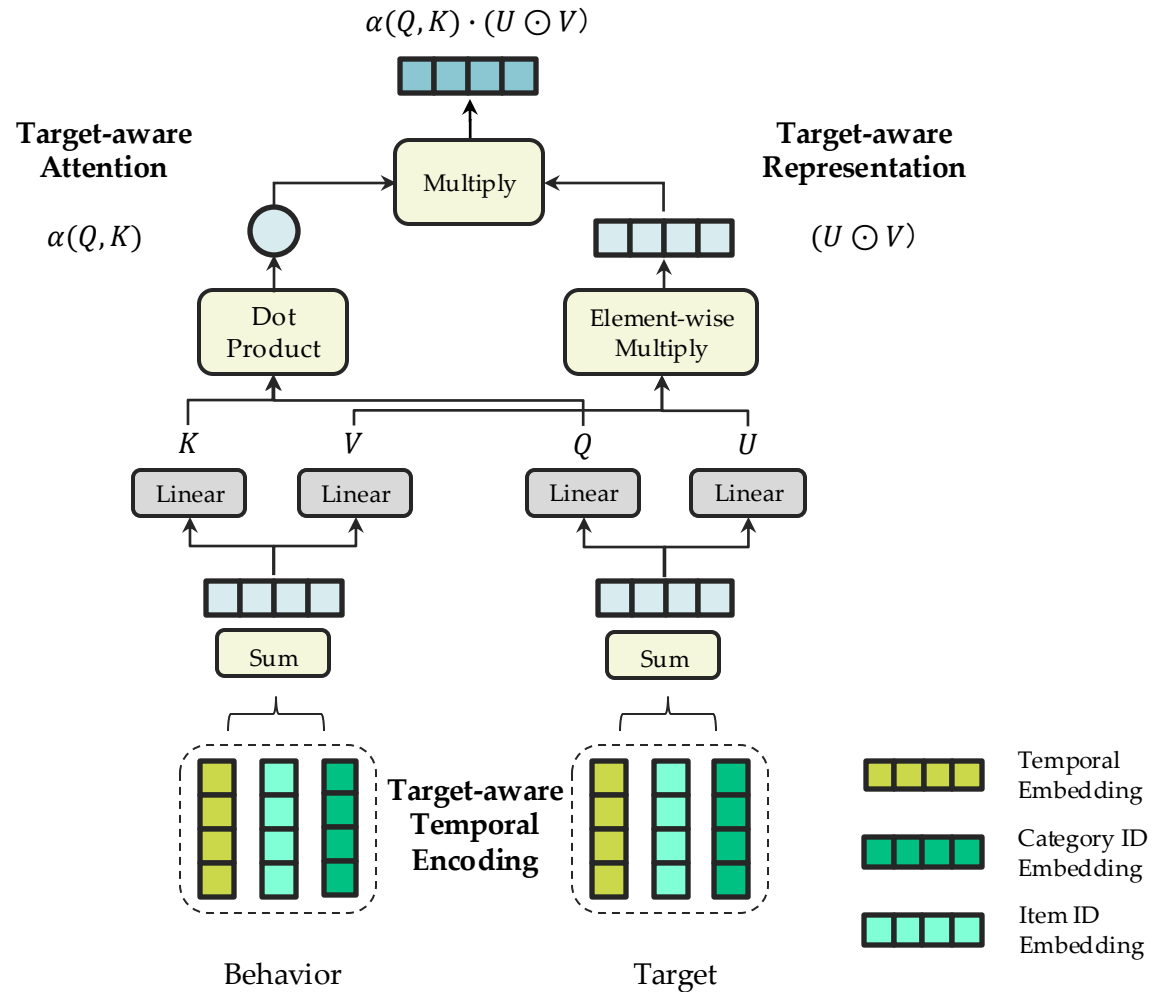
Formulation of TIN

$$U(X), V(X), Q(X), K(X) = \text{Split}(\phi_1(f_1(X)))$$

$$A(X)V(X) = \phi_2 \left( Q(X)K(X)^T + \text{rab}^{p,t} \right) V(X)$$

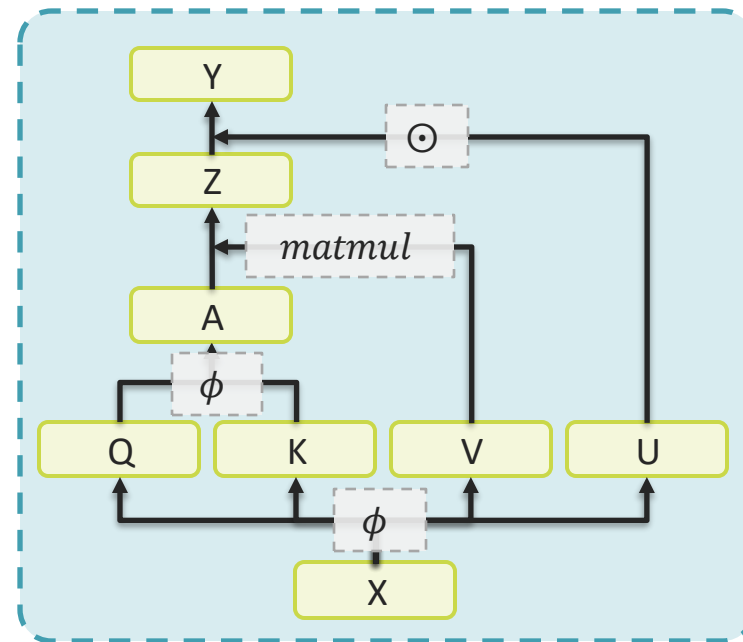
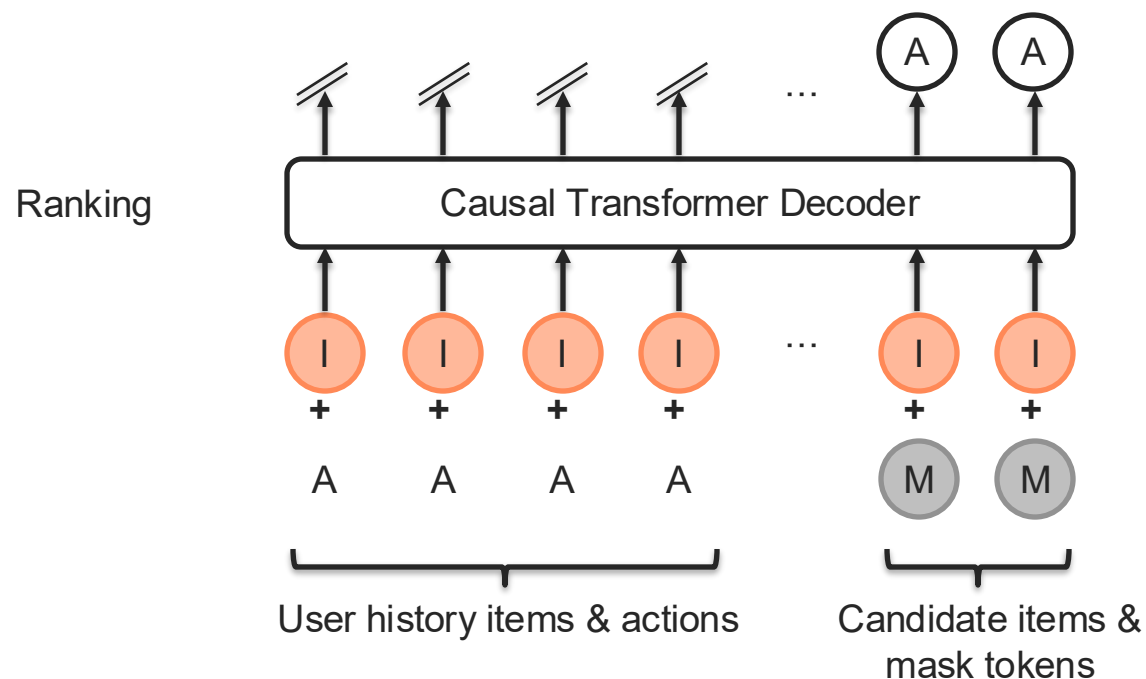
$$Y(X) = f_2(\text{Norm}(A(X)V(X)) \odot U(X))$$

Formulation of HSTU



# GenRank

- Similar architecture with HSTU: SwiGLU activation function



$$X = I + A \text{ or } X = I + M$$

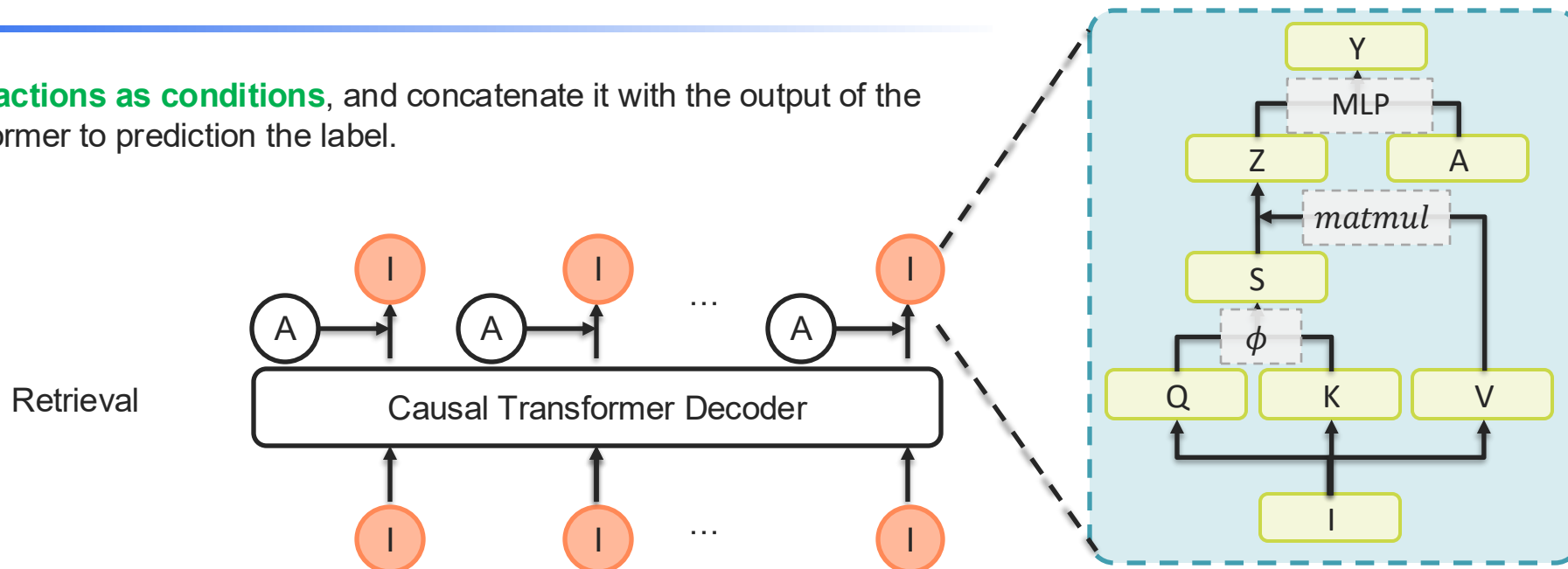
$$Q = \phi(W^Q X), K = \phi(W^K X), V = \phi(W^V X), U = \phi(W^U X)$$

$$Z = \phi(QK^T + rab)V$$

$$Y = W^Y(\text{Norm}(Z) \odot U)$$

# PinRec

- **Treat actions as conditions**, and concatenate it with the output of the transformer to prediction the label.



$$Q = \phi(W^Q I), K = \phi(W^K I), V = \phi(W^V I)$$

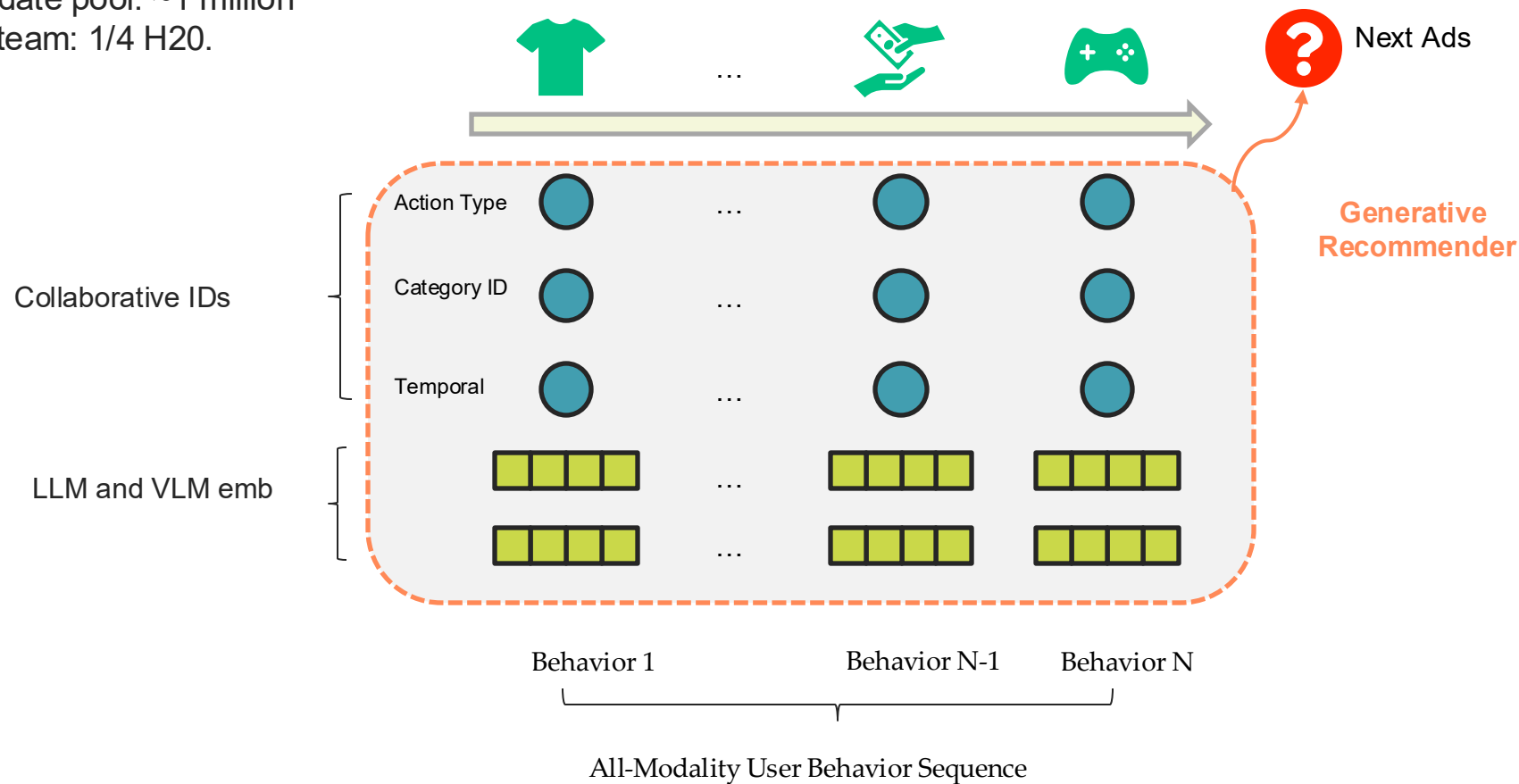
$$Z = \phi(QK^T)V \quad Y = \phi(W^Y([Z; A]))$$

# Tencent Advertising Algorithm Competition

## All-Modality Generative Recommendation

### Setting

- 1 millions of user sequences
- Candidate pool: ~1 million
- Each team: 1/4 H20.

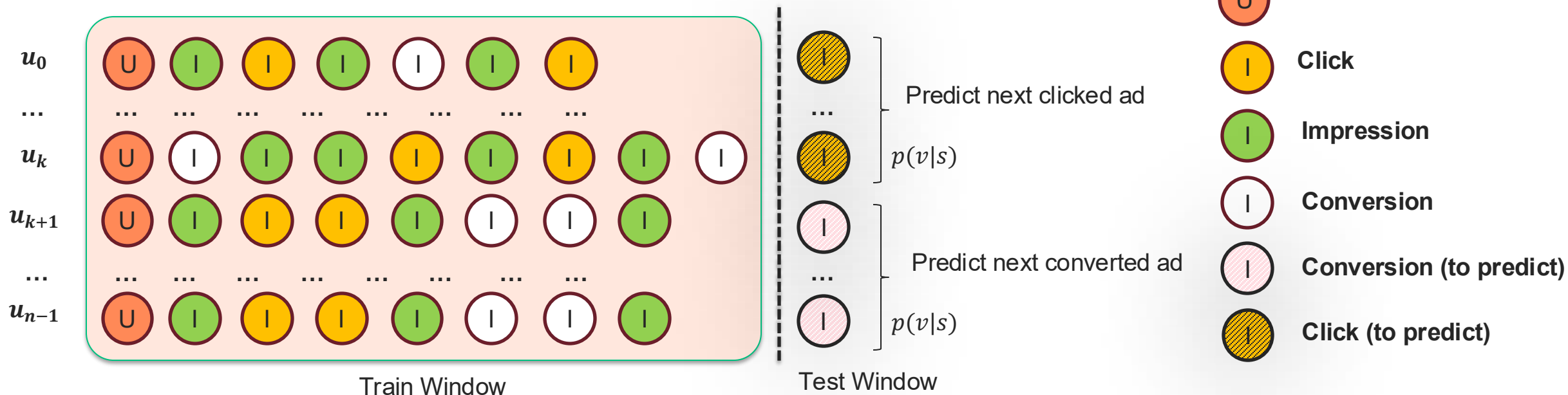


# Tencent Advertising Algorithm Competition

## All-Modality Generative Recommendation

### Setting

- **10 millions** of user sequences
- Include **conversions** in both features and target
- Each team: **7 H20 GPU**
- In evaluation, conversions have larger scores (2.5x) than the clicks





# TAAC

## AWARD

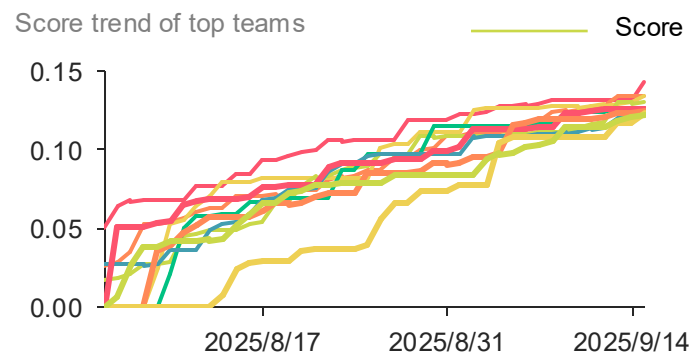
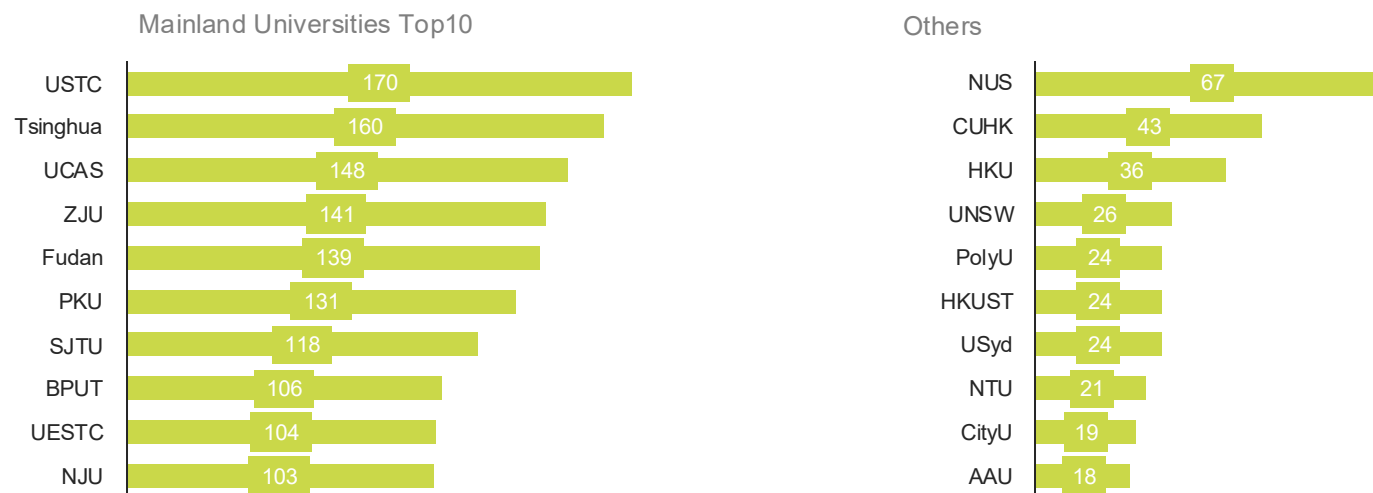
Total Prize Pool of 3.6 Million RMB Ready!

Note: Based on comprehensive evaluation by the competition, the award may go unassigned. Disbursement of the prize is subject to the rules established post-award.

			
<b>Champion</b> 1 team	<b>Runner-Up</b> 1 team	<b>Second Runner-Up</b> 1 team	<b>4<sup>th</sup>-10<sup>th</sup> Places</b> per team
¥ <b>2,000,000</b> RMB	¥ <b>600,000</b> RMB	¥ <b>300,000</b> RMB	¥ <b>100,000</b> RMB

# TAAC

8400+      4600+      2800+  
Registration      Enroll      #Teams





# Overview

## Generative Recommendation

- Token Organization
- Action Handling

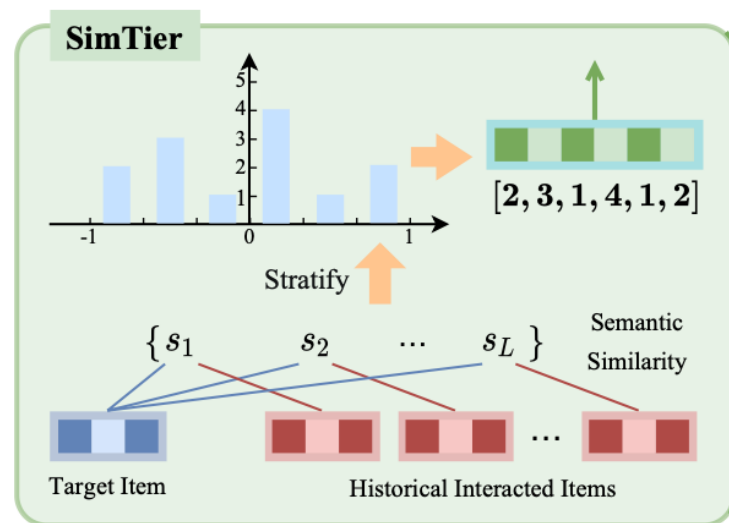
## Multimodal Recommendation

- Distance Transfer
- Alignment
- Semantic IDs

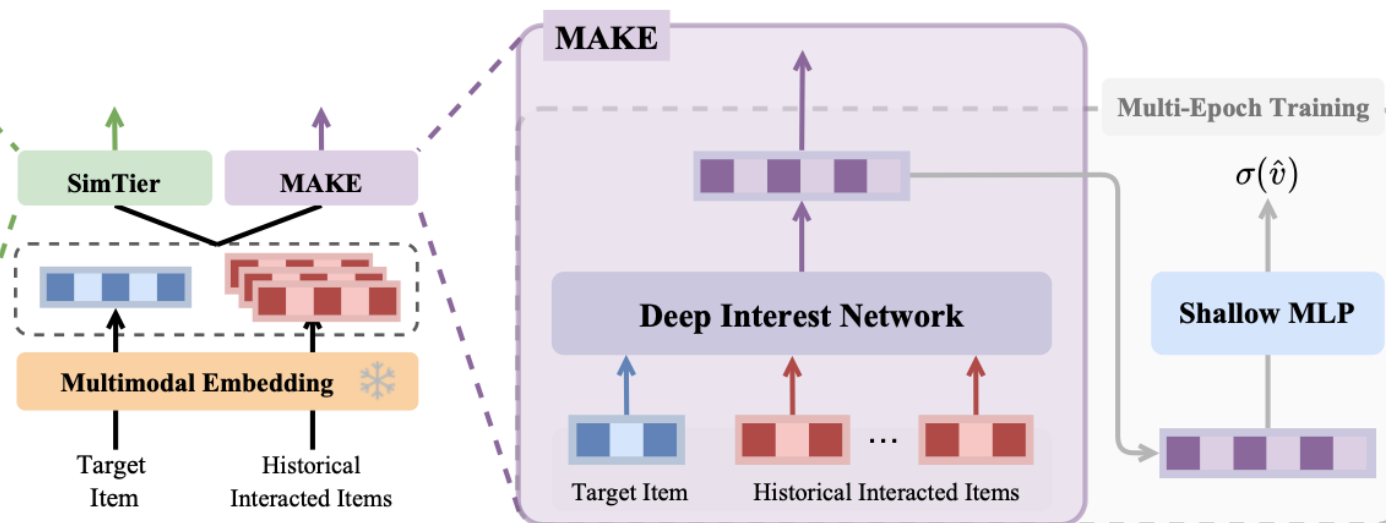


# SimTier

- Calculate the **similarity score** between multi-modal representation of behaviors and the target
- **Histogram** of the scores in a N pre-defined buckets
- Use the N-dimension vector as a new representation



(a). SimTier



(b). MAKE

# MNSE

- Calculate the **distance** between two features based on the source embedding
- Encoding the distance with **n-ary**
- Train the encoded embeddings in the target task, together with other embeddings in the target space

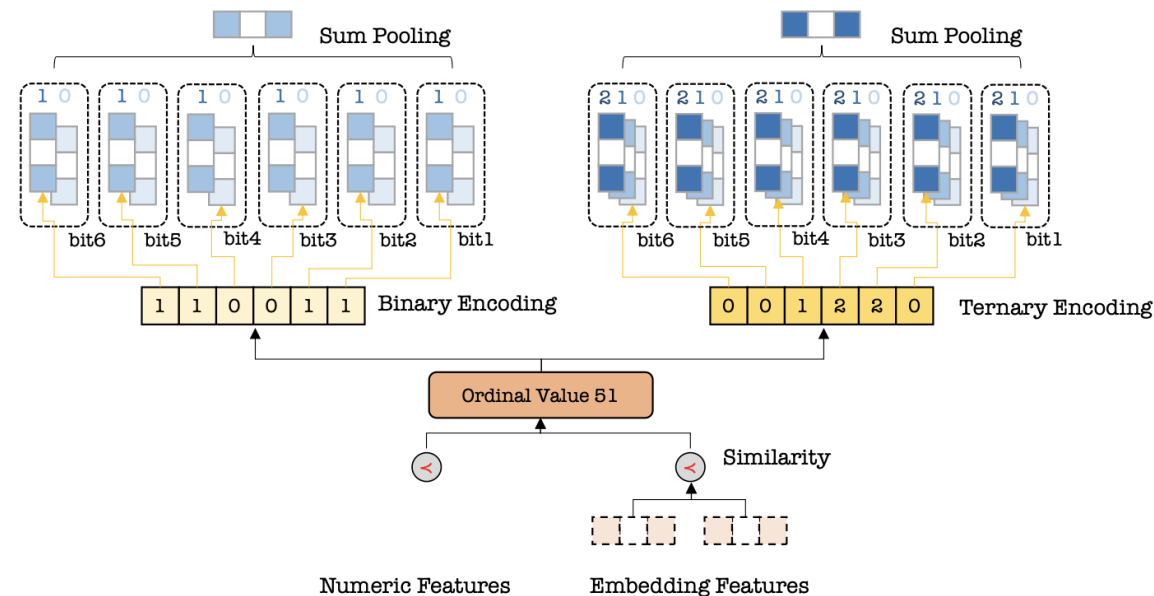
$$f_{\text{MNS}}(v) = \left[ \sum_{k=1}^{K_2} \mathbf{X}_{2k+\mathbb{B}_k}^{(2)}, \sum_{k=1}^{K_3} \mathbf{X}_{3k+\mathbb{C}_k}^{(3)}, \dots, \sum_{k=1}^{K_n} \mathbf{X}_{nk+\mathbb{N}_k}^{(n)} \right]$$

$\mathbb{B} = \text{func\_binary}(v)$ ,  $\mathbb{C} = \text{func\_ternary}(v), \dots$

Numerical Feature (Decimal)	Binary	
45	0000101101	<div>Continuity</div> <div>Discriminability</div>
46	0000101110	
957	1110111101	

Numerical Feature (Decimal)	Binary	Ternary	
63	0000111111	000002100	<div>Carry</div> <div>Modulo</div> <div>Continuity</div> <div>Discriminability</div>
64	0001000000	000002101	
575	1000111111	000210022	



Key properties of n-ary: continuity, discriminability

# Overview

## Generative Recommendation

- Token Organization
- Action Handling

## Multimodal Recommendation

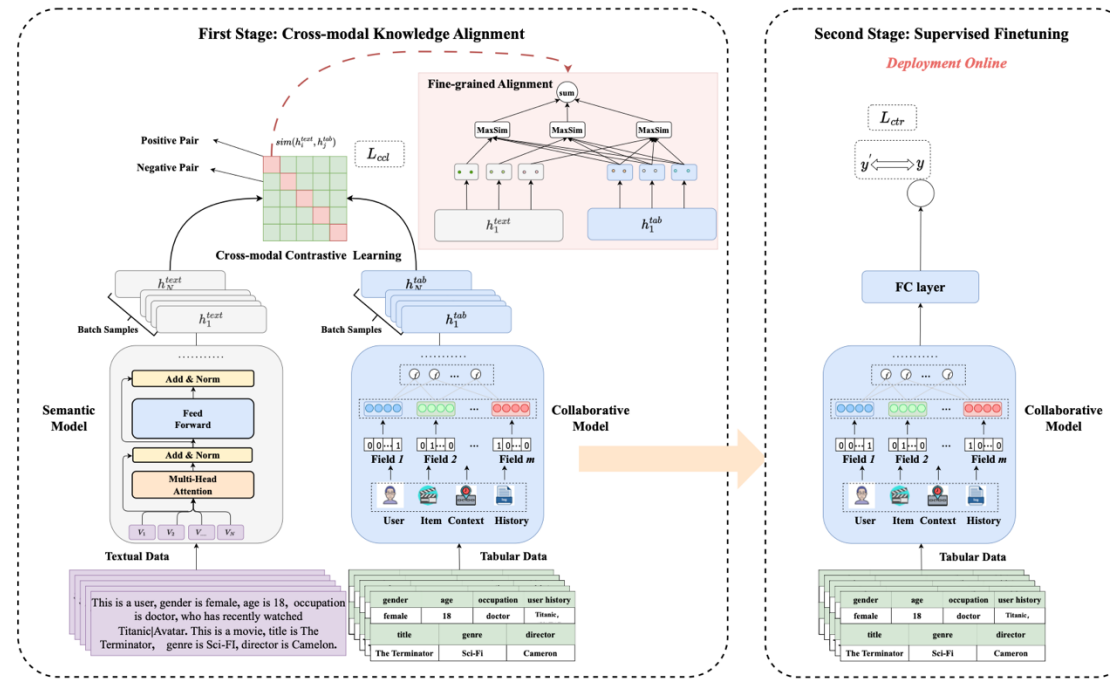
- Distance Transfer
- Alignment
- Semantic IDs



# CTRL

- Clip-like Textual-to-Tabular **contrastive loss**

$$\mathcal{L}^{textual2tabular} = -\frac{1}{N} \sum_{k=1}^N \log \frac{\exp(\text{sim}(\mathbf{h}_k^{text}, \mathbf{h}_k^{tab})/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{h}_k^{text}, \mathbf{h}_j^{tab})/\tau)},$$



# PAD (Pre-train, Align and Disentangle)

- Adopt **MK-MMD** (multi-kernel maximum mean discrepancy) as the alignment loss to capture **all information about the distribution**
- Combine the alignment with the BCE loss to avoid **catastrophic forgetting**

Avoid **catastrophic forgetting**

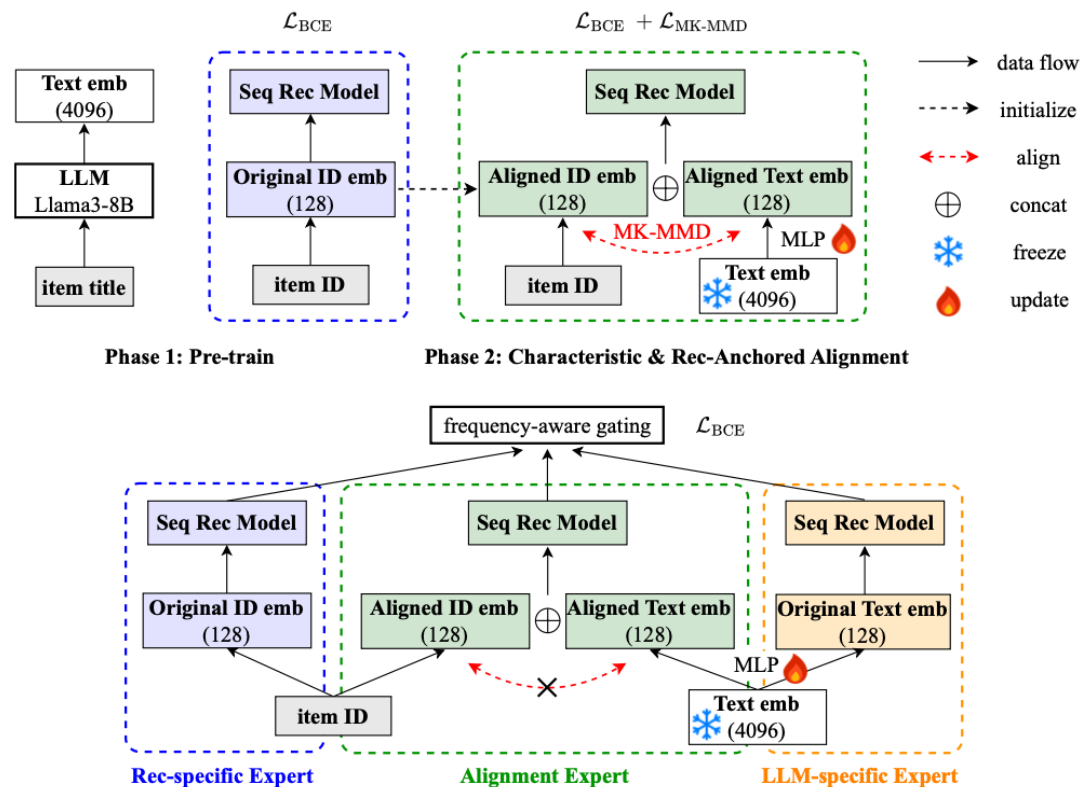
**Space Alignment**

$$\mathcal{L} = \mathcal{L}_{\text{REC}} + \gamma \cdot \mathcal{L}_{\text{MK-MMD}}$$

$$\mathcal{L}_{\text{MK-MMD}} = D_k^2(f_{\text{MLP}}(\text{SG}(\mathcal{D}_{\text{text}}), \mathbf{w}), \mathcal{D}_{\text{rec}})$$

$$\mathcal{L}_{\text{REC}} = \frac{1}{n} \sum_{i=1}^n \text{BCE}(f_{\theta}(\{\mathbf{h}_i^s\}, \{\mathbf{h}_i^c\}, \mathbf{x}_i^s, \mathbf{x}_i^c), y_i)$$

$$\text{MK-MMD}^2(X_s, X_t) = \left\| \frac{1}{n} \sum_{i=1}^n \phi_k(x_s^i) - \frac{1}{m} \sum_{j=1}^m \phi_k(x_t^j) \right\|_{\mathcal{H}_k}^2$$



# Overview

## Generative Recommendation

- Token Organization
- Action Handling

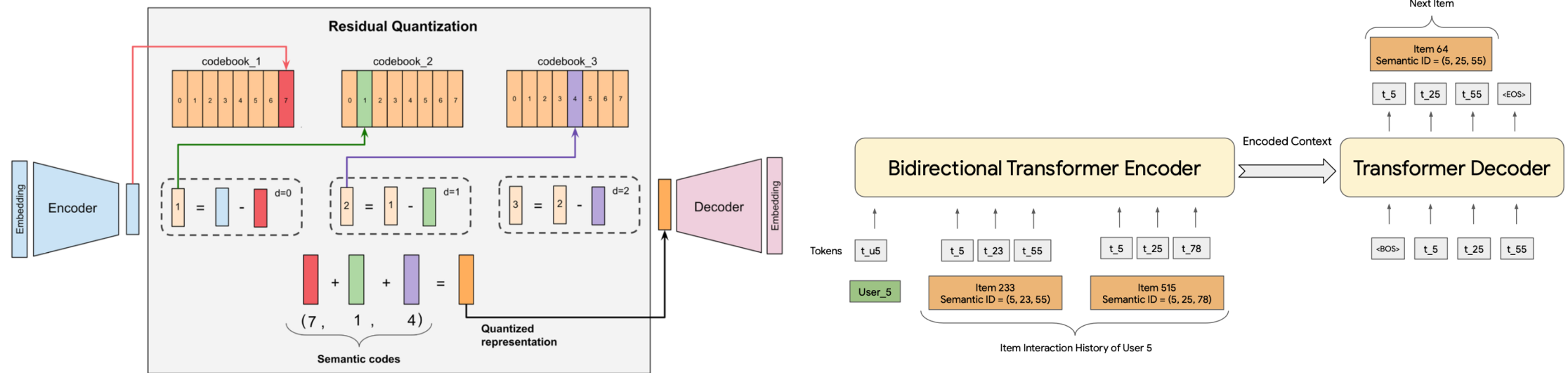
## Multimodal Recommendation

- Distance Transfer
- Alignment
- Semantic IDs



# Tiger

- Use **RQ-VAE** on the LLM representation to get **semantic IDs**
- Use semantic IDs in the downstream recommendation models



Why it works?

- **RQ-VAE** and **Semantic IDs** to capture the **source** space structures.
- **Semantic ID Embeddings** to align with the **target** space.



CIKM 2025  
November 10 - 14

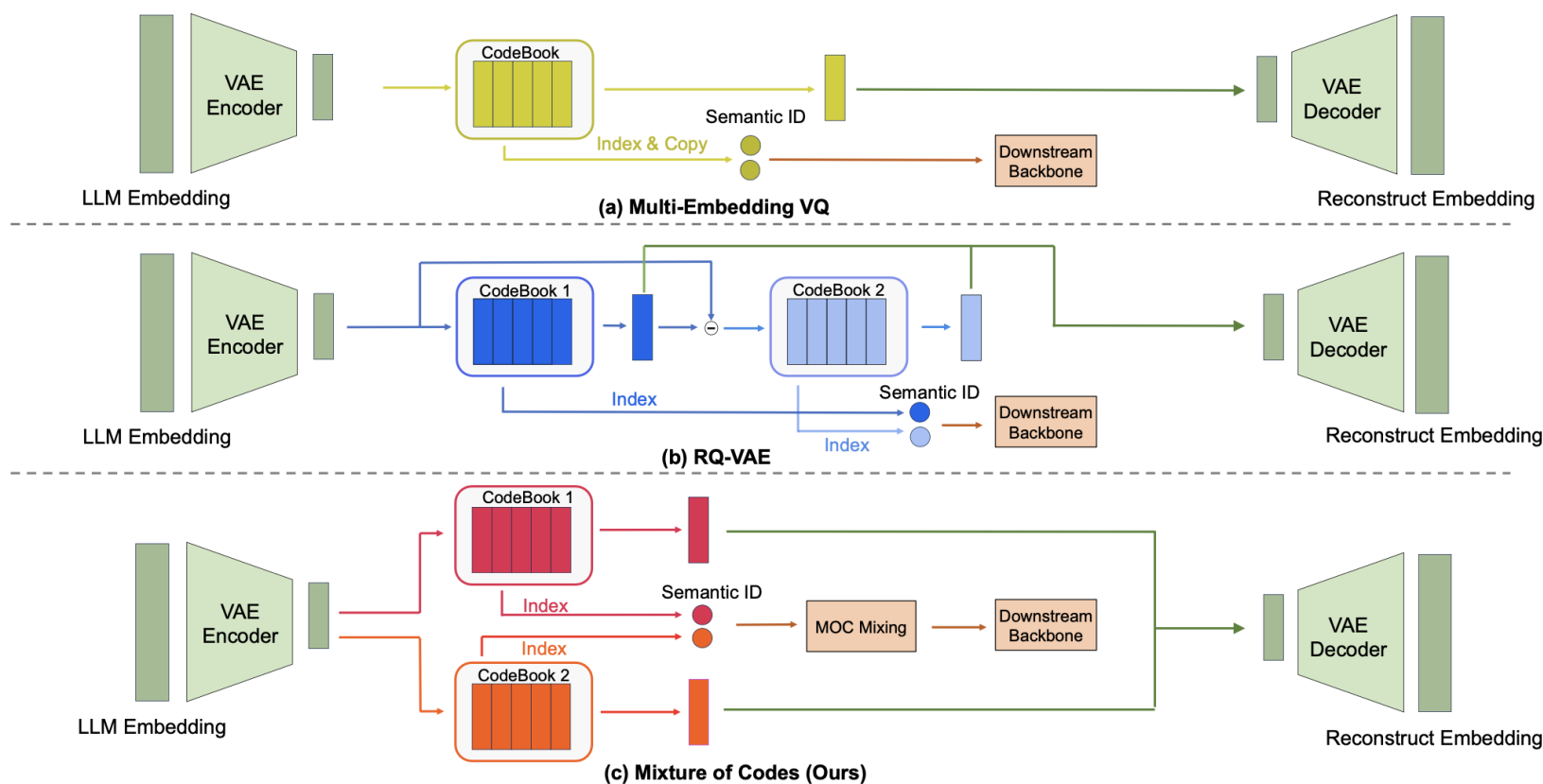
Recommender Systems with Generative Retrieval, NIPS 2023.

Reference



# Parallel Semantic IDs - MoC

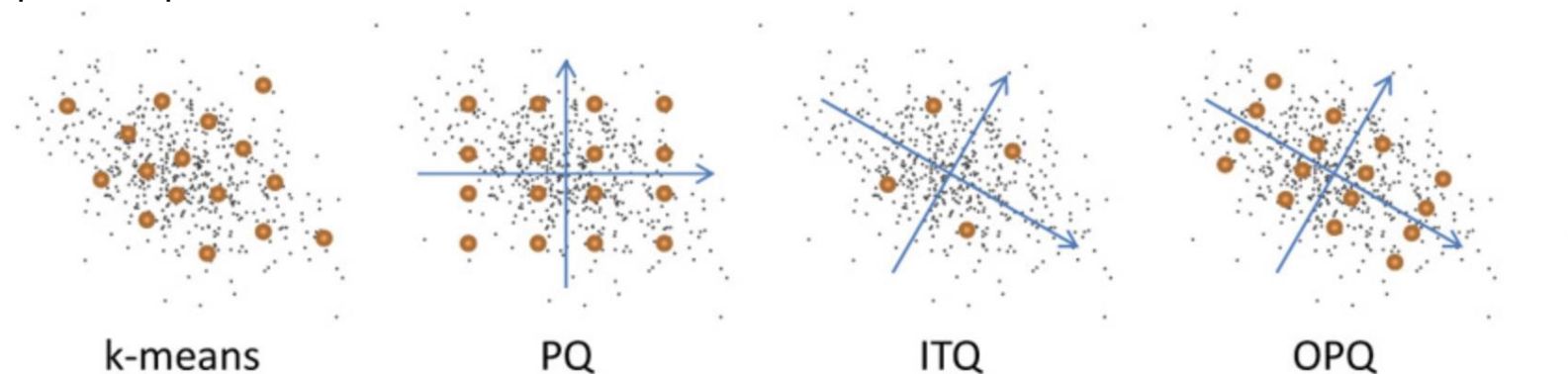
- RQ-VAE scales semantic IDs in a cascading way
- MoC scales semantic IDs in a **parallel way**



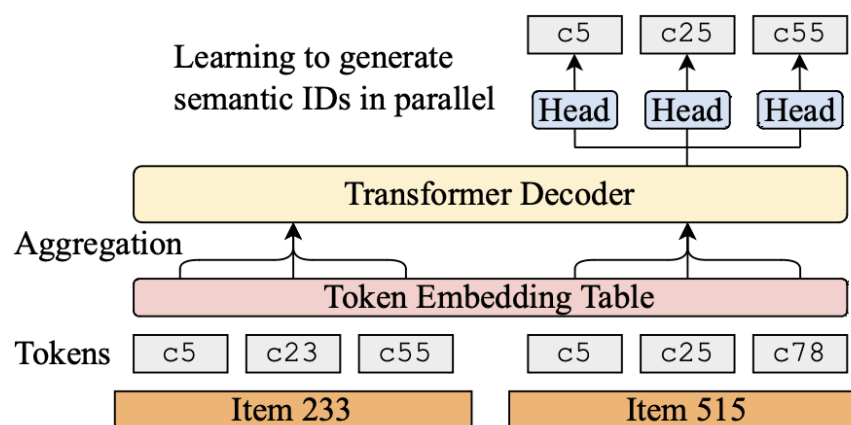
# Parallel Semantic IDs

- Generates parallel VQs with a **multi-token prediction loss**
- Employ **OPQ** to split sub-space

$$\mathcal{L} = - \sum_{j=1}^m \log \mathbb{P}^{(j)}(c_{t,j}|s) = - \sum_{j=1}^m \log \frac{\exp(\mathbf{e}_{c_{t,j}}^\top \cdot \mathbf{g}_j(s)/\tau)}{\sum_{c \in C^{(j)}} \exp(\mathbf{e}_c^\top \cdot \mathbf{g}_j(s)/\tau)},$$



*Training w/ Multi-token Prediction*



# Parallel Semantic IDs - MMQ

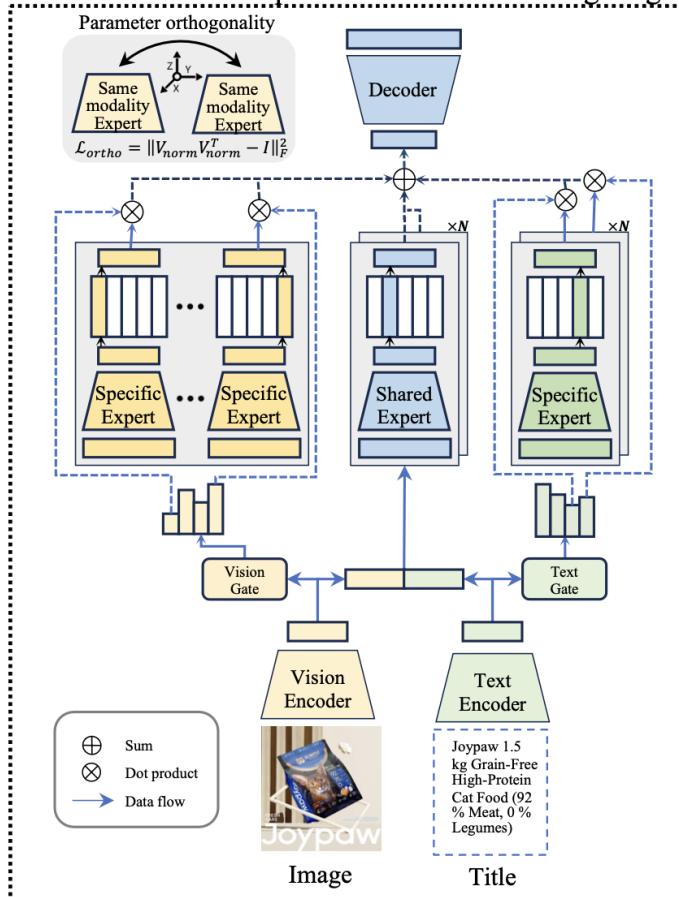
- Modality-shared and –specific experts
- Orthogonality constraints** on experts

$$\mathcal{L}_{ortho\_shared} = \|\mathbf{V}_{norm} \mathbf{V}_{norm}^T - \mathbf{I}\|_F^2,$$

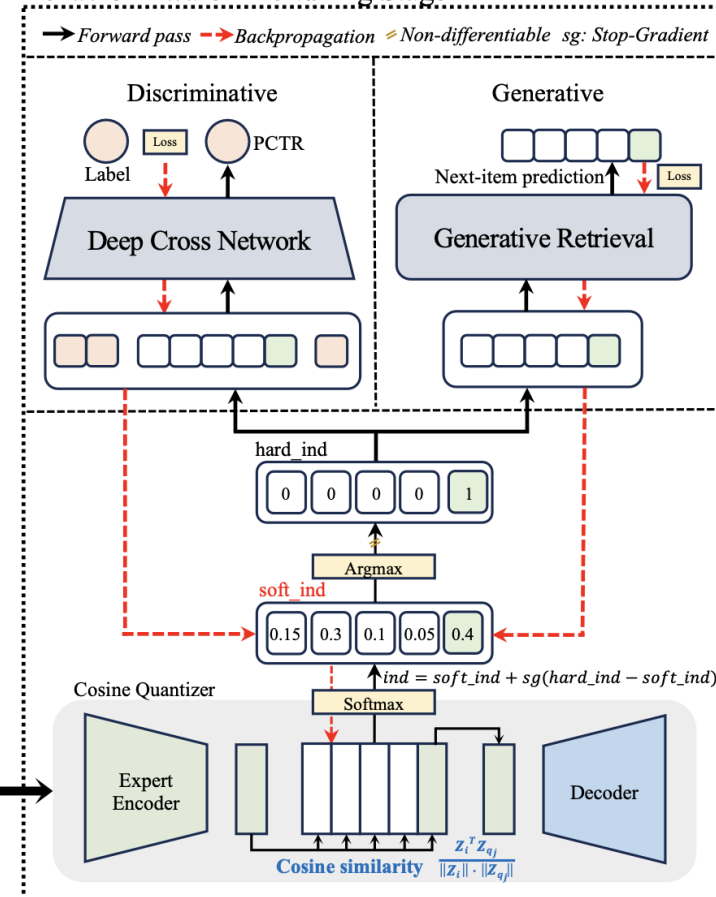
$$\mathcal{L}_{ortho\_specific} = \|\mathbf{V}'_{norm} \mathbf{V}'_{norm}^T - \mathbf{I}\|_F^2,$$

- Similar to Redundancy Reduction or Expert De-correlation Principle

Multimodal Shared-Specific Tokenizer Training Stage



Behavior-Aware Fine-tuning Stage



# Collaborative-aware - LETTER

- LETTER **aligns** the semantic IDs with the **downstream recommendation models**
- Diversity loss** to promote the uniform distribution of the code embeddings

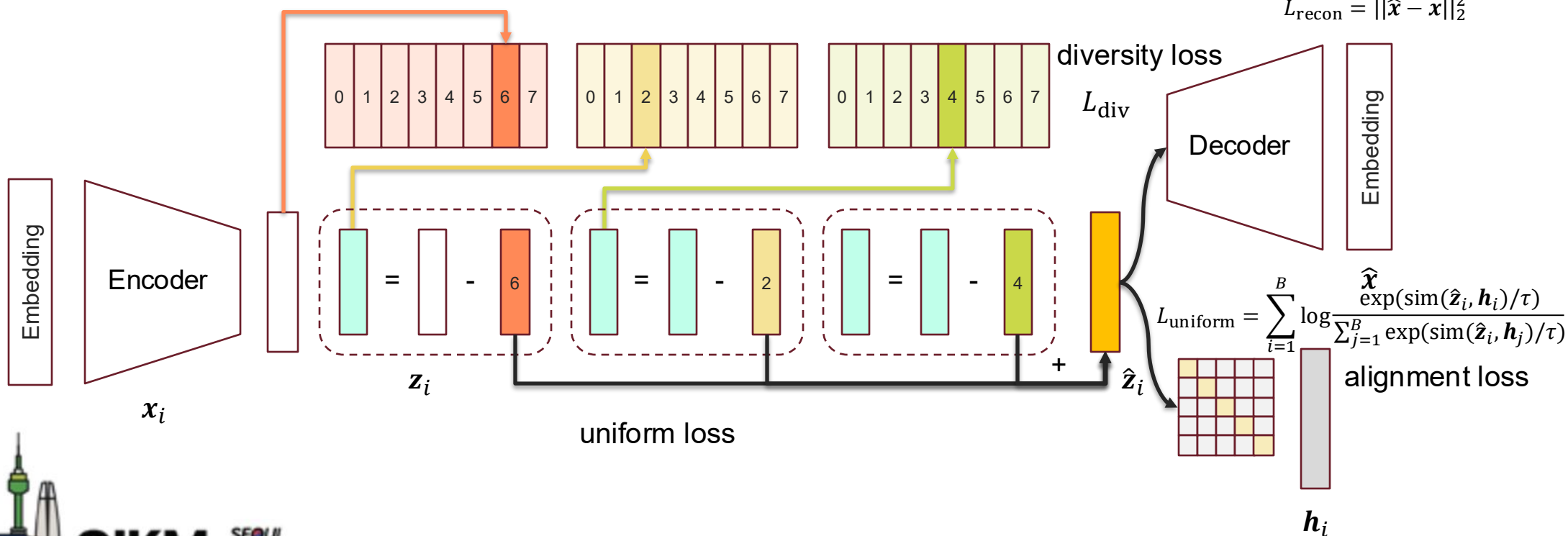
$$\mathcal{L}_{CF} = -\frac{1}{B} \sum_{i=1}^B \frac{\exp(\langle \hat{\mathbf{z}}_i, \mathbf{h}_i \rangle)}{\sum_{j=1}^B \exp(\langle \hat{\mathbf{z}}_i, \mathbf{h}_j \rangle)},$$

RQ-VAE loss

$$L_{\text{commit}} = \|\text{sg}[\mathbf{r}_{l-1}] - \mathbf{e}_{c_l}\|_2^2 + \beta \|\mathbf{r}_{l-1} - \text{sg}[\mathbf{e}_{c_l}]\|_2^2$$

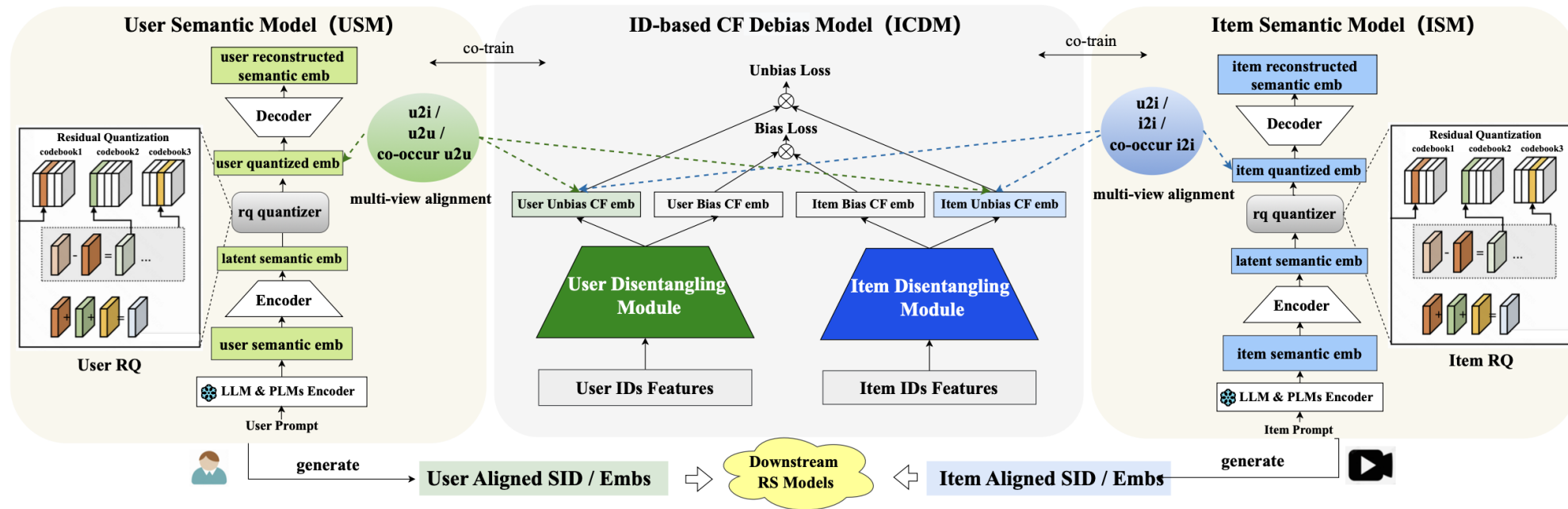
reconstruction loss

$$L_{\text{recon}} = \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2$$



# Collaborative-aware - DAS

- **Various contrastive/alignment losses** between users and items across the semantic and collaborative embeddings



# Q & A

[jonaspan@tencent.com](mailto:jonaspan@tencent.com)