

# Semantic ID

-based Generative Recommendation

### **Semantic IDs**

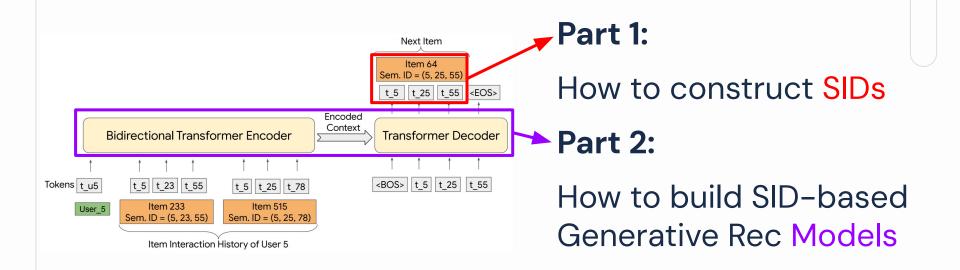
(also called: SemID or SID)

A few tokens that jointly index one item.

t3, t321, t643, t1011



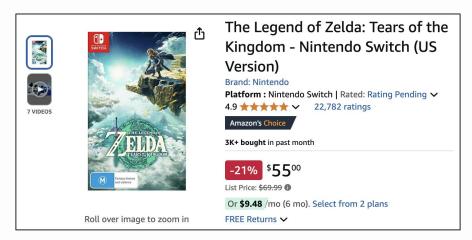
### SemID-based Generative Recommendation



### **Part 1: Semantic ID Construction**

### Semantic ID Construction

**Input:** all data associated with the item (description, title, interactions, features, ...)

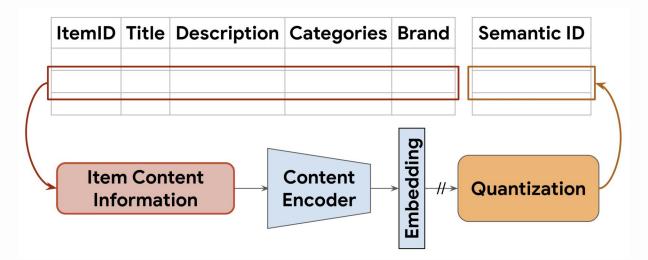


Output: mapping between items ⇔ Semantic IDs

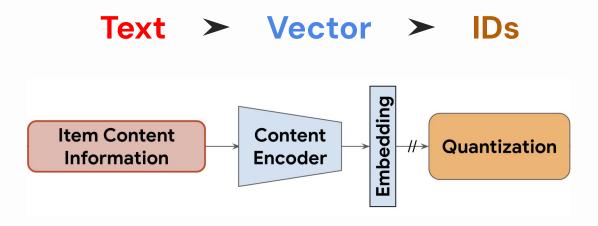
### **Part 1: Semantic ID Construction**

(i) First example: TIGER and RQ-VAE-based SemIDs;

Input: concatenated text features



Output: mapping between items ⇔ Semantic IDs



1. Item Content Information (Text)

ItemID Title Description

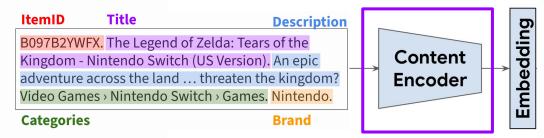
B097B2YWFX. The Legend of Zelda: Tears of the Kingdom - Nintendo Switch (US Version). An epic adventure across the land ... threaten the kingdom? Video Games > Nintendo Switch > Games. Nintendo.

**Categories** 

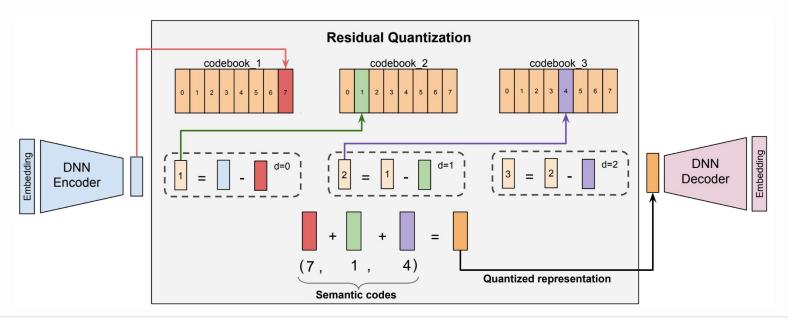
**Brand** 

2. Content Encoder + Embedding (Text ➤ Vector)

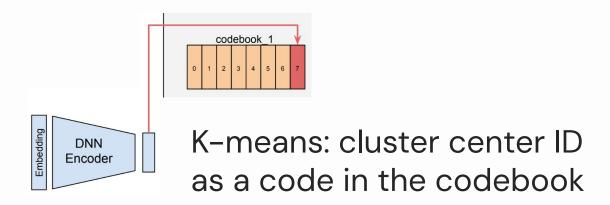
Pre-trained (fixed) sentence embedding model (SentenceT5)



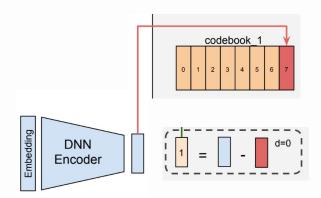
### 3. RQ-VAE Quantization (Vector ➤ IDs)



3. RQ-VAE Quantization (Vector ➤ IDs)

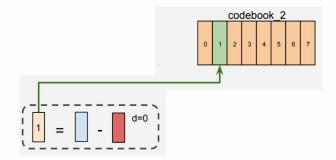


3. RQ-VAE Quantization (Vector ➤ IDs)



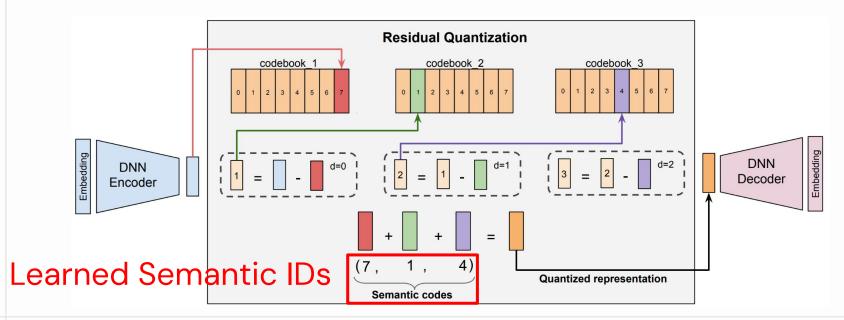
Residual of "input vector" and "clustering center vector"

3. RQ-VAE Quantization (Vector ➤ IDs)



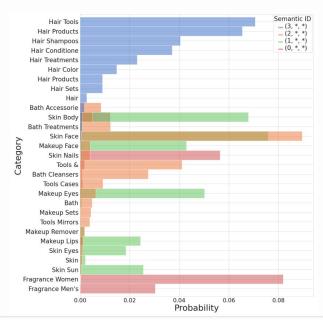
Residual as next level's input

### 3. RQ-VAE Quantization (Vector ➤ IDs)



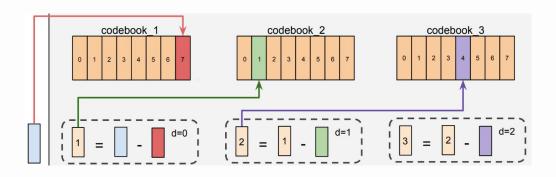
### **Properties** of **RQ-VAE**-based SemIDs

1. Semantic;

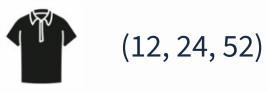


### Properties of RQ-VAE-based SemIDs

- 1. Semantic;
- 2. Ordered / sequential dependent;



#### **Collisions**





(12, 24, 52)

#### **Collisions**



(12, 24, 52, <mark>0</mark>)



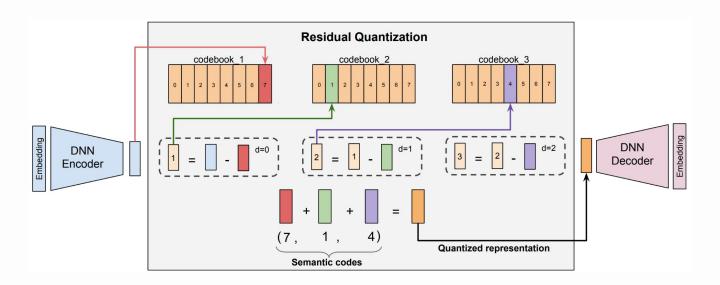
(12, 24, 52, **1**)

One extra token to avoid conflicts

### **Part 1: Semantic ID Construction**

- (i) First example: TIGER and RQ-VAE-based SemIDs;
- (ii) Techniques to construct SemIDs;

### Residual Quantization: RQ-VAE



Residual Quantization: RQ-VAE

#### Issues:

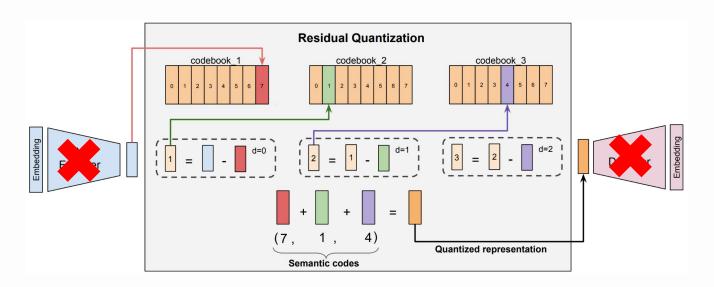
- Enc-Dec Training Unstable
- Unbalanced IDs

Residual Quantization: RQ-VAE

#### Issues:

- Enc-Dec Training Unstable
- Unbalanced IDs

Variants of Residual Quantization: RQ

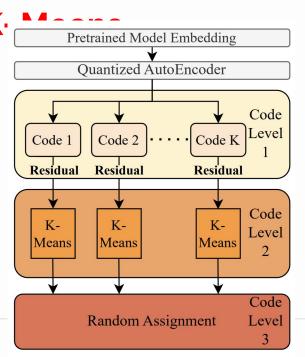


Residual Quantization: RQ-VAE

#### Issues:

- Enc-Dec Training Unstable
- Unbalanced IDs

Variants of Residual Quantization: VQ-VAE +



MBGen [CIKM'24]:

Layer 1: VQ-VAE

Layer 2: K-Means on residuals

Layer 3: Random hash

### Variants of Residual Quantization: RQ + K-Means

```
Algorithm 1: Balanced K-means Clustering
   Input: Item set \mathcal{V}, number of clusters K
1 Compute w \leftarrow |\mathcal{V}|/K
<sup>2</sup> Initialize centroids C_l = \{c_1^l, \dots, c_K^l\} with random selection;
 3 repeat
         Initialize unassigned set \mathcal{U} \leftarrow \mathcal{V}
         for each cluster k \in \{1, ..., K\} do
              Sort \mathcal{U} by ascending distance from centroid c_{\iota}^{l};
 6
              Assign \mathcal{V}_k \leftarrow \mathcal{U}[0:w-1];
 7
              Update centroid c_k^l \leftarrow \frac{1}{w} \sum_{r^l \in V_k} r^l;
 8
              Remove assigned items \mathcal{U} \leftarrow \mathcal{U} \setminus \mathcal{V}_k;
         end
11 until Assignment convergence;
    Output: Optimized codebook C_l = \{c_1^l, \dots, c_K^l\}
```

OneRec:

Limit the max #items per cluster

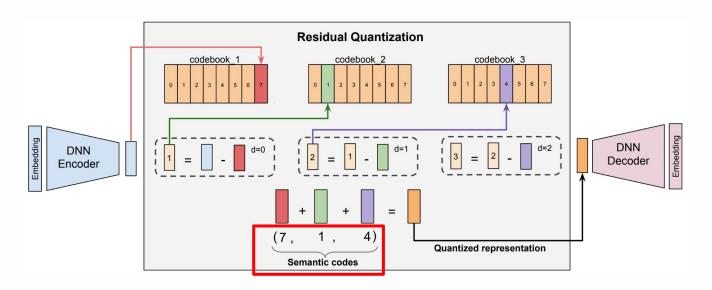
### Variants of Residual Quantization

Table 1: Performance of GR models with by SIDs generated by different tokenization algorithms.

Methods	Recommendation Performance (Recall@5/Recall@10/NDCG@5/NDCG@10)											
	Beauty				Toys				Sports			
RK-Means	0.0422	0.0639	0.0277	0.0347	0.0376	0.0577	0.0243	0.0308	0.0236	0.0353	0.0153	0.0191
R-VQ	0.0422	0.0638	0.0282	0.0351	0.0327	0.0493	0.0209	0.0262	0.0234	0.0352	0.0151	0.0189
RQ-VAE	0.0404	0.0593	0.0268	0.0329	0.0342	0.0514	0.0224	0.0280	0.0205	0.0312	0.0132	0.0166

### https://github.com/snap-research/GRID

Are the orders among tokens necessary?



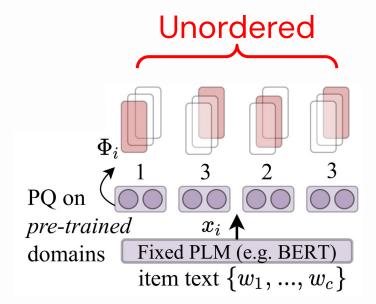
Are the orders among tokens necessary?

#### **Orders**

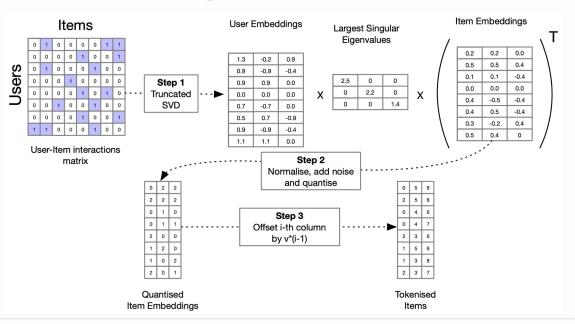
<-> autoregressive generation

<-> SemIDs have to be short

#### **Product Quantization**



### **Product Quantization**



Early exploration on PQ-based SIDs, but performance is not ideal

Autoregressive generation needs orders

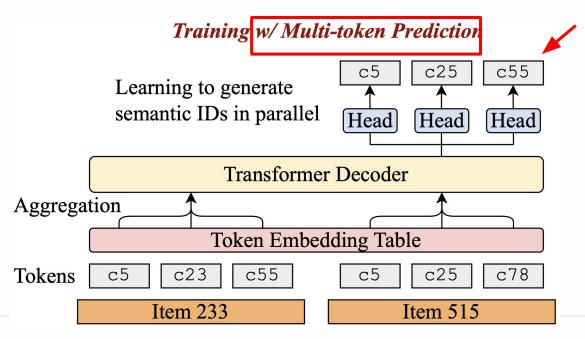
**Product Quantization: RPG [KDD'25]** 

Long semantic IDs + parallel generation

Original (#digit=4): †34, †392, †600, †891

**RPG** (#digit=64): †34, †392, †600, †891, †1221, †1495, †1556, †1924, †2296, †2511, †2715, †3060, †3255, †3439, †3655, †3984, †4167, †4400, †4736, †4939, †5278, †5426, †5669, †6057, †6385, †6451, †6837, †7134, †7329, †7528, †7924, †8162, †8325, †8479, †8711, †9007, †9420, †9472, †9980, †10154, †10364, †10662, †10784, †11105, †11377, †11642, †11848, †12261, †12334, †12585, †12963, †13306, †13367, †13722, †13973, †14143, †14506, †14696, †14995, †15331, †15406, †15813, †16034, †16251

**Product Quantization: RPG [KDD'25]** 



Using *OPQ* to produce semIDs w/o sequential dependencies

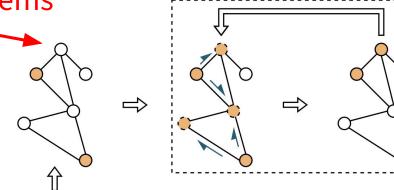
### **Product Quantization: RPG [KDD'25]**

Inference w/ Graph-constrained Decoding beam size b=2

(4) Iterate for q steps

Graph: connecting similar SIDs/items

(offline)



(1) Sample the initial beam

(2) **Propagate** on decoding graph

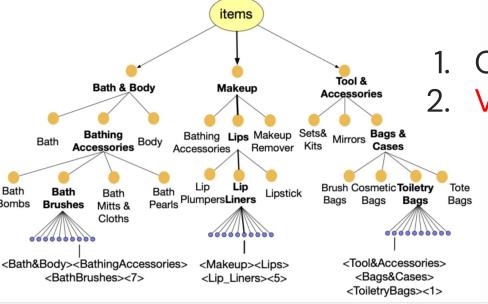
(3) Keep the best *b* nodes

Final recommendations

- C5
  - c5 c
    - :25 c5

- 0
- с3
- 23 c6
- O Candidate SID
- SID in beam
- Similar SIDs
- Propagation

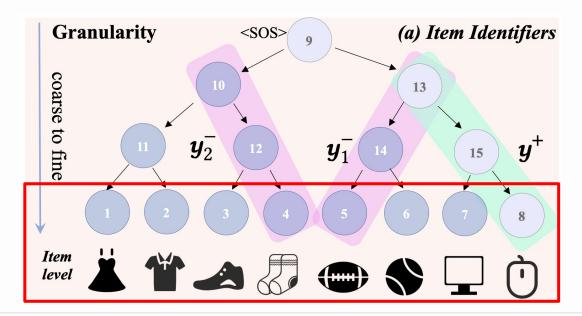
Hierarchical Clustering (Heuristics-based)



- 1. Ordered;
- 2. Variable-length SemIDs;

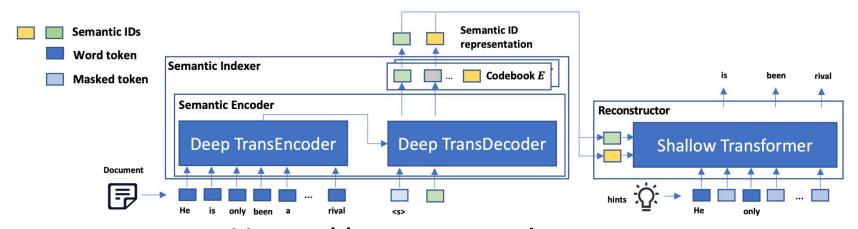
### **Techniques to Construct SemIDs**

Hierarchical Clustering (Latent-based)



### **Techniques to Construct SemIDs**

#### Language Model-based ID Generator



Natural language as inputs; SemIDs as outputs

# **Techniques to Construct SemIDs**

#### Language Model-based ID Generator

Generated ID for Item A

Generated ID for Item B

jessica simpson perfume women

broadway performer buffalo dance lessons

LLM-based ID Generator

Words as SemIDs (like tagging)

title: jessica simpson fancy love eau de parfum spray, 1.7 ounce; brand: jessica simpson; description: buy jessica simpson womens perfumes fancy love by jessica simpson for women 1.7 oz eau de parfum spray; categories: beauty, fragrance, womens, eau de parfum; price: 23.54; salesrank: beauty: 132446

Item A plain text from Beauty

title: stepping out vhs; brand: na; description: a hasbeen broadway performer moves to buffalo and starts teaching tap dance lessons to a group of misfits who, through their dance classes, bond and realize what they can achieve. their newfound selfconfidence changes their lives forever.; categories: movies tv, movies; price: 19.99; salesrank: movies tv: 244008

Item B plain text from Movie

# **Summary of Techniques to Construct SemIDs**

- Residual Quantization (ordered)
- Product Quantization (unordered)
- Hierarchical Clustering
- LM-based ID Generator

#### **Part 1: Semantic ID Construction**

- (i) First example: TIGER and RQ-VAE-based SemIDs; (ii) Techniques to construct SemIDs;
- (iii) Inputs for SemID construction;

#### Input: all data associated with the item



Input: all data associated with the item

What exactly does "all data" mean? 🤷



#### **Text or Multimodal Features**

Text/Visual/Acoustic ➤ Vector ➤ IDs

**Brand** 

Pretrained Encoder Quantization

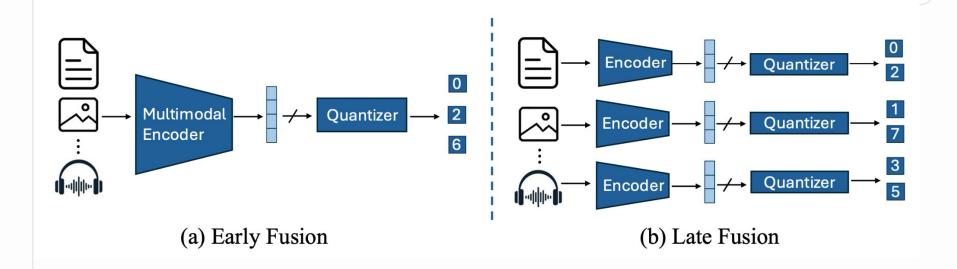
B097B2YWFX. The Legend of Zelda: Tears of the Kingdom - Nintendo Switch (US Version). An epic adventure across the land ... threaten the kingdom? Video Games > Nintendo Switch > Games. Nintendo.

Text

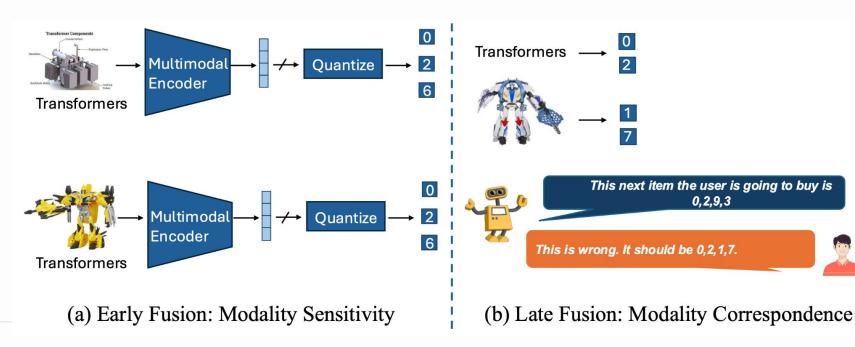
**Categories** 



#### **Text or Multimodal Features**



#### **Text or Multimodal Features**

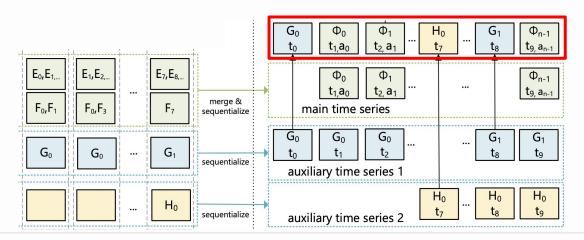


0, 2, 9, 3

#### **Categorical Features**

#### **Categorical Features** > IDs

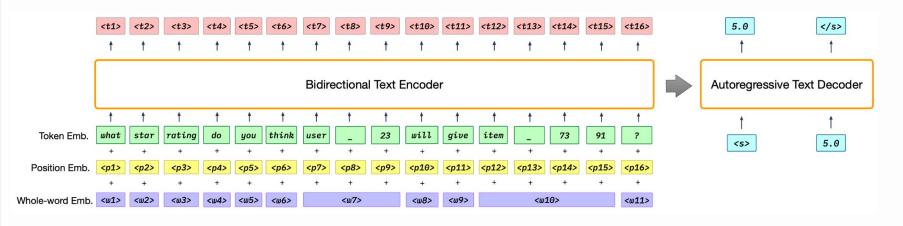
Merge & Sequentialize



#### No Features

#### Item ID ➤ IDs

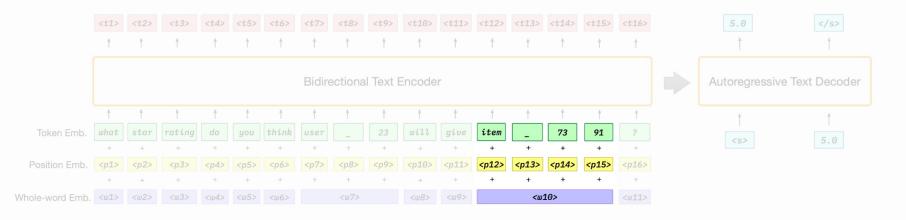
Text Tokenizer



#### **No Features**

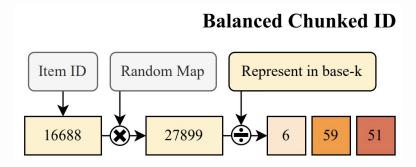
Item ID ➤ IDs

Text Tokenizer



**No Features** 

#### Random IDs



Input: all data associated with the item

#### (1) Item Metadata

Text / Multimodal / Categorical / No Features

Input: all data associated with the item

(1) Item Metadata

Text / Multimodal / Categorical / No Features

(2) Item Metadata + Behaviors

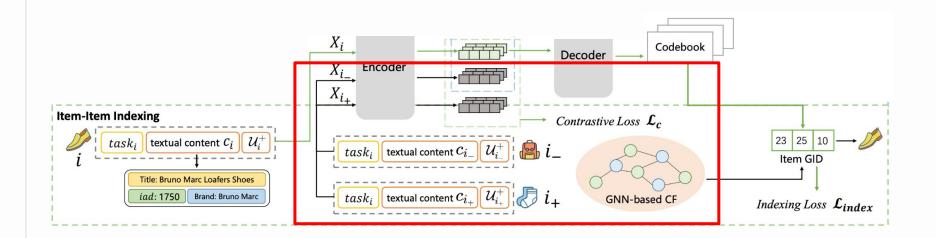
Input: all data associated with the item

(1) Item Metadata

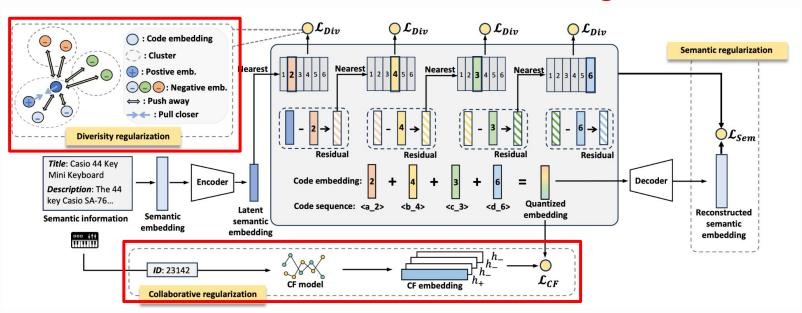
Text / Multimodal / Categorical / No Features

(2) Item Metadata + Behaviors
But how?

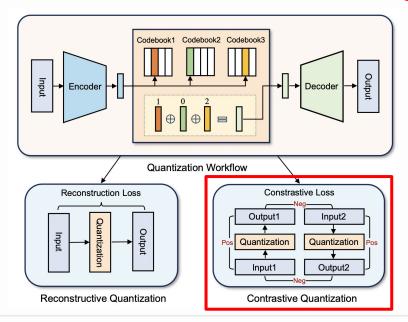
Residual Quantization + Item-level Regularization



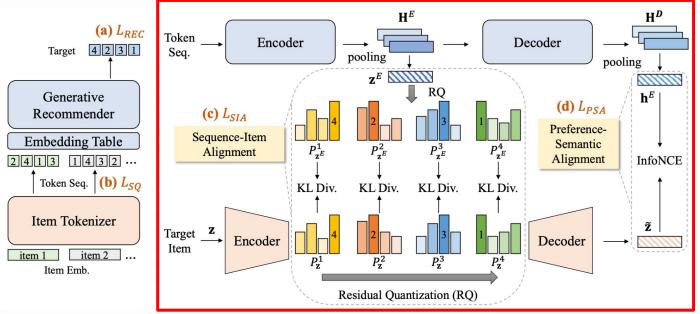
#### Residual Quantization + Item-level Regularization



Residual Quantization + Item-level Regularization

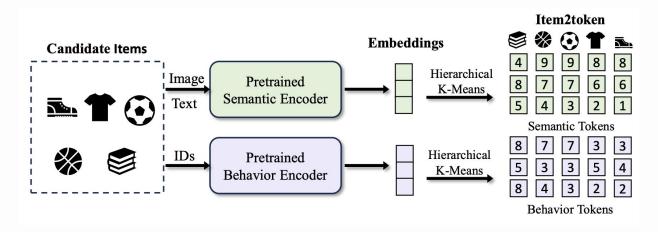


Residual Quantization + Recommendation Loss



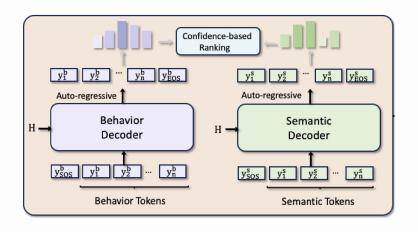
Item Metadata + Behaviors

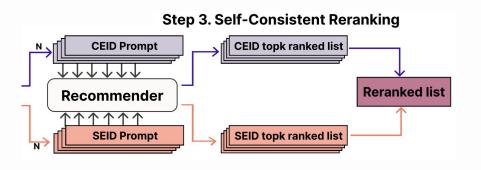
#### **Fused Semantic IDs**



**Item Metadata + Behaviors** 

Fused Semantic IDs + Two-stream Generation

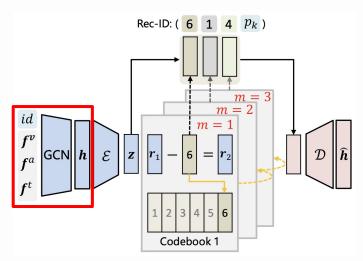




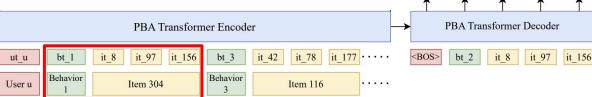
**Item Metadata + Behaviors** 

**Fused Representations** 

User-Item Graph + Semantic Features



Unified Multi-Task Framework Item Metadata + Behaviors **Target Behavior Item Prediction Multi-Behavior Recommendation Behavior-Specific Item Prediction Behavior-Item Prediction Behavior Prediction** Semantic IDs fused with behavior types Next Behavior



it?

it?

it?

Next Item

it 97

it?

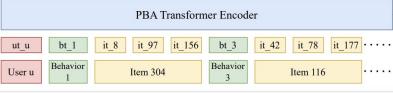
it 156

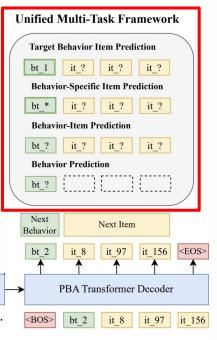
**Item Metadata + Behaviors** 

**Multi-Behavior Recommendation** 

Next Token Prediction as natural multi-task learning

(prompted by behavior type)



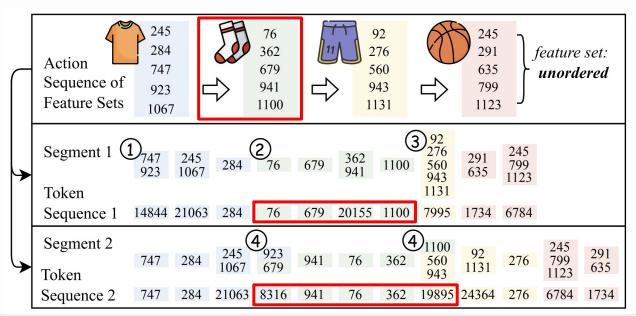


#### Context-independent

Action Tokenization	Example	Contextual	Unordered
Product Quantization	VQ-Rec (Hou et al., 2023)	X	✓
Hierarchical Clustering		X	×
Residual Quantization	TIGER (Rajput et al., 2023)	X	X
<b>Text Tokenization</b>	LMIndexer (Jin et al., 2024)	X	×
Raw Features	HSTU (Zhai et al., 2024)	X	×
SentencePiece	SPM-SID (Singh et al., 2024)	X	×

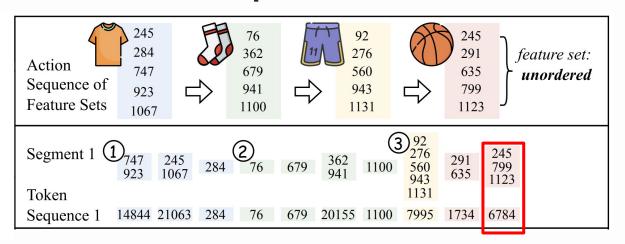
Same item ⇒ fixed semIDs in all sequences

#### Context-independent ⇒ Context-aware



Same item ⇒
different semIDs
based on context

#### **Context-independent** ⇒ **Context-aware**



#### Core Idea:

Merge frequently co-occurring features as new tokens

(ActionPiece: "WordPiece" tokenization for generative rec)

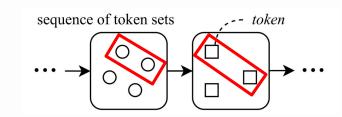
#### **Context-independent** ⇒ **Context-aware**

#### **Algorithm 1** ActionPiece Vocabulary Construction

**input** Sequence corpus S', initial tokens  $V_0$ , target size Q**output** Merge rules  $\mathcal{R}$ , constructed vocabulary  $\mathcal{V}$ 

- 1: Initialize vocabulary  $V \leftarrow V_0$  # each initial token corresponds to one unique item feature
- 2:  $\mathcal{R} \leftarrow \emptyset$
- 3: while  $|\mathcal{V}| < Q$  do
- # Count: accumulate weighted token co-occurrences
- $count(\cdot, \cdot) \leftarrow Count(\mathcal{S}', \mathcal{V}) \# Algorithm 2$
- # Update: Merge a frequent token pair into a new token
- Select  $(c_u, c_v) \leftarrow \arg\max_{(c_i, c_j)} \operatorname{count}(c_i, c_j)$
- $S' \leftarrow \text{Update}(S', \{(c_u, c_v) \rightarrow c_{\text{new}}\}) \text{ # Algorithm 3}$
- $\mathcal{R} \leftarrow \mathcal{R} \cup \{(c_u, c_v) \rightarrow c_{\text{new}}\}$  # new merge rule
- $\mathcal{V} \leftarrow \mathcal{V} \cup \{c_{\text{new}}\}\$ # add new token to the vocabulary
- 11: end while

return  $\mathcal{R}, \mathcal{V}$ 



$$P(\bigcirc,\bigcirc) = \frac{|\bigcirc -\bigcirc|}{|<\bigcirc,\bigcirc>|} = \frac{4-1}{\binom{4}{2}}$$

$$P(\square,\square) = \frac{|\square - \square|}{|<\square,\square>|} = \frac{3-1}{\binom{3}{2}}$$

 $P(O,O) = \frac{|O-O|}{|<O,O>|} = \frac{4-1}{\binom{4}{2}}$  Features co-occurring within

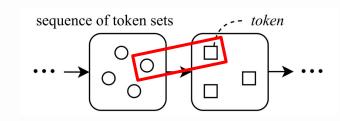
#### **Context-independent** ⇒ **Context-aware**

#### Algorithm 1 ActionPiece Vocabulary Construction

**input** Sequence corpus S', initial tokens  $V_0$ , target size Q **output** Merge rules  $\mathcal{R}$ , constructed vocabulary V

- 1: Initialize vocabulary  $V \leftarrow V_0$  # each initial token corresponds to one unique item feature
- 2:  $\mathcal{R} \leftarrow \emptyset$
- 3: while  $|\mathcal{V}| < Q$  do
- 4: # Count: accumulate weighted token co-occurrences
- 5:  $\operatorname{count}(\cdot, \cdot) \leftarrow \operatorname{Count}(\mathcal{S}', \mathcal{V}) \# \operatorname{Algorithm} 2$
- 6: # *Update*: Merge a frequent token pair into a new token
- 7: Select  $(c_u, c_v) \leftarrow \arg\max_{(c_i, c_j)} \operatorname{count}(c_i, c_j)$
- 8:  $S' \leftarrow \text{Update}(S', \{(c_u, c_v) \rightarrow c_{\text{new}}\}) \# \text{Algorithm 3}$
- 9:  $\mathcal{R} \leftarrow \mathcal{R} \cup \{(c_u, c_v) \rightarrow c_{\text{new}}\}\$ # new merge rule
- 10:  $\mathcal{V} \leftarrow \mathcal{V} \cup \{c_{\text{new}}\}\$ # add new token to the vocabulary
- 11: end while

return  $\mathcal{R}, \mathcal{V}$ 



$$P(\bigcirc,\bigcirc) = \frac{|\bigcirc -\bigcirc|}{|<\bigcirc,\bigcirc>|} = \frac{4-1}{\binom{4}{2}}$$

$$P(\bigcirc,\bigcirc) = \frac{|\bigcirc -\bigcirc|}{|\bigcirc| \times |\bigcirc|} = \frac{1}{4 \times 3}$$

$$P(\square,\square) = rac{|\square - \square|}{|<\square,\square>|} = rac{3-1}{{3 \choose 2}}$$

Features
co-occurring
within or
across items

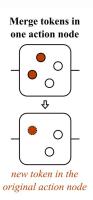
#### **Context-independent** ⇒ **Context-aware**

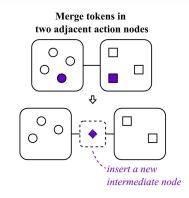
#### Algorithm 1 ActionPiece Vocabulary Construction

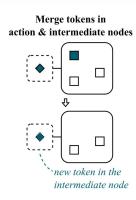
**input** Sequence corpus S', initial tokens  $V_0$ , target size Q **output** Merge rules R, constructed vocabulary V

- 1: Initialize vocabulary  $V \leftarrow V_0$  # each initial token corresponds to one unique item feature
- 2:  $\mathcal{R} \leftarrow \emptyset$
- 3: while  $|\mathcal{V}| < Q$  do
- 4: # Count: accumulate weighted token co-occurrences
- 5:  $\operatorname{count}(\cdot, \cdot) \leftarrow \operatorname{Count}(\mathcal{S}', \mathcal{V}) \# \operatorname{Algorithm} 2$
- 6: # *Update*: Merge a frequent token pair into a new token
  - Select  $(c_u, c_v) \leftarrow \arg\max_{(c_i, c_j)} \operatorname{count}(c_i, c_j)$
- 8:  $S' \leftarrow \text{Update}(S', \{(c_u, c_v) \rightarrow c_{\text{new}}\}) \# \text{Algorithm 3}$
- 9:  $\mathcal{R} \leftarrow \mathcal{R} \cup \{(c_u, c_v) \rightarrow c_{\text{new}}\} \# \text{ new merge rule}$
- 10:  $\mathcal{V} \leftarrow \mathcal{V} \cup \{c_{\text{new}}\}\$ # add new token to the vocabulary
- 11: end while

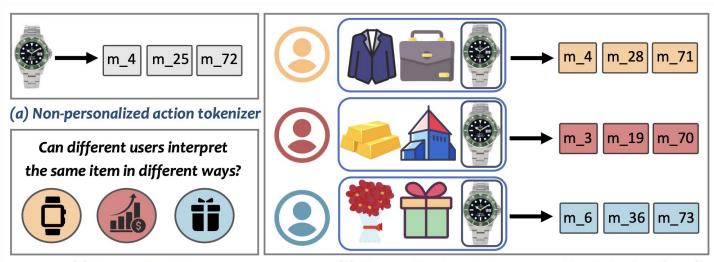
return  $\mathcal{R}, \mathcal{V}$ 







# Context-independent ⇒ Context-aware ⇒ Personalized



(c) Our motivation

(b) Personalized context-aware action tokenizer (ours)

Input: all data associated with the item

(1) Item Metadata

Text / Multimodal / Categorical / No Features

(2) Item Metadata + Behaviors

Regularization / Fusion

Context-independent ⇒ Context-aware

### Part 1 Summary - SemID Construction

- (1) First Example: TIGER
- (2) Construction Techniques

(3) Inputs

### Part 1 Summary - SemID Construction

- (1) First Example: TIGER
- (2) Construction Techniques

RQ, PQ, Clustering, LM-based generator

(3) Inputs

# Part 1 Summary - SemID Construction

- (1) First Example: TIGER
- (2) Construction Techniques

RQ, PQ, Clustering, LM-based generator

(3) Inputs

Item Metadata (Text, Multimodal)

+ Behaviors (Regularization, Fusion, Context)

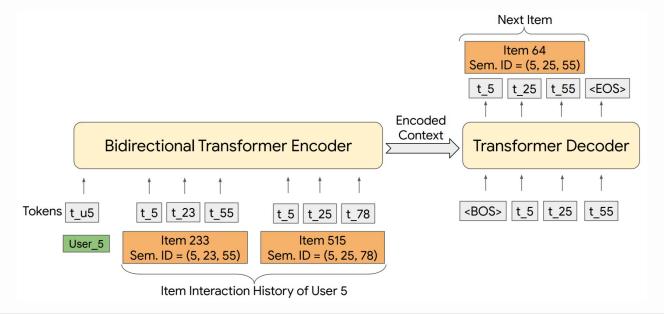
# Part 2: SemID-based Generative Recommendation Model Architecture

Recommendation as a seq-to-seq generation problem

```
Input: user interacted items \{c_{11'}, c_{12'}, c_{13'}, c_{14'}, c_{21'}, c_{22'}, ...\}
```

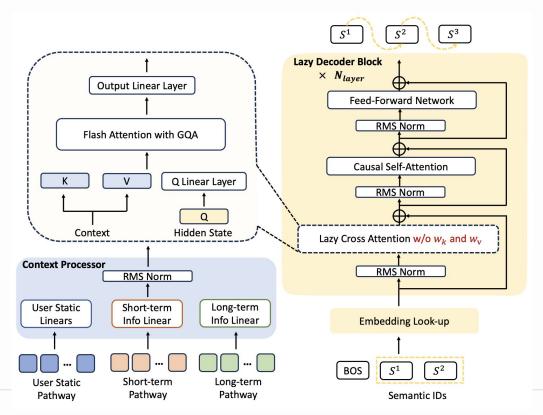
Output: next item  $\{c_{t1'}, c_{t2'}, c_{t3'}, c_{t4}\}$ 

Architecture: Encoder-Decoder

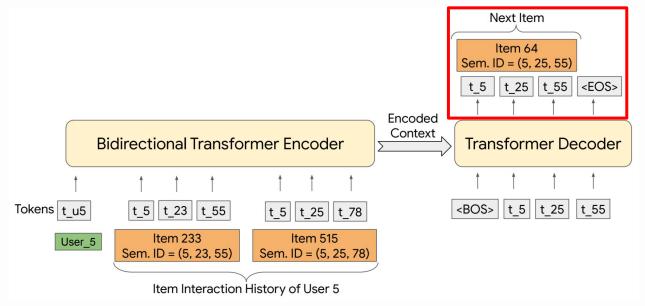


#### **Architecture:**

Decoder-Only?

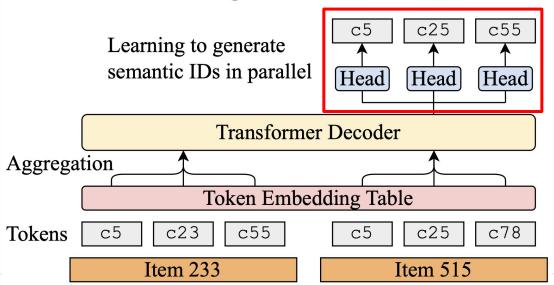


**Objective:** Next-Token Prediction (w/ RQ)



**Objective:** Multi-Token Prediction (w/ PQ)

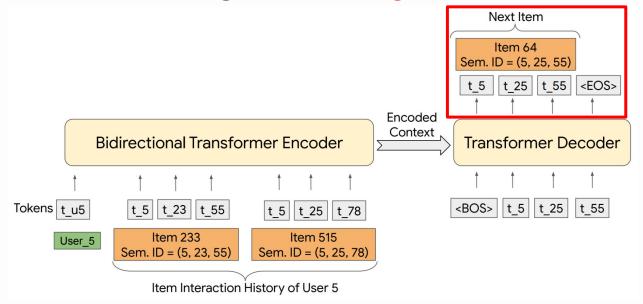
#### Training w/ Multi-token Prediction



Objective: RL

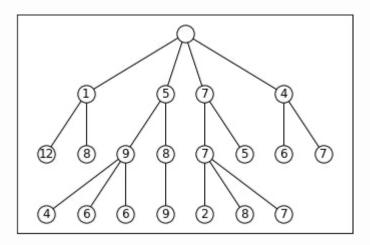
Will introduce later at "Align w/ LLMs"

**Inference:** How to get a ranking list?

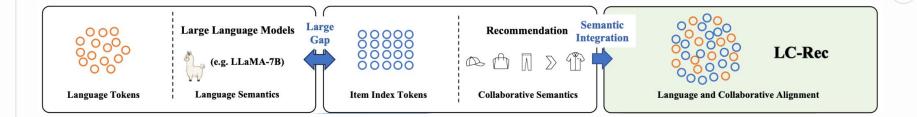


**Inference:** How to get a ranking list?

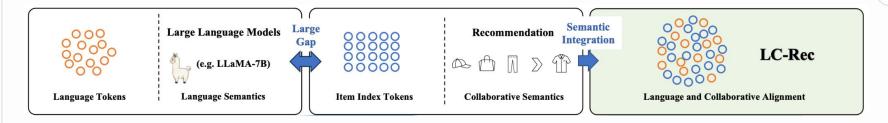
(Constrained) Beam Search



**Inference:** Can we do dense retrieval-like grounding?

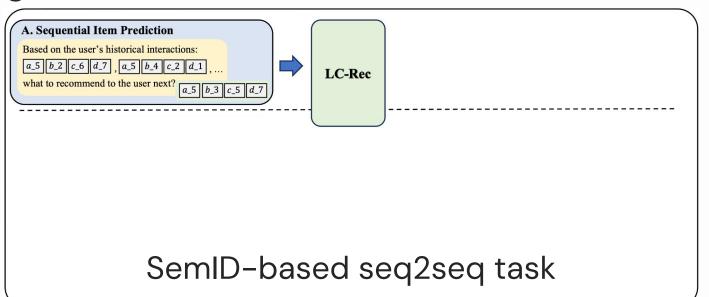


Align with LLMS - LC-Rec

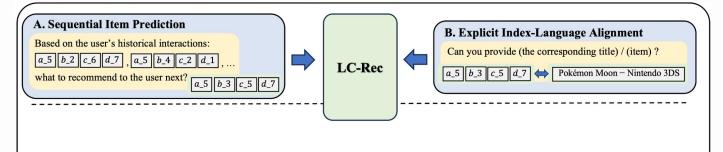


Core Idea:

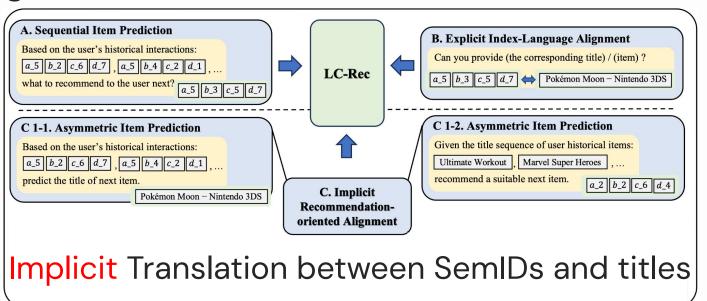
Construct instructions containing both Semantic IDs and language tokens

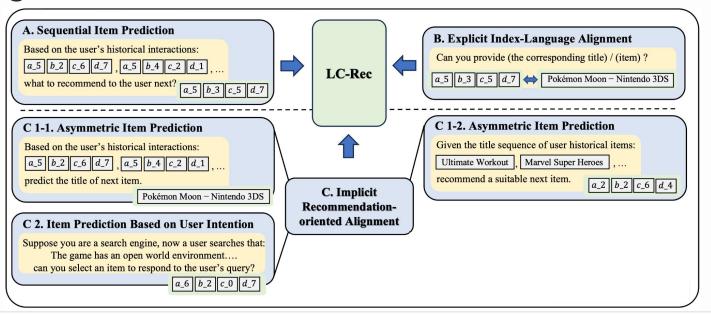


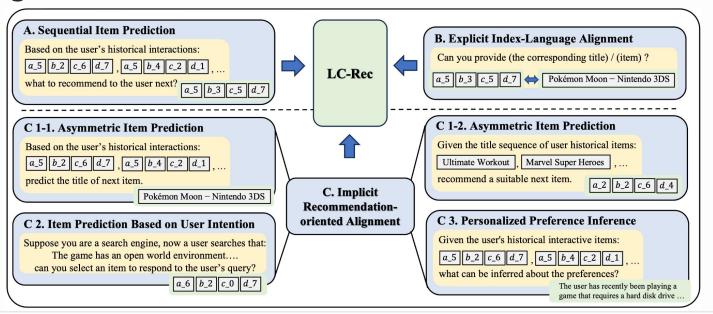
Align with LLMS - LC-Rec



Translation between SemIDs and titles

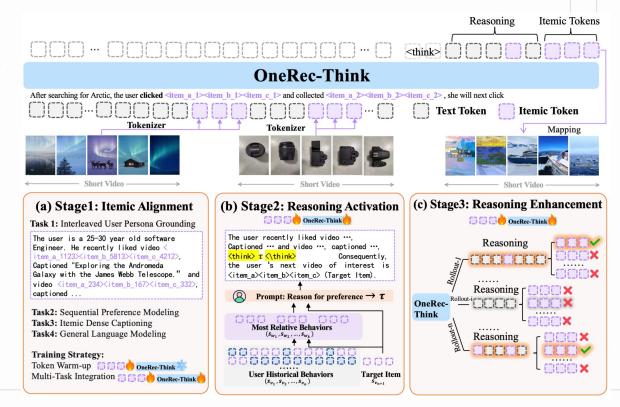




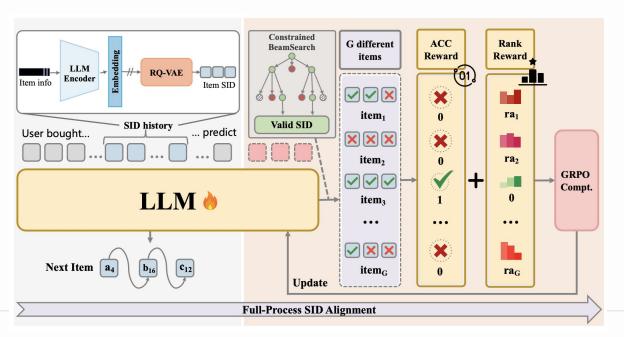


Align with LLMS

OneRec-Think



#### Align with LLMS: MiniOneRec



# Part 2 Summary - Architecture

(1) Train from Scratch

(2) Align with LLMs

# Part 2 Summary - Architecture

(1) Train from Scratch

Objective (NTP, MTP, RL)

Inference (Beam Search)

(2) Align with LLMs

# Part 2 Summary - Architecture

(1) Train from Scratch

Objective (NTP, MTP, RL)

Inference (Beam Search)

(2) Align with LLMs

LC-Rec / OneRec-Think / MiniOneRec