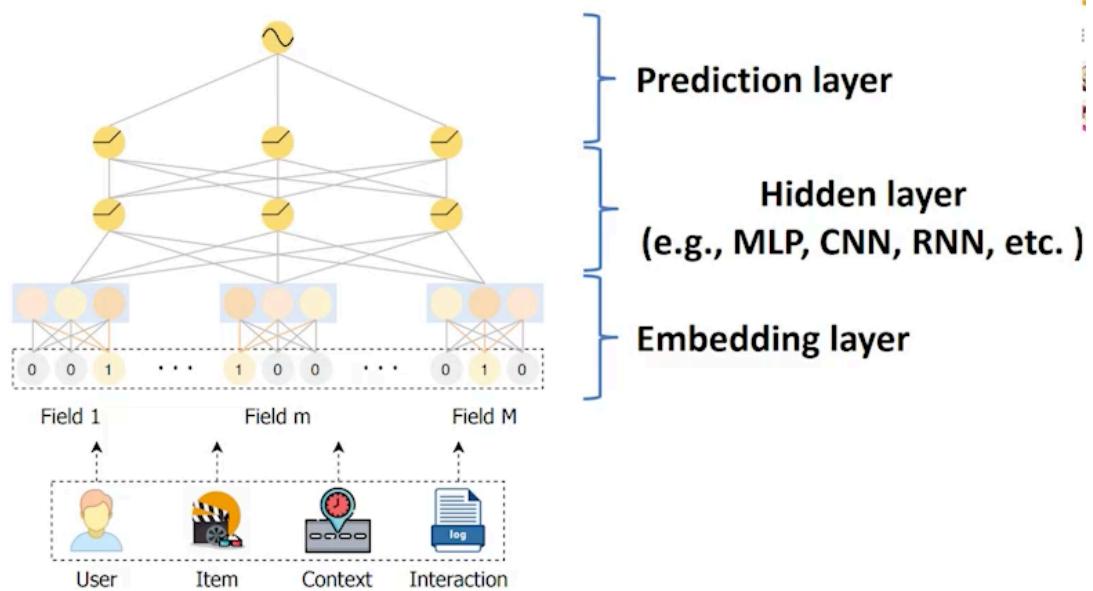


A General Architecture of Deep Recommender System



输入特征-文本嵌入-隐藏层（RNN，深度学习网络），得到更精准的预测表征-预测层（通过查询与物品的相似性度量，点击，对相关物品进行排序，得到预测）

现在：检索/召回一系列可行的候选对象，然后用排序模型对其进行排序

检索：通过矩阵分解，在同一空间学习查询和候选的嵌入，为了更好地捕获数据中的非线性关系，近年来，采用内积查询和候选嵌入到同一空间的双塔编码架构（即一个塔用于查询另一个塔用于候选）成为主流。

为了在推理期间使用这些模型，使用候选塔创建一个存储所有物品嵌入的索引。对于给定查询，通过在同一个空间内嵌入查询和候选项来执行大规模检索，然后近似最近邻搜索（ANN）来选择给定查询嵌入的最佳候选项。

推荐系统存在的局限与挑战

- 召回 -> 排序 -> 重排。流程复杂，需要分别优化。
 - 对召回阶段的依赖（庞大候选集快速筛选相关物品）要是召回都没召回用户感兴趣的东西，排序更不可能找到。
 - 对排序阶段的局限（基于学习的排序模型//神经网络模型//对候选物品排序）
 - 对反馈循环的影响（基于用户的历史交互行为进行预测，产生反馈循环）
- 冷启动推荐 新商品被反馈循环影响
- 推荐的“长尾”现象 热门更热，冷门更冷

- 可解释性，个性化
- 推荐多样性
- 依赖于用户和物品的协同过滤信号和特征

本文工作：提出全新的推荐系统框架（TIGER）基于生成式检索范式并应用到序列推荐中： 1. 我们创建语义上有意义的token元组，作为每个物品的语义ID。通过生成的用户交互序列中物品的语义ID，训练基于transformer的序列到序列模型来“生成”用户将与之交互的下一个物品的语义ID。 2. 提出TIGER推荐模型在各种数据集上的性能显著优于当前的SOTA模型。 3. 展现出两个关键能力：其一是冷启动推荐能力，能够推荐此前从未出现过的物品；其二是推荐多样性控制能力，可通过调节生成参数实现推荐内容的丰富性。 4. 生成式推荐新范式，为构建更具泛化能力、可解释性和灵活性的推荐系统开辟了新方向。

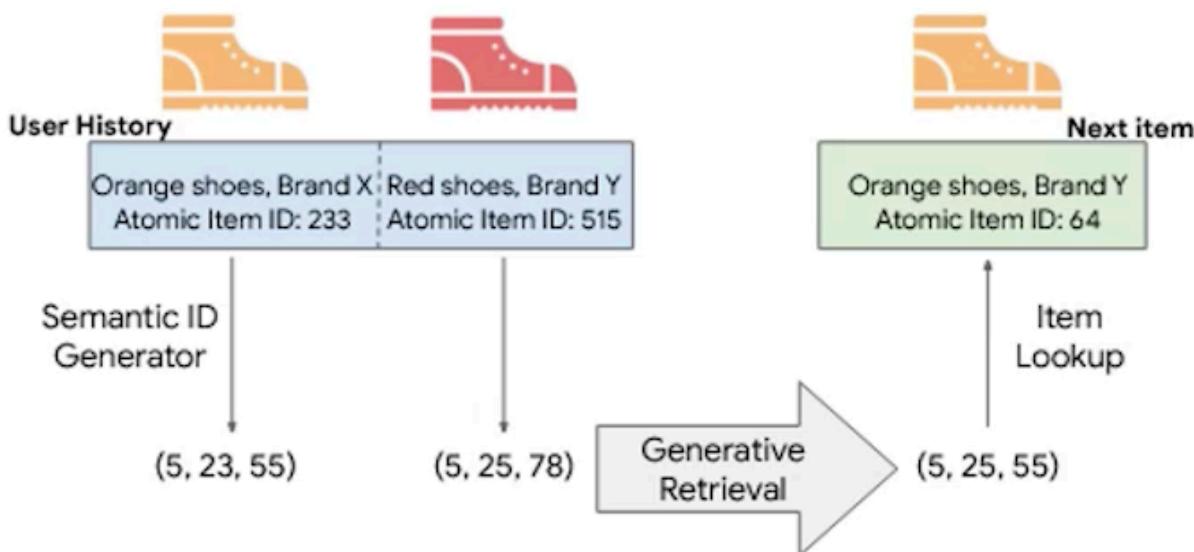


Figure 1: Overview of the *Transformer Index for GEnenerative Recommenders* (TIGER) framework. With TIGER, sequential recommendation is expressed as a generative retrieval task by representing each item as a tuple of discrete semantic tokens.

提出了一种构建序列推荐生成式检索模型的新范式

与传统的检索-排序（候选匹配）方法不同，我们的方法使用直接预测候选物品ID的端到端生成模型。我们提出利用Transformer内存（参数）作为推荐系统中检索的端到端索引引擎。我们将提出的方法称为Transformer Index for GEnenerative Recommenders（TIGER）。图1展示了TIGER框

架的整体工作流程，说明了如何将序列推荐任务转化为一个生成式检索任务。

TIGER [Transformer Index for GEnerative RecommenDers] 推荐系统的建模流程，包括两个主要阶段：(a) 语义ID的生成；(b) 基于Transformer的生成式推荐模型训练

具体来说，给定物品的文本特征，我们使用预训练的文本编码器（例如SentenceT5）来生成密集的内容嵌入。然后将量化方案应用于物品的嵌入，形成一组离散有序的token，我们称之为物品的语义ID。最终，这些语义ID用于在序列推荐任务上来训练Transformer模型。

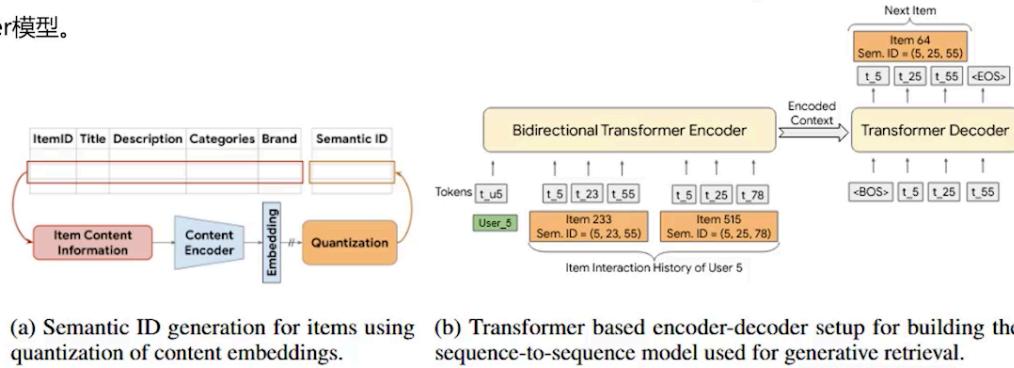


Figure 2: An overview of the modeling approach used in TIGER.

语义ID eg (2, 3, 55) 每一个数字都代表着一个稠密向量

将物品/item 表示为语义ID序列具有许多优点：

- ✓ 在具有语义意义的数据上训练Transformer允许在相似的物品之间知识共享，避免在推荐模型中使用原子、随机的物品ID作为物品特征，从而支持知识迁移和泛化。
- ✓ 使用物品的语义ID，能够缓解模型受到的推荐系统中固有反馈循环的影响，减少推荐系统对热门物品的依赖，并允许模型泛化到语料库中的新物品，更好地支持对新物品的推荐。
- ✓ 通常，物品的规模可以在数十亿的数量级。使用有限数量的代码词组合生成的语义ID来表示物品可以显著降低了物品表示的存储成本和参数规模，提升了系统的可扩展性。

本工作的主要贡献总结如下：

1. 我们提出了一个新的基于生成检索的推荐框架 TIGER，它为每个物品分配语义ID，并训练检索模型来预测生成给定用户可能交互的物品的语义ID。
2. 通过召回率和 NDCG 指标，我们证明 TIGER 在多个数据集上优于现有的 SOTA 推荐系统。
3. 我们发现这种生成式检索的新范式在序列推荐系统中带来了两个额外能力：
 - 能够推荐新的和冷门的物品，从而改进冷启动推荐。
 - 能够使用可调参数生成多样的推荐结果。

相关工作

SASRec

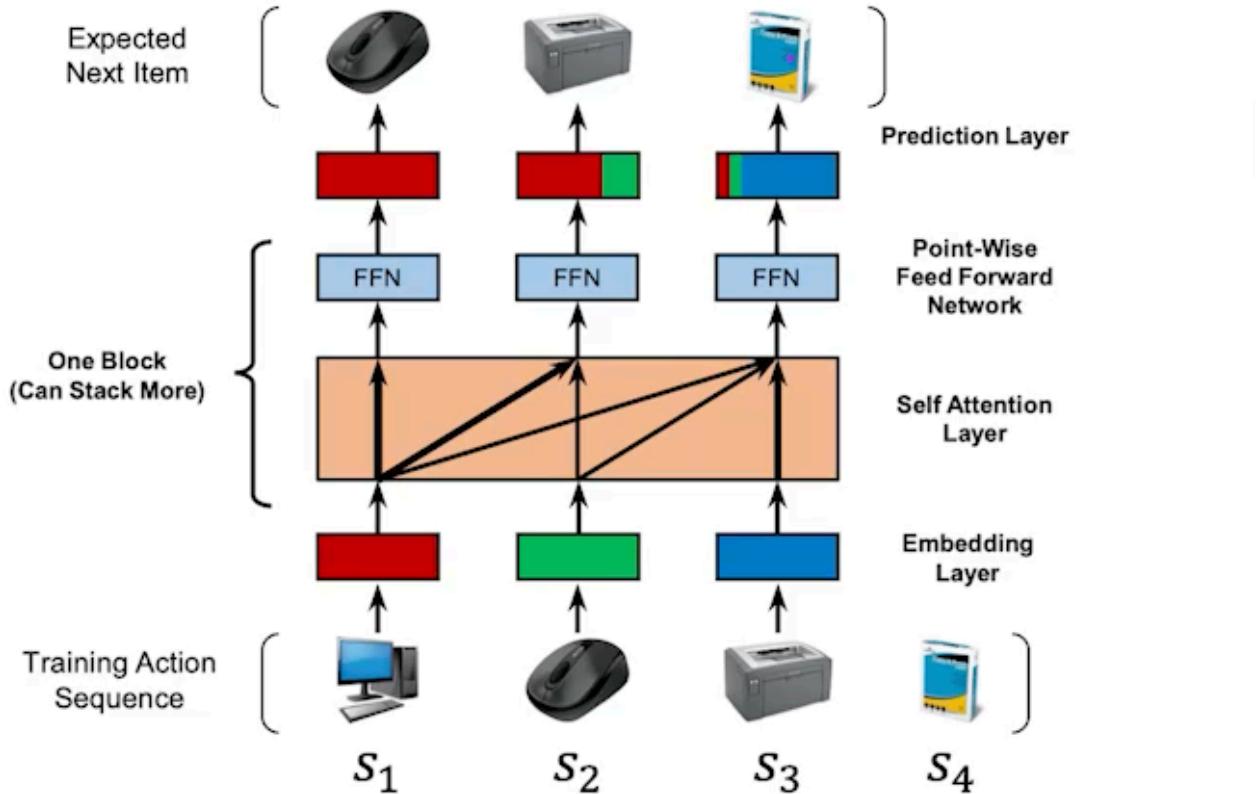


Figure 1: A simplified diagram showing the training process of SASRec. At each time step, the model considers all previous items, and uses attention to ‘focus on’ items relevant to the next action.

SASRec 将 Transformer 中的自注意力机制应用在序列推荐中，基于给定的物品序列来预测下一个最可能出现的物品。

采用自注意力机制来对用户的历史行为信息建模，从而提取更有价值的信息。最后将得到的用户表征分别与所有的物品 Embedding 内容做内积，根据相关性大小排序、筛选，得到 Top-k 个推荐。

该方法在稀疏和稠密数据集上都有较为突出的性能。

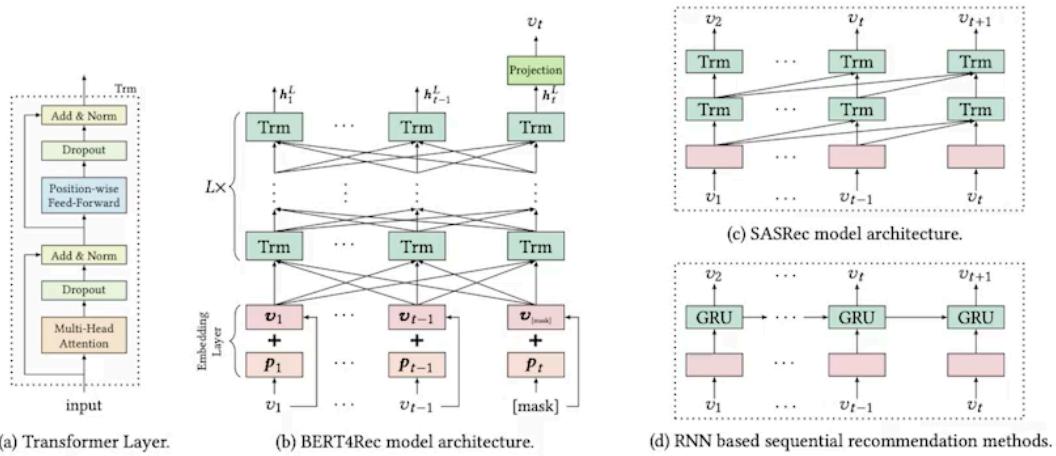


Figure 1: Differences in sequential recommendation model architectures. BERT4Rec learns a bidirectional model via Cloze task, while SASRec and RNN based methods are all left-to-right unidirectional model which predict next item sequentially.

CSDN @frostjy

之前的算法使用序列神经网络从左向右地编码用户的历史交互信息，只利用了单向的信息进行建模。尽管它们是有效的，但是存在两点问题：

- 单向的模型结构限制了用户行为序列的隐式表征
- 它们通常假定一个严格有序的序列，而这并不总是实用的

为了解决这些问题，BERT4Rec 模型采用深层的双向 **Self-Attention** 来对用户行为序列进行建模。

因为用户的兴趣不一定是单向的，有可能转变

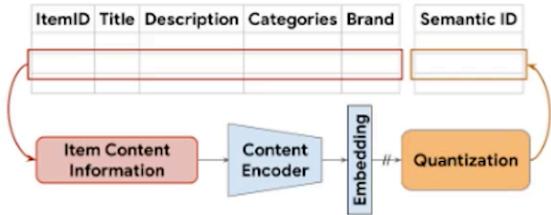
- **GRU[2015]**: 首个将定制化 GRU（门控循环单元）用于序列推荐任务的基于 RNN 的方法。
- **Caser[2018]**: 该模型采用 CNN 架构，通过应用横向和垂直卷积操作来捕捉高阶马尔可夫链，以实现序列推荐。
- **HGN[2019]**: 层次门控网络利用一种新的门控架构以捕捉用户的长期与短期兴趣。
- **SASRec[2018]**: 自注意力序列推荐模型采用因果掩码 Transformer，为用户的序列交互行为建模。
- **BERT4Rec[2019]**: BERT4Rec 采用双向自注意力 Transformer，克服了单向模型的缺陷，并将其应用于推荐任务。
- **FDSA[2019]**: 特征级深度自注意力网络将物品特征与物品嵌入一同引入，作为 Transformer 的输入序列。
- **S³-Rec[2020]**: 序列推荐的自监督学习技术，引入 4 个自监督任务的双向 Transformer 预训练，以提升模型性能。
- **P5[2022]**: P5 是一项新近提出的技术，它通过使用预训练大语言模型，将多种推荐任务整合到一个统一的框架中。

P5 模型依赖 LLM tokenizer 从随机分配的物品 ID 中生成 token，形成语义 ID。

我们将语义ID定义为长度为m的token元组。

元组中的每个token来自不同的codebook。因此，语义id可以唯一表示的项数等于每层codebook大小的乘积。虽然生成语义ID的不同技术导致ID具有不同的语义属性，但我们希望它们至少具有以下属性：

相似的项（具有相似内容特征或语义嵌入接近的项）应该具有重叠的语义id。



例如，语义ID [10,21,35] 的项应该与语义ID [10,21,40] 的项更相似，而不是具有ID [10,23,32] 的项。

(a) Semantic ID generation for items using quantization of content embeddings.

语义ID与大语言模型的tokenizer的区别



1. 它是离散化的语义表示，不依赖具体物品ID，而是基于内容语义生成；
2. 它的组成单位（codewords）来自固定大小的代码簿，大小可控、可重复利用；
3. 它支持组合式表达：比如每个ID由3个token组成，每个token选自1000个候选，那总共可表达10亿个组合，足以覆盖大规模推荐系统中的物品空间。



通过语义ID的设计，使得推荐系统能够像语言模型那样“输出”目标物品，但又规避了推荐领域中item ID/token空间过大、实体变化频繁等根本难题。可以说，语义ID是TIGER成功的核心支柱之一。

RQ-VAE: 多层次量化机制



残差量化变分自动编码器 (RQ-VAE) 是一个多级矢量量化器，它对残差进行量化以生成语义id。通过更新量化codebook和DNN编码器参数对自编码器进行联合训练。

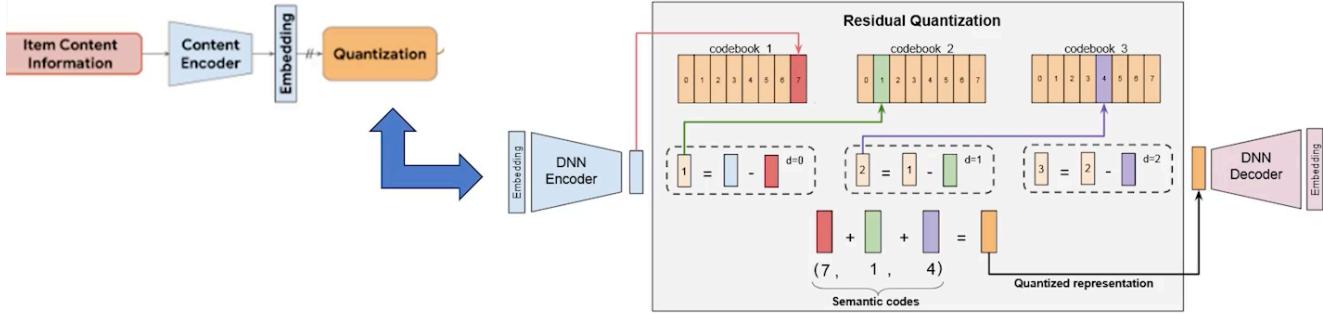


Figure 3: RQ-VAE: In the figure, the vector output by the DNN Encoder, say r_0 (represented by the blue bar), is fed to the quantizer, which works iteratively. First, the closest vector to r_0 is found in the first level codebook. Let this closest vector be e_{c_0} (represented by the red bar). Then, the residual error is computed as $r_1 := r_0 - e_{c_0}$. This is fed into the second level of the quantizer, and the process is repeated: The closest vector to r_1 is found in the second level, say e_{c_1} (represented by the green bar), and then the second level residual error is computed as $r_2 = r_1 - e_{c_1}$. Then, the process is repeated for a third time on r_2 . The semantic codes are computed as the indices of e_{c_0}, e_{c_1} , and e_{c_2} in their respective codebooks. In the example shown in the figure, this results in the code (7, 1, 4).

基于RQ-VAE的semantic ID



RQ-VAE first encodes the input x via an encoder \mathcal{E} to learn a latent representation $z := \mathcal{E}(x)$.

At the zero-th level ($d = 0$), the initial residual is simply defined as $r_0 := z$. At each level d , we have a codebook $\mathcal{C}_d := \{e_k\}_{k=1}^K$, where K is the codebook size.

Then, r_0 is quantized by mapping it to the nearest embedding from that level's codebook. The index of the closest embedding e_{c_d} at $d = 0$, i.e., $c_0 = \arg \min_k \|r_0 - e_k\|$, represents the zero-th codeword.

For the next level $d = 1$, the residual is defined as $r_1 := r_0 - e_{c_0}$.

Then, similar to the zero-th level, the code for the first level is computed by finding the embedding in the codebook for the first level which is nearest to r_1 .

This process is repeated recursively m times to get a tuple of m codewords that represent the Semantic ID.

Once we have the Semantic ID (c_0, \dots, c_{m-1}) , a quantized representation of \mathbf{z} is computed as $\hat{\mathbf{z}} := \sum_{d=0}^{m-1} e_{c_i}$.

Then $\hat{\mathbf{z}}$ is passed to the decoder, which tries to recreate the input \mathbf{x} using $\hat{\mathbf{z}}$. The RQ-VAE loss is defined as

$$\mathcal{L}(\mathbf{x}) := \mathcal{L}_{\text{recon}} + \mathcal{L}_{\text{rqvae}}$$

where $\mathcal{L}_{\text{recon}} := \|\mathbf{x} - \hat{\mathbf{x}}\|^2$, and $\mathcal{L}_{\text{rqvae}} := \sum_{d=0}^{m-1} \|\text{sg}[r_i] - e_{c_i}\|^2 + \beta \|r_i - \text{sg}[e_{c_i}]\|^2$.

Here $\hat{\mathbf{x}}$ is the output of the decoder, and sg is the stop-gradient operation. This loss jointly trains the encoder, decoder, and the codebook.

To prevent RQ-VAE from a codebook collapse, where most of the input gets mapped to only a few codebook vectors, we use k-means clustering-based initialization for the codebook. Specifically, we apply the k-means algorithm on the first training batch and use the centroids as initialization.

语义ID生成

其它量化方法：

LSH: LSH (Locality Sensitive Hashing) : 一种简单但粗糙的离散化方式，速度快但精度低消融实验证实RQ-VAE性能更优

VQ-VAE: 不具有语义层次粒度

基于k-means层次聚类：丢失ID之间的语义关系；

冲突处理：

附加标识位：增加一位token 表示语义重复。

查找表：维护一个映射表，将语义ID映射回具体的物品ID，避免模型误判。

只需训练后处理一次：冲突检测和修复仅在训练好RQ-VAE之后执行一次，不影响训练效率。



Table 3: The entropy of the category distribution predicted by the model for the Beauty dataset. A higher entropy corresponds more diverse items predicted by the model.

Temperature	Entropy@10	Entropy@20	Entropy@50
T = 1.0	0.76	1.14	1.70
T = 1.5	1.14	1.52	2.06
T = 2.0	1.38	1.76	2.28

Table 4: Recommendation diversity with temperature-based decoding.

Target Category	Most-common Categories for top-10 predicted items	
	T = 1.0	T = 2.0
Hair Styling Products	Hair Styling Products	Hair Styling Products, Hair Styling Tools, Skin Face
Tools Nail	Tools Nail	Tools Nail, Makeup Nails
Makeup Nails	Makeup Nails	Makeup Nails, Skin Hands & Nails, Tools Nail
Skin Eyes	Skin Eyes	Hair Relaxers, Skin Face, Hair Styling Products, Skin Eyes
Makeup Face	Tools Makeup Brushes, Makeup Face	Tools Makeup Brushes, Makeup Face, Skin Face, Makeup Sets, Hair Styling Tools
Hair Loss Products	Hair Loss Products, Skin Face, Skin Body	Skin Face, Hair Loss Products, Hair Shampoos, Hair & Scalp Treatments, Hair Conditioners

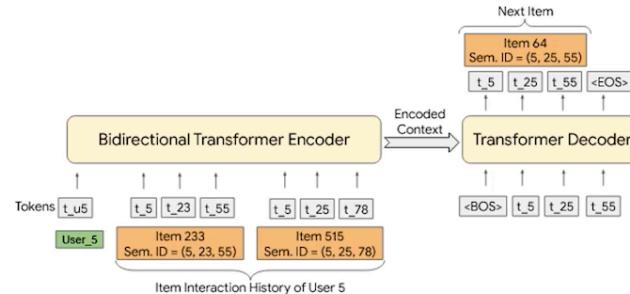
- Entropy@K，计算top-K物品的ground-truth的类别分布；
- 解码器过程中基于温度系数的采样 temperature-based sampling可以应用于任何现有的推荐模型来控制结果多样性。但由于RQ-VAE语义id的特性，TIGER允许跨不同层次的采样。

基于语义ID的生成式检索

我们通过按时间顺序排序的用户历史交互的物品来为每个用户构建物品序列。

然后，给定一个形式为 $(item_1, \dots, item_n)$ 的序列，推荐系统的任务是预测下一个物品 $item_{n+1}$ 。我们提出了一种直接预测下一个物品的语义ID的生成方法。

形式上，设 $(c_{i,0}, \dots, c_{i,m-1})$ 为 $item_i$ 的 m 长度的语义ID。然后我们将物品序列转换为序列 $(c_{1,0}, \dots, c_{1,m-1}, c_{2,0}, \dots, c_{2,m-1}, \dots, c_{n,0}, \dots, c_{n,m-1})$ 。然后训练序列到序列模型来预测物品的语义ID，即 $(c_{n+1,0}, \dots, c_{n+1,m-1})$ 。考虑到我们框架的生成特性，从解码器生成的语义ID可能与推荐语料库中的物品不匹配，但此类事件发生的概率很低。



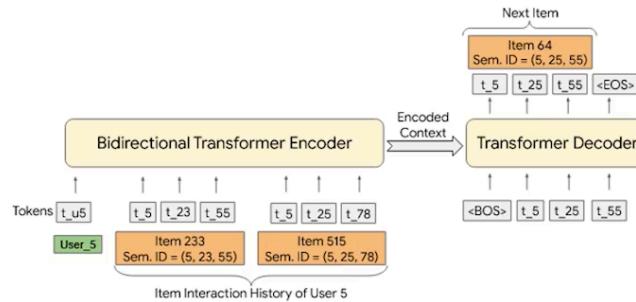
(b) Transformer based encoder-decoder setup for building the sequence-to-sequence model used for generative retrieval.

我们通过按时间顺序排序的用户历史交互的物品来为每个用户构建物品序列。

然后，给定一个形式为 $(item_1, \dots, item_n)$ 的序列，推荐系统的任务是预测下一个物品 $item_{n+1}$ 。我们提出了一种直接预测下一个物品的语义ID的生成方法。

形式上，设 $(c_{i,0}, \dots, c_{i,m-1})$ 为 $item_i$ 的 m 长度的语义ID。然后，我们将物品序列转换为序列 $(c_{1,0}, \dots, c_{1,m-1}, c_{2,0}, \dots, c_{2,m-1}, \dots, c_{n,0}, \dots, c_{n,m-1})$ 。然后训练序列到序列模型来预测物品的语义ID，即 $(c_{n+1,0}, \dots, c_{n+1,m-1})$ 。考虑到我们框架的生成特性，从解码器生成的语义ID可能与推荐语料库中的物品不匹配，但此类事件发生的概率很低。

通过语义ID的设计，TIGER框架成功地将推荐系统引入了一个**生成式检索**的新范式，使得系统能够像语言模型那样直接“输出”目标物品。



(b) Transformer based encoder-decoder setup for building the sequence-to-sequence model used for generative retrieval.

实验设置

- **基座模型:** Sentence-T5 (获得物品的语义嵌入：768维)
- **数据集:** Amazon Product Reviews (1996.5-2014.7)
 - Beauty
 - Sports and Outdoors
 - Toys and Games
- **评估指标**
 - Recall@k
 - NDCG@k (K=5, or 10)

Table 6: Dataset statistics for the three real-world benchmarks.

Dataset	# Users	# Items	Sequence Length	
			Mean	Median
Beauty	22,363	12,101	8.87	6
Sports and Outdoors	35,598	18,357	8.32	6
Toys and Games	19,412	11,924	8.63	6

- **GRU[2015]** : 首个将定制化GRU（门控循环单元）用于序列推荐任务的基于RNN的方法
- **Caser[2018]** : 该模型采用CNN架构，通过应用横向和垂直卷积操作来捕捉高阶马尔可夫链，以实现序列推荐。
- **HGN[2019]** : 层次门控网络利用一种新的门控架构以捕捉用户的长期与短期兴趣。
- **SASRec[2018]** : 自注意力序列推荐模型采用因果掩码Transformer，为用户的序列交互行为建模
- **BERT4Rec[2019]** : BERT4Rec采用双向自注意力Transformer，克服了单向模型的缺陷，并将其应用于推荐任务
- **FDSA[2019]** : 特征级深度自注意力网络将物品特征与物品嵌入一同引入，作为Transformer的输入序列
- **S³-Rec[2020]** : 序列推荐的自监督学习技术引入基于自监督任务的双向Transformer预训练，以提升模型性能。
- **P5[2022]** : P5是一项新近提出的技术，它通过使用预训练大语言模型，将多种推荐任务整合到一个统一的框架中。

实验结果

Table 1: Performance comparison on sequential recommendation. The last row depicts % improvement with TIGER relative to the best baseline. Bold (underline) are used to denote the best (second-best) metric.

Methods	Sports and Outdoors				Beauty				Toys and Games			
	Recall @5	NDCG @5	Recall @10	NDCG @10	Recall @5	NDCG @5	Recall @10	NDCG @10	Recall @5	NDCG @5	Recall @10	NDCG @10
P5 [8]	0.0061	0.0041	0.0095	0.0052	0.0163	0.0107	0.0254	0.0136	0.0070	0.0050	0.0121	0.0066
Caser [33]	0.0116	0.0072	0.0194	0.0097	0.0205	0.0131	0.0347	0.0176	0.0166	0.0107	0.0270	0.0141
HGN [25]	0.0189	0.0120	0.0313	0.0159	0.0325	0.0206	0.0512	0.0266	0.0321	0.0221	0.0497	0.0277
GRU4Rec [11]	0.0129	0.0086	0.0204	0.0110	0.0164	0.0099	0.0283	0.0137	0.0097	0.0059	0.0176	0.0084
BERT4Rec [32]	0.0115	0.0075	0.0191	0.0099	0.0203	0.0124	0.0347	0.0170	0.0116	0.0071	0.0203	0.0099
FDSA [42]	0.0182	0.0122	0.0288	0.0156	0.0267	0.0163	0.0407	0.0208	0.0228	0.0140	0.0381	0.0189
SASRec [17]	0.0233	0.0154	0.0350	0.0192	0.0387	<u>0.0249</u>	0.0605	0.0318	<u>0.0463</u>	<u>0.0306</u>	0.0675	0.0374
S ³ -Rec [44]	0.0251	<u>0.0161</u>	<u>0.0385</u>	<u>0.0204</u>	<u>0.0387</u>	0.0244	<u>0.0647</u>	<u>0.0327</u>	0.0443	0.0294	<u>0.0700</u>	0.0376
TIGER [Ours]	0.0264	0.0181	0.0400	0.0225	0.0454	0.0321	0.0648	0.0384	0.0521	0.0371	0.0712	0.0432
	+5.22%	+12.55%	+3.90%	+10.29%	+17.31%	+29.04%	+0.15%	+17.43%	+12.53%	+21.24%	+1.71%	+14.97%

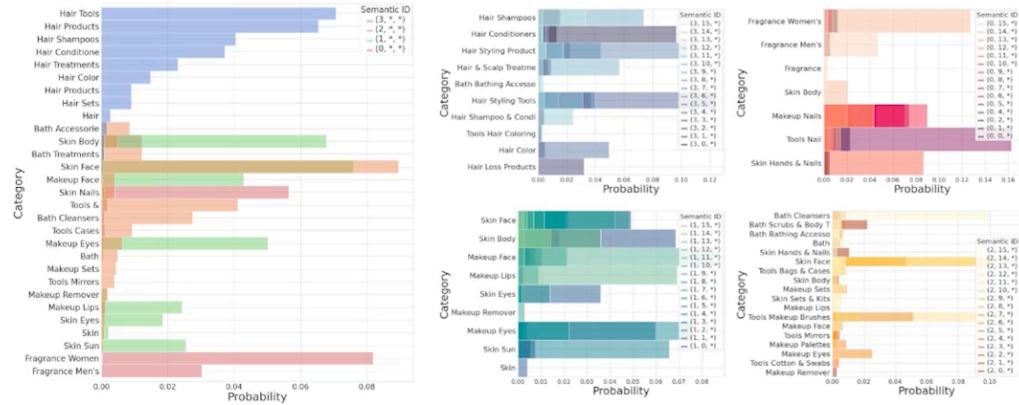
- TIGER在各个数据集和评测指标上取得了最佳效果，特别是在 Beauty数据集上急NDCG指标上

Table 9: The mean and stand error of the metrics for different dataset (computed using 3 runs with different random seeds)

Datasets	Recall@5	NDCG@5	Recall@10	NDCG@10
Beauty	0.0441 ± 0.00069	0.0309 ± 0.00062	0.0642 ± 0.00092	0.0374 ± 0.00061
Sports and Outdoors	0.0278 ± 0.00069	0.0189 ± 0.00043	0.0419 ± 0.0010	0.0234 ± 0.00048
Toys and Games	0.0518 ± 0.00064	0.0375 ± 0.00039	0.0698 ± 0.0013	0.0433 ± 0.00047

- TIGER在三轮使用不同随机数种子的实验中，指标的均值和标准差的结果

商品表示：语义ID质量分析



(a) The ground-truth category distribution for all the items in the dataset colored by the value of the first codeword c_1 .

(b) The category distributions for items having the Semantic ID as $(c_1, *, *)$, where $c_1 \in \{1, 2, 3, 4\}$. The categories are color-coded based on the second semantic token c_2 .

Figure 4: Qualitative study of RQ-VAE Semantic IDs (c_1, c_2, c_3, c_4) on the Amazon Beauty dataset. We show that the ground-truth categories are distributed across different Semantic tokens. Moreover, the RQVAE semantic IDs form a hierarchy of items, where the first semantic token (c_1) corresponds to coarse-level category, while second/third semantic token (c_2/c_3) correspond to fine-grained categories.

消融实验：不同ID生成方式

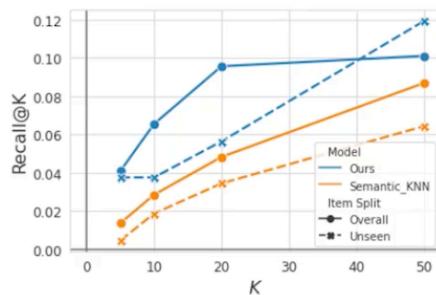


Table 2: Ablation study for different ID generation techniques for generative retrieval. We show that RQ-VAE Semantic ID (SID) perform significantly better compared to hashing SIDs and Random IDs.

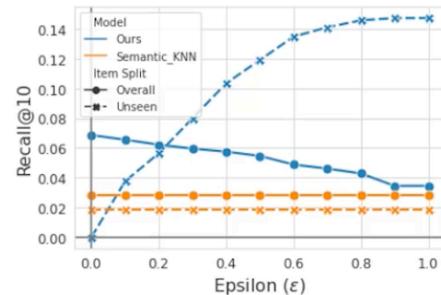
Methods	Sports and Outdoors				Beauty				Toys and Games			
	Recall @5	NDCG @5	Recall @10	NDCG @10	Recall @5	NDCG @5	Recall @10	NDCG @10	Recall @5	NDCG @5	Recall @10	NDCG @10
Random ID	0.007	0.005	0.0116	0.0063	0.0296	0.0205	0.0434	0.0250	0.0362	0.0270	0.0448	0.0298
LSH SID	0.0215	0.0146	0.0321	0.0180	0.0379	0.0259	0.0533	0.0309	0.0412	0.0299	0.0566	0.0349
RQ-VAE SID	0.0264	0.0181	0.0400	0.0225	0.0454	0.0321	0.0648	0.0384	0.0521	0.0371	0.0712	0.0432

- RQ-VAE基本优于局部敏感哈希 (LSH)。
- 给定相同的基于内容的语义嵌入，通过非线性深度神经网络 (DNN) 架构学习语义ID比使用随机投影产生更好的量化。
- 语义ID始终优于随机ID基线，突出了利用基于内容的语义信息的重要性。

冷启动推荐



(a) Recall@K vs. K, ($\epsilon = 0.1$).



(b) Recall@10 vs. ϵ .

Figure 5: Performance in the cold-start retrieval setting.

- 基于 “Beauty” 数据集上的对比
- 将测试物品中的 5 % 从训练集中移除，作为 “未见过” 物品； ϵ 表示TIGER框架中能够选择到的“未见过”物品比例
- 与KNN相比，在 $\epsilon=0.1$ ，TIGER在所有K值上的Recall性能更优；针对不同 ϵ ，也显示出性能更优。
- 相比之下，TIGER框架可以很容易地执行冷启动推荐，因为它在预测下一个物品时利用了物品语义。

Table 3: The entropy of the category distribution predicted by the model for the Beauty dataset. A higher entropy corresponds more diverse items predicted by the model.

Temperature	Entropy@10	Entropy@20	Entropy@50
T = 1.0	0.76	1.14	1.70
T = 1.5	1.14	1.52	2.06
T = 2.0	1.38	1.76	2.28

Table 4: Recommendation diversity with temperature-based decoding.

Target Category	Most-common Categories for top-10 predicted items	
	T = 1.0	T = 2.0
Hair Styling Products	Hair Styling Products	Hair Styling Products, Hair Styling Tools, Skin Face
Tools Nail	Tools Nail	Tools Nail, Makeup Nails
Makeup Nails	Makeup Nails	Makeup Nails, Skin Hands & Nails, Tools Nail
Skin Eyes	Skin Eyes	Hair Relaxers, Skin Face, Hair Styling Products, Skin Eyes
Makeup Face	Tools Makeup Brushes, Makeup Face	Tools Makeup Brushes, Makeup Face, Skin Face, Makeup Sets, Hair Styling Tools
Hair Loss Products	Hair Loss Products, Skin Face, Skin Body	Skin Face, Hair Loss Products, Hair Shampoos, Hair & Scalp Treatments, Hair Conditioners

- Entropy@K，计算top-K物品的ground-truth的类别分布；
- 解码器过程中基于温度系数的采样 temperature-based sampling可以应用于任何现有的推荐模型来控制结果多样性，但由于RQ-VAE语义id的特性，TIGER允许跨不同层次的采样。

消融实验：模型层数及用户信息的作用

Table 5: Recall and NDCG metrics for different number layers.

Number of Layers	Recall@5	NDCG@5	Recall@10	NDCG@10
3	0.04499	0.03062	0.06699	0.03768
4	0.0454	0.0321	0.0648	0.0384
5	0.04633	0.03206	0.06596	0.03834

Table 8: The effect of providing user information to the recommender system

Recall@5	NDCG@5	Recall@10	NDCG@10
No user information	0.04458	0.0302	0.06479
With user id (reported in the paper)	0.0454	0.0321	0.0648

- 当网络变大（层数变多），模型性能有轻微的提升。
- 为语言模型提供用户信息对性能有一定的帮助

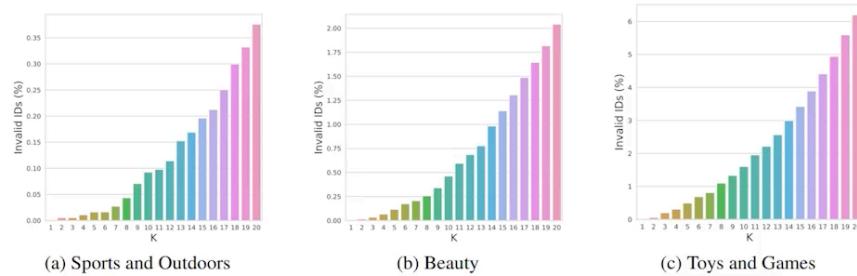


Figure 6: Percentage of invalid IDs when generating Semantic IDs using Beam search for various values of K . As shown, $\sim 0.3\% - 6\%$ of the IDs are invalid when retrieving the top-20 items.

- 由于模型是自回归地解码生成目标语义ID，因此模型可能会预测出无效的ID（即在推荐数据集中没有对应物品的ID）。我们观察到模型几乎总是预测出有效的ID。在图6中，对于前10个预测结果，三个数据集中无效ID的比例在约0.1%至1.6%之间变化。可以增加beam search波束搜索的大小并过滤掉无效ID。
- 尽管会生成无效ID，但与用于顺序推荐的其他流行方法相比，TIGER仍然实现了最优的性能。一种处理无效token的方法是，当模型生成无效token时进行前缀匹配。RQ-VAE token具有层次结构，前缀匹配可以看作是模型预测物品类别而不是物品索引，这样的扩展可以进一步提高召回率/NDCG指标。

其它讨论

Table 10: Testing scalability by generating the Semantic IDs on the combined dataset vs generating the Semantic IDs on only the Beauty dataset.

	Recall@5	NDCG@5	Recall@10	NDCG@10
Semantic ID [Combined datasets]	0.04355	0.3047	0.06314	0.03676
Semantic ID [Amazon Beauty]	0.0454	0.0321	0.0648	0.0384

- 语义ID长度和码本大小的影响：**通过试验（尝试6-元组的语义ID，每层64个codewords），TIGER的推荐指标对语义ID的长度和codebook大小的变化具有鲁棒性。然而随着ID变长，输入序列的长度也会增加，这会使基于Transformer的序列到序列模型的计算成本更高。
- 可扩展性：**将所有三个数据集合并，与仅从Beauty数据集生成语义ID的原始实验结果进行比较的结果如表10所示。我们发现，性能仅有略有下降。
- 推理成本：**由于使用了自回归解码的beam search，模型在推理期间的计算成本可能比基于ANN的模型更高。论文开辟了一个新的研究领域：基于生成检索的推荐系统。我们将考虑使模型更小的方法或探索其他提高推理效率的方法。
- 查找表的内存成本：**我们为TIGER维护了两个查找哈希表：一个是从物品ID到语义ID的表，另一个是从语义ID到物品ID的表。它们在基于RQ-VAE的语义ID生成模型训练后生成，每个语义ID由一个4个整数的元组组成，每个查找表的大小将为64N位左右，其中N是数据集中的物品数量。
- 嵌入表的内存成本：**与传统推荐系统相比，TIGER使用的嵌入表要小得多。这是因为传统推荐系统为每个物品存储一个嵌入，而TIGER只为每个语义码字存储一个嵌入。

- 本文提出了一种新的推荐范式，称为TIGER，使用生成检索模型在序列推荐中“生成”下一个可能的交互对象。
- 支撑这种方法的是一种新的物品语义ID表示，它在内容嵌入上使用层次量化器（RQ-VAE）来生成形成语义ID。
- 基于Transformer的序列生成模型：将用户的历史行为表示为语义ID序列，使用Encoder-Decoder结构的Transformer模型学习用户偏好，并生成下一个物品的语义ID。
- 嵌入表的基数不会随着物品空间的基数线性增长，这与需要在训练期间创建大型嵌入表或为每个单独物品生成索引的系统相比，更为有利。
- 通过在三个数据集上的实验，证明了我们的模型可以达到由于SOTA检索的性能，同时可以泛化到新的和未见过的物品。

不足

论文中的语义ID建模是静态的，没有与user-item交互相关，也就是没有融合协作信息，后续的生成式模型有不少改进空间。

基于生成式检索的推荐：发展前沿



<i>Recommender Systems with Generative Retrieval</i>	2023 NeurIPS	谷歌提出TIGER框架，首次将生成式检索用于推荐系统，通过语义ID序列，利用Transformer自回归生成候选集，
<i>Actions Speak Louder than Words: Trillion-Parameter Sequential Transducers for Generative Recommendations</i>	2024 ICML	Meta提出生成式推荐新范式，设计HSTU架构处理高基数动态数据，将推荐任务转化为序列变换问题，
<i>MTGR: Industrial-Scale Generative Recommendation Framework in Meituan</i>	2025 CIKM	美团设计工业级生成框架，提升多语义建模，
<i>Sparse Meets Dense: Unified Generative Recommendations with Cascaded Sparse-Dense Representations</i>	2025 arXiv	+8618631625879 百度提出COBRA框架，在广告场景显著优化业务指标
<i>OneRec: Unifying Retrieve and Rank with Generative Recommender and Preference Alignment</i>	2025 arXiv	快手提出端到端单阶段生成框架，采用MoE扩展模型容量 在线实验提升总观看时长1.68%

基于生成式检索的推荐：发展前沿



<i>Learnable Item Tokenization for Generative Recommendation</i>	对齐推荐
<i>DAS: Dual-Aligned Semantic IDs Empowered Industrial Recommender System</i>	对齐推荐
<i>OneRec Technical Report - RQ-KMeans Towards Scalable Semantic Representation for Recommendation</i>	多套semantic ID
<i>Generating Long Semantic IDs in Parallel for Recommendation</i>	多套semantic ID
<i>MMQ: Multimodal Mixture-of-Quantization Tokenization for Semantic ID Generation and User Behavioral Adaptation</i>	+8618631625879 多套semantic ID