# Tokenize Once, Recommend Anywhere: Unified Item Tokenization for Multi-domain LLM-based Recommendation

**Yu Hou[1], Won-Yong Shin[1*]**

[1]Yonsei University
{houyu, wy.shin}@yonsei.ac.kr

## Abstract

Large language model (LLM)-based recommender systems have achieved high-quality performance by bridging the discrepancy between the item space and the language space through item tokenization. However, existing item tokenization methods typically require training separate models for each item domain, limiting generalization. Moreover, the diverse distributions and semantics across item domains make it difficult to construct a unified tokenization that preserves domain-specific information. To address these challenges, we propose **UniTok**, a **Uni**fied item **Tok**enization framework that integrates our own mixture-of-experts (MoE) architecture with a series of codebooks to convert items into discrete tokens, enabling scalable tokenization while preserving semantic information across *multiple item domains*. Specifically, items from different domains are first projected into a unified latent space through a shared encoder. They are then routed to *domain-specific* experts to capture the *unique* semantics, while a *shared* expert, which is always active, encodes common knowledge transferable across domains. Additionally, to mitigate *semantic imbalance* across domains, we present a mutual information calibration mechanism, which guides the model towards retaining similar levels of semantic information for each domain. Comprehensive experiments on wide-ranging real-world datasets demonstrate that the proposed UniTok framework is **(a) highly effective:** achieving up to 51.89% improvements over strong benchmarks, **(b) theoretically sound:** showing the analytical validity of our architectural design and optimization; and **(c) highly generalizable:** demonstrating robust performance across diverse domains without requiring per-domain retraining, a capability not supported by existing baselines.

**Code** — https://github.com/jackfrost168/UniTok

## Introduction

Large language models (LLMs) have recently become a promising paradigm for generative recommendation (Rajput et al. 2023; Hua et al. 2023), leveraging their strong generalization, language understanding, and world knowledge to support personalized recommendation beyond traditional language processing tasks. To effectively use LLMs

for recommendation, items must be indexed using identifiers, a process known as item tokenization (Rajput et al. 2023). Item tokenization converts items into *discrete tokens*, such as ID-based representations (Hua et al. 2023), textual descriptors (Zhang et al. 2021), or codebook-based identifiers (Rajput et al. 2023). This bridges the gap between the item space and the language space, enabling LLMs to process items as part of natural language sequences and making generative recommendations feasible.

Existing item tokenization methods (Rajput et al. 2023; Wang et al. 2024) are primarily tailored to items in single-domain settings, necessitating the training of separate tokenizers for each item domain (hereafter, "domain" refers to "item domain" for brevity). In practice, this domain-specific design aligns with the fact that recommender systems are often deployed independently per domain; thus, users rarely perceive quality issues. However, as recommendation tasks increasingly span multiple domains, such as diverse item categories or services, this siloed approach leads to inefficiencies in training, deployment, and maintenance, ultimately hindering scalability. In contrast, other machine learning fields have seen a growing shift towards building *unified models* for multi-domain learning, driven by the need to reduce redundant training, improve parameter efficiency, and facilitate knowledge sharing across domains. Notable advances in this direction have been achieved in language processing (Gururangan et al. 2020; Raffel et al. 2020) and computer vision (Ullah et al. 2022; Jain et al. 2023), demonstrating the feasibility and value of such generalization.

Inspired by this, a natural question arising is: "How can we design a unified item tokenization framework for LLM-based recommendation that can be effectively generalized across multiple domains with minimal computational overhead?" To answer this question, we would like to outline the following two design challenges:

- **C1. Training overhead**: *Repeatedly training* domain-specific tokenizers is inefficient and resource-intensive. As shown in Figure 1a, when applied to 10 distinct domains, our UniTok method reduces the total number of trainable parameters by 9.63× compared to codebook-based item tokenization methods (Rajput et al. 2023; Wang et al. 2024), which require training a separate set of codebooks for each dataset to quantize items.

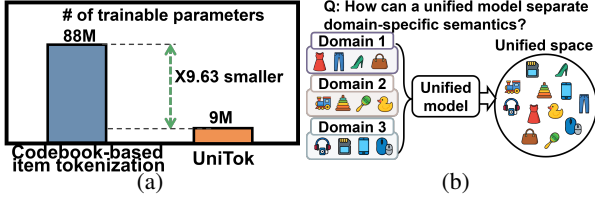- **C2. Semantic alignment**: The tokenizer must capture

Figure 1: Examples illustrating (a) a comparison of the number of trainable parameters between codebook-based item tokenization methods and our method, UniTok, when applied to 10 distinct item domains, and (b) the inherent challenge of item tokenization across multiple domains.

*rich semantics* from diverse domains; however, naïvely using a shared token space across domains can cause semantic mixing and biased token assignments. Figure 1b exemplifies this challenge.

To address these aforementioned challenges, we make the first attempt towards developing a **Uni**fied item **Tok**enization framework designed to work effectively across multiple domains, named **UniTok**.

(**Idea 1**): Different domains often have exhibit distinct data distributions, requiring models to be trained separately to capture domain-specific patterns. To move beyond this limitation, we aim to design a unified item tokenization model capable of handling multiple domains without losing domain-specific knowledge. Achieving this goal requires the model to internally disentangle domain-specific learning from shared representations. To this end, we propose a new mixture-of-experts (MoE) architecture, dubbed Token-MoE, wherein domain-specific experts specialize in modeling patterns unique to each domain, while a shared expert captures the common knowledge across multiple domains. This architectural design enables the unified model to retain domain specialization without sacrificing global knowledge sharing (solving **C1** and partially contributing to **C2**).

(**Idea 2**): To preserve item semantics during tokenization, we adopt a codebook-based approach (Rajput et al. 2023) embedded within our MoE architecture, allowing each expert to specialize in distinct semantic patterns. However, the central challenge in multi-domain settings lies in ensuring semantic balance across diverse domains (Ma et al. 2022). To address this, we introduce a mutual information (MI) calibration mechanism that explicitly encourages latent embeddings from each domain to retain sufficient and consistent informativeness. By minimizing the variance of MI across item domains, this mechanism attenuates inter-domain performance variability, enabling more stable and consistent generalization across diverse semantic spaces (solving **C2**).

Our main contributions are summarized as follows:

- **New methodology:** We propose UniTok, a unified item tokenization framework that integrates our customized MoE with codebooks to extract the semantic tokens for items across multiple domains while maintaining semantic balance.

- **Extensive evaluations:** Through comprehensive experimental evaluations on diverse real-world datasets, we

demonstrate (a) the superiority of UniTok in multi-domain scenarios, achieving substantial improvements of up to 51.89% in NDCG@10, (b) the efficiency of UniTok, with a 9.63× reduction in model size compared to competitors, and (c) the strong generalization capability of UniTok, exhibiting robust performance in zero-shot settings without additional retraining.

- **Theoretical justifications:** We theoretically prove that UniTok (a) induces a higher entropy token space, (b) achieves a lower quantization error, and (c) ensures semantic consistency across domains by reducing the MI variance, fostering stable and balanced performance.

We refer readers to the supplementary materials for technical details omitted due to page limits.

## Preliminaries

### Mixture-of-Experts

The mixture-of-experts (MoE) architecture (Jacobs et al. 1991; Shazeer et al. 2017; Fedus, Zoph, and Shazeer 2022; Dai et al. 2024) consists of multiple expert networks, each specializing in handling different parts of the input space. Formally, given an input $\mathbf{x}$, the output of an MoE model is a weighted combination of expert modules: $\text{MoE}(\mathbf{x}) = \sum_{k=1}^{K} G_k(\mathbf{x})E_k(\mathbf{x})$, where $K$ is the number of experts, $E_k(\mathbf{x})$ denotes the output of the $k$-th expert, and $G_k(\mathbf{x})$ is a softmax-based router function that assigns a probability to the $k$-th expert for given input. By dynamically routing input to specialized experts, MoE models can achieve greater model capacity and computational efficiency.

The MoE architecture is not only effective for model scaling but also well-suited for training on a large mixture of datasets (Jain et al. 2023). Its gating mechanism enables adaptive expert selection for each dataset, allowing the model to generalize across diverse sources while minimizing interference between distributions.

### Codebook-based Identifiers

Codebook-based identifiers (Rajput et al. 2023) are generated using residual quantization (RQ), which encodes item metadata into hierarchical code sequences by sequentially applying codebooks to residuals across multiple stages. Given an input vector $\mathbf{x}$, the process starts with an initial residual $\mathbf{r}^{(0)} = \mathbf{x}$, and RQ recursively quantizes it using a series of $L$ codebooks $\{C_1, C_2, \ldots, C_L\}$, where $C_\ell \triangleq \{\mathbf{c}_t\}_{t=1}^{T}$ contains $T$ code vectors at level $\ell$. The process is defined as

$$\mathbf{c}_\ell = \arg\min_{\mathbf{c} \in C_\ell} \left\| \mathbf{r}^{(\ell-1)} - \mathbf{c} \right\|^2, \tag{1}$$

$$\mathbf{r}^{(\ell)} = \mathbf{r}^{(\ell-1)} - \mathbf{c}_\ell, \tag{2}$$

where $\mathbf{r}^{(\ell)}$ is the residual at level $l$, and $\mathbf{c}_\ell$ is the nearest code vector selected from codebook $C_\ell$. The final approximation of $\mathbf{x}$ after $L$ quantization stages is given by $\hat{\mathbf{x}} = \sum_{\ell=1}^{L} \mathbf{c}_\ell$.

Each selected code vector $\mathbf{c}_\ell$ corresponds to a discrete index $z_\ell \in \{1, \cdots, T\}$, resulting in a token sequence $(z_1, \cdots, z_L)$ that compactly represents the original
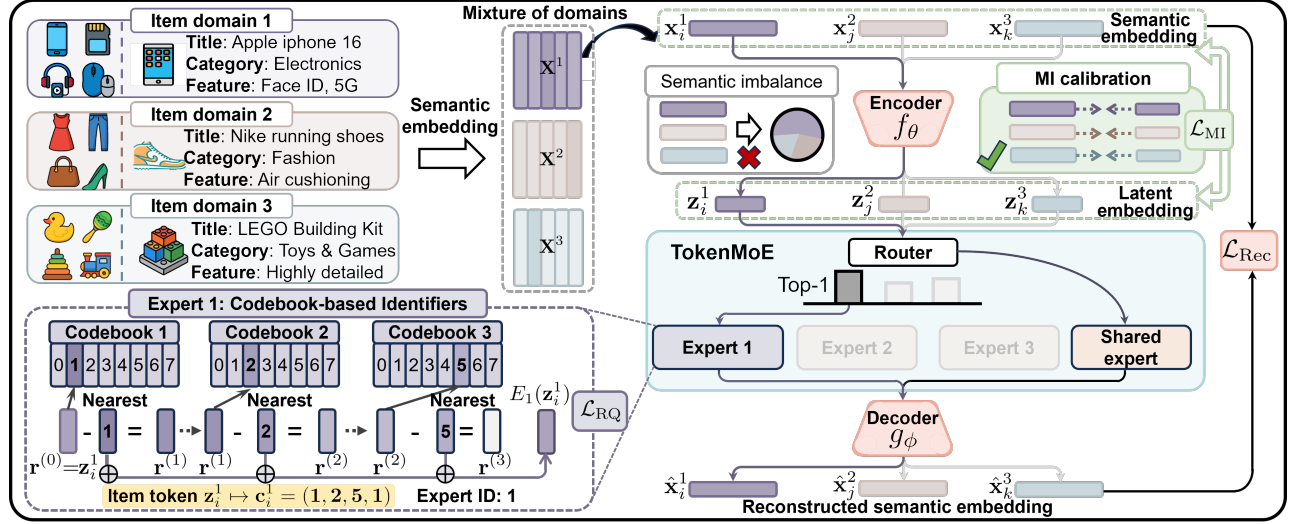
Figure 2: The schematic overview of the proposed UniTok framework.

input. This hierarchical quantization mechanism provides the foundation for our codebook-based identifiers, offering compact and semantically meaningful tokens well-suited for item tokenization. In recommender systems, such discrete tokens can then be directly used as input to LLMs (Rajput et al. 2023; Wang et al. 2024), bridging the gap between item and language spaces while preserving semantic structure.

## Methodology

### Task Formulation

**Multi-domain setting.** Let $\mathcal{D} = \{\mathcal{D}_1, \ldots, \mathcal{D}_K\}$ be a mixture of recommendation datasets for $K$ distinct item domains, where each domain $\mathcal{D}_k$ contains an item set $\mathcal{I}_k$ with associated textual metadata, such as titles, categories, features. Unlike typical single-domain scenarios, multi-domain settings pose significant challenges due to distributional inconsistencies (Jain et al. 2023), making it crucial to handle such variability in item tokenization. Addressing this issue is essential for developing a scalable and unified item tokenization method that supports multiple domains, particularly in LLM-based recommender systems. Although collaborative signals can enhance recommendation performance, they inherently rely on user–item interactions, introducing computational overhead and limiting generalization when *shared users are absent* across domains. Instead, we aim to develop a unified tokenization method that operates independently of user data, enabling scalable and general-purpose recommendations in multi-domain LLM-based systems.

**Item tokenization task.** Given raw items $\mathcal{I}_k$ from any domain-specific dataset $\mathcal{D}_k$, we assume that we have access to a pre-trained content encoder to generate the semantic embeddings for domain $k$, denoted as $\mathbf{X}^k \in \mathbb{R}^{|\mathcal{I}_k| \times d}$, where $|\mathcal{I}_k|$ is the number of items in domain $k$ and $d$ is the embedding dimensionality. The $i$-th embedded item is represented as $\mathbf{x}_i^k \in \mathbf{X}^k$. The objective is to learn a mapping function $\mathcal{F} : \mathbb{R}^d \to \mathcal{C}$ that projects each continuous embedding $\mathbf{x}_i^k$

into a discrete codeword $\mathbf{c}_i^k \in \mathcal{C}$, where $\mathcal{C}$ denotes the shared space of discrete item tokens across all domains.

## Overview of UniTok

We recall that recent item tokenization methods for LLM-based recommendations focus merely on a single-domain setting (Rajput et al. 2023; Wang et al. 2024; Zheng et al. 2024). Howerer, real-world systems recently operate across multiple domains (Jiang et al. 2022; Ning et al. 2023), leading to repeated training and semantic inconsistency across domains—an issue largely overlooked by prior studies.

To tackle challenges **C1** and **C2** in Section 1, UniTok leverages a shared autoencoder to project items in the mixture of domains into a unified latent space. To achieve effective tokenization across diverse domains, we present TokenMoE with codebook-based identifiers, an MoE architecture composed of domain-specific experts that capture specialized semantics and a shared expert that encodes generalized knowledge. Additionally, we present an MI-based loss that enforces consistent semantics across multiple domains by regulating the informativeness of latent embeddings.

### Architectural Details

As illustrated in Figure 2, our UniTok consists of four key components: a shared autoencoder, TokenMoE, codebook-based identifiers, and an MI calibration mechanism.

**Shared autoencoder.** Given semantic embeddings from the mixture of domains, we first employ a shared autoencoder composed of an encoder $f_\theta$, to project items into a unified latent space, and decoder $g_\phi$ to reconstruct the semantic embeddings (see the coral pink region in Figure 2). This establishes a unified representation space across multiple domains, which captures common structural patterns while retaining essential information for reconstruction.

Formally, for each input item $\mathbf{x}_i^k \in \mathbf{X}^k$ from domain $\mathcal{D}_k$, the encoder produces a latent embedding $\mathbf{z}_i^k = f_\theta(\mathbf{x}_i^k)$, and

the decoder reconstructs the input item as $\hat{\mathbf{x}}_i^k = g_\phi(\hat{\mathbf{z}}_i^k)$, where $\hat{\mathbf{z}}_i^k = \text{TokenMoE}\left(\mathbf{z}_i^k\right)$ (to be specified in Eq. (5)). The model is optimized using the reconstruction loss:

$$\mathcal{L}_{\text{Rec}} = \sum_{k=1}^{K} \sum_{\mathbf{x}_i^k \in \mathbf{X}^k} \left\| \mathbf{x}_i^k - \hat{\mathbf{x}}_i^k \right\|^2, \qquad (3)$$

where $\|\cdot\|$ denotes the $L_2$ norm. This reconstruction loss shapes the latent space into a compact and informative representation that retains core semantics for item tokenization.

**TokenMoE.** Conventional tokenization models applied to a unified item space often fail to capture domain-specific nuances, as they treat diverse domains uniformly, potentially leading to the loss of specialized semantic information. To address this limitation, we present TokenMoE, a more generalizable MoE architecture, which routes items to both domain-specific experts and a shared expert. This design enables domain-specific experts to learn specialized patterns, while a shared expert with always-active routing facilitates efficient knowledge transfer across multiple domains. Distinct from earlier approaches that utilize MoE within the feedforward layers of transformers (Dai et al. 2024), our key contribution lies in uniquely integrating MoE into the tokenization module to better enable domain-aware tokenization. The TokenMoE module is illustrated in the light blue region of Figure 2 (with one of the domain-specific expert highlighted in purple and the shared expert in orange). This design addresses **C1** and partially mitigates **C2**.

Specifically, after encoding, the item latent embedding $\mathbf{z}_i^k$ passes through a router function $G(\cdot)$, which produces a softmax distribution over $K$ domain-specific experts: $G(\mathbf{z}_i^k) = \{G_1, G_2, \ldots, G_K\}$, where each $G_k$ is computed as

$$G_k = \frac{\exp(s_i^{(k)})}{\sum_{j=1}^{K} \exp(s_i^{(j)})}, s_i = h(\mathbf{z}_i^k) \in \mathbb{R}^K, \qquad (4)$$

and $h(\cdot)$ is a learnable linear transformation in the router producing the router logits $s_i$, and $s_i^{(k)}$ denotes the logit corresponding to the $k$-th domain-specific expert. Here, $K$ is the total number of domain-specific experts, which is typically aligned with the number of domains.

The item is then routed to the top-$N$ domain-specific experts[1] based on the highest values of $G(\mathbf{z}_i^k)$, while it is also deterministically assigned to a shared expert. Each expert, including the shared one, is implemented as a codebook-based identifier (to be specified later). The final $\hat{\mathbf{z}}_i^k$ is computed as a weighted combination of the selected domain-specific experts and the shared expert, which is formulated as follows:

$$\begin{cases} \hat{\mathbf{z}}_i^k = \text{TokenMoE}\left(\mathbf{z}_i^k\right) = \sum_{k=1}^{K} G_k E_k(\mathbf{z}_i^k) + E_{\text{share}}(\mathbf{z}_i^k), \\ G_k = \begin{cases} G_k, & \text{if } k \in \text{Top}_N(G(\mathbf{z}_i^k)) \\ 0, & \text{otherwise} \end{cases}, \end{cases}$$
$$\tag{5}$$

---

[1] $N$ is typically set as either 1 or 2 to promote sparsity in expert activation, which significantly reduces computational overhead while preserving model capacity (Lepikhin et al. 2021; Fedus, Zoph, and Shazeer 2022).

where, $E_k(\cdot)$ denotes the $k$-th expert module, $E_{\text{share}}(\cdot)$ denotes the shared expert module, and $\text{Top}_N\left(G\left(\mathbf{z}_i^k\right)\right) = \{k_1, k_2, \ldots, k_N\}$ is the set of indices corresponding to the top-$N$ experts selected by the router. Figure 2 shows an example where the top-1 expert is selected.

TokenMoE routes each item to domain-specific experts via its learned router, while a shared expert captures common knowledge through a deterministic path. To encourage expert specialization, each expert is initialized with the mean feature of a specific domain (Wang et al. 2024), providing a strong inductive bias for domain-aware tokenization *without requiring explicit supervision*. By activating only a subset of experts per item, TokenMoE enhances scalability, maintains domain specificity, and supports better generalization.

**Codebook-based identifiers.** We adopt RQ (Rajput et al. 2023) in each expert to discretize each item into compact token sequences. Given an item latent embedding $\mathbf{z}_i^k \in \mathbb{R}^d$ produced by the shared encoder, RQ approximates it through a sequence of codebooks $\{C_1, C_2, \ldots, C_L\}$, where $L$ is the codebook size, and each codebook $C_\ell \triangleq \{\mathbf{c}_t\}_{t=1}^{T}$ contains $T$ code vectors. As shown in the bottom-left of Figure 2, at each level $\ell$, the residual $\mathbf{r}^{(\ell)}$ from the previous step is encoded using the nearest code $\mathbf{c}_\ell$ in $C_\ell$. The sum of all selected codes reconstructs the original latent embedding:

$$E_k(\mathbf{z}_i^k) \approx \sum_{\ell=1}^{L} \mathbf{c}_\ell, \quad \text{where } \mathbf{c}_\ell \in C_\ell. \qquad (6)$$

This hierarchical quantization enables fine-grained and memory-efficient tokenization, as each item is tokenized by a discrete codeword:

$$\mathbf{z}_i^k \mapsto \mathbf{c}_i^k = (z_1, \ldots, z_L, e_1, \ldots, e_N), \qquad (7)$$

where $z_\ell \in \{1, \ldots, T\}$ denotes the index of selected code vector $\mathbf{c}_\ell$ from the $\ell$-th codebook $C_\ell$ and $e_n \in \{1, \ldots, K\}$ indicates the expert ID of the $n$-th top-$N$ expert chosen by the router.

To train this quantization process, we adopt the RQ loss:

$$\mathcal{L}_{\text{RQ}} := \sum_{\ell=1}^{L} \left\| \text{sg}[\mathbf{r}^{(\ell)}] - \mathbf{c}_\ell \right\|^2 + \alpha \left\| \mathbf{r}^{(\ell)} - \text{sg}[\mathbf{c}_\ell] \right\|^2, \quad (8)$$

where $\mathbf{r}^{(\ell)}$ is the residual vector at level $\ell$, $\mathbf{c}_\ell$ is the selected code vector from $C_\ell$, $\text{sg}[\cdot]$ is the stop-gradient operator, and $\alpha$ is a balancing hyperparameter. In Eq. (8), the first term aligns the code vector to the target residual (codebook learning), and the second term forces the encoder and router to commit to the selected quantized code vector.

**MI calibration.** As the shared encoder will project all items into a unified latent embedding space, it is not straightforward for the model to precisely capture domain-specific features due to semantic imbalance. This occurs when the quality of learned latent embeddings varies significantly across diverse domains—particularly between simple and complex domains—causing semantically similar items to be assigned inconsistent tokens (Ma et al. 2022).

As another key contribution aimed at mitigating this issue, we introduce an MI mechanism to ensure that the latent space retains sufficient information from each domain.
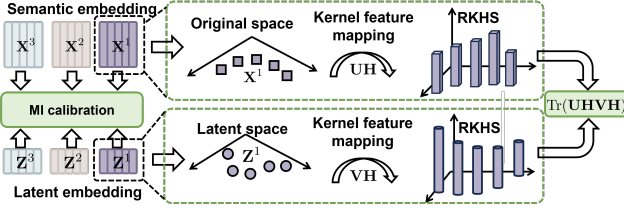
Figure 3: The illustration of MI calibration.

Specifically, we adopt the Hilbert-Schmidt independence criterion (HSIC) (Gretton et al. 2005; Li et al. 2021)[2], serving as a proxy for the MI, with a higher HSIC value indicates stronger dependence. As illustrated in Figure 3, for domain $k$, HSIC measures the dependence between the input semantic embeddings $\mathbf{X}^k = \{\mathbf{x}_1^k, \ldots, \mathbf{x}_{|\mathcal{I}_k|}^k\}$ and their latent embeddings $\mathbf{Z}^k = \{\mathbf{z}_1^k, \ldots, \mathbf{z}_{|\mathcal{I}_k|}^k\}$ in a reproducing kernel Hilbert space (RKHS), computed as:

$$\widehat{\text{HSIC}}(\mathbf{X}^k, \mathbf{Z}^k) = \frac{1}{(|\mathcal{I}_k| - 1)^2} \text{Tr}(\mathbf{UHVH}), \quad (9)$$

where $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{|\mathcal{I}_k| \times |\mathcal{I}_k|}$ are Gaussian kernel matrices computed over $\mathbf{X}^k$ and $\mathbf{Z}^k$, respectively[3]. The centering matrix $\mathbf{H} = \mathbf{I} - \frac{1}{|\mathcal{I}_k|}\mathbf{11}^\top$ ensures zero-mean embeddings in RKHS; and $\text{Tr}(\mathbf{UHVH})$ computes the Hilbert-Schmidt norm of the cross-covariance between the two RKHSs. This design addresses **C2**.

To enforce semantic balance across multiple domains (see Figure 3), we characterize the MI calibration loss as:

$$\mathcal{L}_{\text{MI}} = \text{Var}\left[\widehat{I}^{(k)}\right] - \beta\mathbb{E}\left[\widehat{I}^{(k)}\right], \quad (10)$$

where $\widehat{I}^{(k)} = \widehat{\text{HSIC}}(\mathbf{X}^k, \mathbf{Z}^k)$ and $\beta$ is a weighting hyperparameter. The first term penalizes high variance of MI across domains to mitigate semantic imbalance, while the second term enforces each domain to retain sufficient domain-specific information.

**Optimization.** We train the model using the overall loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{Rec}} + \lambda_{\text{RQ}}\mathcal{L}_{\text{RQ}} + \lambda_{\text{MI}}\mathcal{L}_{\text{MI}}, \quad (11)$$

where $\lambda_{\text{RQ}}$ and $\lambda_{\text{MI}}$ are hyper-parameters to control the strength of RQ and MI, respectively.

To instantiate UniTok in LLM-based recommender systems, we first train UniTok on all items using Eq. (11). The trained UniTok then tokenizes each item into discrete semantic tokens, enabling LLMs to operate in the token space. Following prior work (Wang et al. 2024), user interaction histories $\mathbf{u}$ are converted into item token sequences $\mathbf{u} = [\tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2, \ldots, \tilde{\mathbf{c}}_P]$, and the LLM-based recommender system learns to predict the next interacted item token $\tilde{\mathbf{c}}_{P+1}$.

---

[2] Other methods, such as MINE (Belghazi et al. 2018) and InfoNCE (Oord, Li, and Vinyals 2018), can also estimate MI using neural networks, but require additional training. To maintain simplicity, we adopt HSIC, which is non-parametric.

[3] The kernel matrices are computed element-wise using the Gaussian kernel: $U_{ij} = u(\mathbf{x}_i^k, \mathbf{x}_j^k) = \exp\left(-\|\mathbf{x}_i^k - \mathbf{x}_j^k\|^2/2\sigma^2\right)$ and $V_{ij} = v(\mathbf{z}_i^k, \mathbf{z}_j^k) = \exp\left(-\|\mathbf{z}_i^k - \mathbf{z}_j^k\|^2/2\sigma^2\right)$ for the bandwidth $\sigma$.

## Theoretical Analyses

In this subsection, we establish the following theorems, which analyze and support the effectiveness of the key components within UniTok.

**Theorem 1.** *The token space induced by UniTok exhibits strictly higher entropy than that of standard codebook-based methods:*

$$H(\mathcal{C}_{\text{UniTok}}) > H(\mathcal{C}_{\text{standard}}), \quad (12)$$

*where $\mathcal{C}_{UniTok}$ and $\mathcal{C}_{standard}$ denote the discrete token distributions generated by UniTok and standard codebook-based methods, respectively.*

This indicates that judiciously incorporating multiple experts increases the overall entropy, thereby expanding the token space capacity.

**Theorem 2.** *Let $\mathbb{E}[\mathcal{L}_{UniTok}]$ and $\mathbb{E}[\mathcal{L}_{standard}]$ denote the expected quantization error of UniTok and standard codebook-based methods using a single shared codebook, respectively. Then, the following inequality holds:*

$$\mathbb{E}[\mathcal{L}_{\text{UniTok}}] \leq \mathbb{E}[\mathcal{L}_{\text{standard}}]. \quad (13)$$

This implies that a lower expected quantization error reflects more precise modeling of item tokenization. The TokenMoE framework further reduces this error by leveraging expert specialization, which effectively compensates for domain-specific inaccuracies to some extent, thereby improving quantization quality across diverse domains.

**Theorem 3.** *Suppose that the loss $\mathcal{L}^{(k)}$ on the $k$-th domain is Lipschitz-continuous with respect to the informativeness of representations. Then, the performance variability across domains is upper-bounded by the variance of MI:*

$$|\mathcal{L}^{(i)} - \mathcal{L}^{(j)}| \leq C\sqrt{\text{Var}\left[\widehat{I}^{(k)}\right]}, \forall i, j \quad (14)$$

*where $Var\left[\widehat{I}^{(k)}\right]$ is the variance of MI estimates across domains and $C$ is a constant.*

This implies that reducing MI variance across domains promotes consistent semantic representations, leading to more stable and reliable downstream performance.

We empirically validate each theorem's technical correctness and practical relevance through extensive experiments.

## Experimental Evaluation

In this section, we carry out comprehensive experiments to empirically validate the effectiveness of UniTok across multiple domains.

## Experimental Settings

**Datasets.** We conduct our experiments on ten real-world datasets spanning ten domains widely adopted for evaluating the performance of recommendations, which include Beauty, Cellphones, Grocery, Instruments, Office, Pet Supplies, Tools, Toys, Games[4], and Yelp[5]. Each item contains metadata, including title, category, and features. Due to page limitations, we report complete results across all datasets to assess overall performance, while employing three representative datasets for the efficiency and ablation analyses.

---

[4] https://nijianmo.github.io/amazon/index.html.

[5] https://www.yelp.com/dataset.

| Method | Beauty | Cellphones | Grocery | Instruments | Office | Pet Supplies | Tools | Toys | Games | Yelp |
|---|---|---|---|---|---|---|---|---|---|---|
| **MF** | 0.0369 | 0.0267 | 0.0216 | 0.0710 | 0.0255 | 0.0268 | 0.0169 | 0.0192 | 0.0366 | 0.0144 |
| **LightGCN** | 0.0285 | 0.0456 | 0.0357 | 0.0781 | 0.0301 | 0.0289 | 0.0257 | 0.0287 | 0.0417 | 0.0195 |
| **SASRec** | 0.0314 | 0.0446 | 0.0376 | 0.0609 | 0.0285 | 0.0301 | 0.0234 | 0.0239 | 0.0412 | 0.0183 |
| **Bert4Rec** | 0.0194 | 0.0268 | 0.0237 | 0.0573 | 0.0274 | 0.0161 | 0.0092 | 0.0177 | 0.0379 | 0.0131 |
| **P5-TID** | 0.0255 | 0.0357 | 0.0316 | 0.0721 | 0.0239 | 0.0243 | 0.0198 | 0.0202 | 0.0388 | 0.0154 |
| **P5-SemID** | 0.0304 | 0.0406 | 0.0351 | 0.0730 | 0.0283 | 0.0282 | 0.0237 | 0.0231 | 0.0432 | 0.0188 |
| **TIGER** | 0.0324 | 0.0446 | 0.0375 | 0.0788 | 0.0295 | 0.0279 | 0.0284 | 0.0268 | 0.0427 | 0.0208 |
| **LC-Rec** | <u>0.0381</u> | 0.0458 | 0.0369 | 0.0802 | 0.0311 | <u>0.0335</u> | <u>0.0307</u> | 0.0279 | 0.0451 | 0.0215 |
| **LETTER** | 0.0364 | <u>0.0473</u> | <u>0.0392</u> | <u>0.0831</u> | <u>0.0326</u> | 0.0307 | 0.0298 | <u>0.0291</u> | <u>0.0469</u> | <u>0.0231</u> |
| **UniTok** | **0.0478** | **0.0647** | **0.0533** | **0.0884** | **0.0432** | **0.0496** | **0.0439** | **0.0442** | **0.0476** | **0.0321** |
| **Gain** | 25.46% | 36.78% | 35.97% | 6.38% | 32.52% | 48.06% | 42.99% | 51.89% | 1.49% | 38.96% |

Table 1: Performance comparison among UniTok and recommendation competitors for the ten benchmark datasets in terms of NDCG@10. Here, the best and second-best performers are highlighted by bold and underline, respectively. The improvements are statistically significant on average ($p = 0.0219 < 0.05$) based on paired t-tests over five runs across all datasets.

**Competitors.** To comprehensively demonstrate the superiority of UniTok, we present nine benchmark recommendation methods, including four widely-used collaborative filtering methods, MF (Rendle et al. 2009), SASRec (Kang and McAuley 2018), LightGCN (He et al. 2020), Bert4Rec (Sun et al. 2019), and five item tokenization-aided recommendations, P5-TID (Hua et al. 2023), P5-SemID (Hua et al. 2023), TIGER (Rajput et al. 2023), LC-Rec (Zheng et al. 2024), and LETTER (Wang et al. 2024).

**Performance metrics.** Following the full-ranking protocol (He et al. 2020), we rank all non-interacted items for each user and evaluate using Recall@$M$ (R@$M$) and NDCG@$M$ (N@$M$), where $M \in \{5, 10\}$.

**Implementation details.** For item tokenization, we use 4-level codebook-based identifiers, where each codebook comprises 256 code vectors with a dimension of 32. We set $\lambda_{\text{RQ}}$ to 1 and $\lambda_{\text{MI}}$ to 0.03. All experiments are conducted on two NVIDIA RTX 3090 GPUs.

## Can One Tokenizer Serve All Domains?

To evaluate the recommendation accuracy of UniTok, we compare it with benchmark item tokenization methods across ten benchmark datasets from diverse domains. Notably, distinct from benchmark methods that train separate tokenizers for each dataset, UniTok trains a *single unified model* that handles all ten datasets jointly. As shown in Table 1, UniTok consistently outperforms competitors across all datasets, validating the effectiveness of our unified tokenization framework, achieving up to 51.89% improvement in terms of NDCG@10 on the Tools dataset. Unlike other approaches that require domain-specific customization, UniTok effectively captures item semantics across diverse domains through a single, shared tokenization process—highlighting the strength and generality of our proposed design.

We observe that item tokenization-aided LLM-based recommender systems, including TIGER, LC-Rec, LETTER, and UniTok, are superior to standard collaborative filtering methods such as MF, LightGCN, SASRec, and Bert4Rec.

| Module | Codebook-based methods | UniTok |
|---|---|---|
| **Codebook** | 0.33M | 0.36M |
| **Autoencoder** | 87.45M | 8.75M |
| **Router** | – | 0.01M |
| **Total** | 87.78M | 9.11M |

Table 2: Comparison of the number of trainable parameters. For competitors, we report the total number of trainable parameters accumulated across all ten datasets.

This improvement benefits from the rich semantic understanding and reasoning capabilities inherent in LLMs, which go beyond conventional user–item interactions. Moreover, integrating carefully designed item tokenization into LLM-based recommender systems yields additional performance gains compared to LLMs that solely on item metadata for tokenization (*e.g.*, P5-TID and P5-SemID). This is because item tokenization serves as a critical bridge between the item space and the language space, allowing the underlying model to represent items in a discrete, language-aligned semantic space that enhances generalization and reasoning.

## Is UniTok More Efficient than Traditional Tokenizers?

To validate the efficiency of UniTok, we compare the total number of trainable parameters between UniTok and traditional codebook-based competitors, including TIGER, LC-Rec, and LETTER, which share the same underlying architecture. While these competitors require training and storing separate tokenization models for each dataset, UniTok employs a single unified model shared across all datasets. For a fair comparison, we report the cumulative trainable parameter count of the codebook-based competitors over all datasets. As shown in Table 2, UniTok achieves approximately a tenfold reduction in the total number of trainable parameters. This efficiency comes primarily from using the shared autoencoder, while the number of additional trainable parameters introduced by the codebook and Token-MoE router remains negligible. By eliminating the need for

|  | Beauty | | Cellphones | | Grocery | |
|---|---|---|---|---|---|---|
| **Method** | **R@10** | **N@10** | **R@10** | **N@10** | **R@10** | **N@10** |
| **TIGER** | 0.0499 | 0.0267 | 0.0661 | 0.0342 | 0.0576 | 0.0273 |
| **LC-Rec** | <u>0.0564</u> | <u>0.0302</u> | 0.0647 | 0.0337 | 0.0584 | 0.0287 |
| **LETTER** | 0.0528 | 0.0288 | <u>0.0678</u> | <u>0.0363</u> | <u>0.0618</u> | <u>0.0315</u> |
| **UniTok** | **0.0934** | **0.0478** | **0.1251** | **0.0647** | **0.1061** | **0.0533** |
| **Gain** | 65.60% | 58.28% | 84.51% | 78.23% | 71.68% | 69.21% |

Table 3: Performance comparison of UniTok and tokenization competitors under a single unified training setup. Here, the best and second-best performers are highlighted by bold and underline, respectively.

|  | Clothing | | Health | | Sports | |
|---|---|---|---|---|---|---|
| **Method** | **R@10** | **N@10** | **R@10** | **N@10** | **R@10** | **N@10** |
| **TIGER** | 0.0501 | 0.0242 | 0.0677 | 0.0342 | 0.0469 | 0.0228 |
| **LC-Rec** | <u>0.0527</u> | <u>0.0266</u> | 0.0694 | 0.0358 | 0.0494 | 0.0246 |
| **LETTER** | 0.0515 | 0.0257 | <u>0.0717</u> | <u>0.0375</u> | <u>0.0510</u> | <u>0.0265</u> |
| **UniTok** | **0.0592** | **0.0288** | **0.0835** | **0.0442** | **0.0591** | **0.0298** |
| **Gain** | 12.33% | 8.27% | 16.46% | 17.87% | 15.88% | 12.45% |

Table 4: Performance comparison among UniTok and recommendation competitors for the three unseen datasets. Here, the best and second-best performers are highlighted by bold and underline, respectively.

domain-specific tokenization learning, UniTok offers substantial advantages in scalability and deployment efficiency.

We further evaluate the performance of competitors under a single unified training setup, where each competitor is trained jointly across ten datasets using a comparable number of trainable parameters to that of UniTok. As shown in Table 3, the competing methods exhibit substantial performance degradation in this setting compared to their single-domain scenario (see Table 1), primarily due to the difficulty of distinguishing items from different domains when using shared tokenization. In contrast, UniTok maintains consistently superior recommendation performance, achieving up to 84.51% improvement in Recall@10 on Cellphones, while using a *similar trainable parameter budget*. This efficiency stems from its modular TokenMoE architecture, where domain-specific experts learn semantics independently, while sharing a unified token space.

### Can UniTok be Generalized to Unseen Domains?

To evaluate the generalization ability of UniTok, we adopt a *zero-shot* setting where our item tokenizer model, UniTok, is trained once on the ten source datasets and is directly tested on unseen datasets from three target domains—Clothing, Health, and Sports—without any additional training or fine-tuning. This setup reflects real-world scenarios where new domains may appear dynamically after the model has been deployed.

As shown in Table 4, UniTok significantly outperforms existing item tokenization-based recommender systems, achieving up to 17.87% improvement in NDCG@10 on Health. While the competitors require retraining on each new dataset to achieve reasonable tokenization results, UniTok maintains robust accuracy without any further adaptation. This demonstrates our model's ability to learn a discrete token space that captures transferable item semantics across diverse domains. The consistently strong zero-shot performance further highlights the robustness and practical utility of our unified tokenization framework in supporting effective generalization across heterogeneous domains.

### What Makes UniTok Effective?

To assess the contribution of each component in UniTok, we perform an ablation study by progressively removing

|  | Beauty | | Cellphones | | Grocery | |
|---|---|---|---|---|---|---|
| **Method** | **R@10** | **N@10** | **R@10** | **N@10** | **R@10** | **N@10** |
| **UniTok-1** | 0.0558 | 0.0304 | 0.0702 | 0.0371 | 0.0633 | 0.0342 |
| **UniTok-2** | 0.0896 | 0.0436 | 0.1194 | 0.0606 | 0.0989 | 0.0497 |
| **UniTok-3** | 0.0915 | 0.0457 | 0.1225 | 0.0622 | 0.1044 | 0.0515 |
| **UniTok** | **0.0934** | **0.0478** | **0.1251** | **0.0647** | **0.1061** | **0.0533** |

Table 5: Ablation study results on the Beauty, Cellphones, and Grocery datasets.

or modifying its core modules: the TokenMoE module, the shared expert, and the MI calibration part. UniTok-1 removes the TokenMoE and MI calibration, using only a single set of codebooks without any MoE structure; UniTok-2 keeps TokenMoE, but removes the shared expert and MI calibration; and UniTok-3 only removes the MI calibration. UniTok includes all components. As shown in Table 8, removing any module leads to a noticeable drop in recommendation accuracy, which confirms that the combination of modules is crucial to UniTok's effectiveness. In particular, comparing UniTok-1 and UniTok-2 highlights the importance of the TokenMoE module, which significantly improves performance across all datasets by capturing domain-specific semantics. Additionally, in comparison with UniTok-3, introducing MI calibration further enhances performance by enforcing the learned latent embeddings to retain essential semantic information from each domain.

## Conclusions and Outlook

We explored an open yet fundamental challenge in multi-domain LLM-based recommendations by building a unified item tokenization framework. To this end, we proposed UniTok, which integrates a customized MoE architecture with codebooks to generate semantically meaningful tokens across diverse domains. Experiments on wide-ranging datasets showed that UniTok (a) achieves up to 51.89% gains in NDCG@10, (b) reduces trainable parameters by 9.63×, and (c) is theoretically validated with respect to the effectiveness of its individual components. Future work includes extending UniTok into a general-purpose tokenization interface for foundation models in recommendation.

## Acknowledgments

## References

Belghazi, M. I.; Baratin, A.; Rajeswar, S.; Ozair, S.; Bengio, Y.; Hjelm, R. D.; and Courville, A. C. 2018. Mutual Information Neural Estimation. In *ICML, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, 530–539.

Dai, D.; Deng, C.; Zhao, C.; Xu, R.; Gao, H.; Chen, D.; Li, J.; Zeng, W.; Yu, X.; Wu, Y.; et al. 2024. DeepSeekMoE: Towards Ultimate Expert Specialization in Mixture-of-Experts Language Models. *arXiv preprint arXiv:2401.06066*.

Fedus, W.; Zoph, B.; and Shazeer, N. 2022. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *Journal of Machine Learning Research*, 23(120): 1–39.

Gretton, A.; Bousquet, O.; Smola, A.; and Schölkopf, B. 2005. Measuring Statistical Dependence with Hilbert-Schmidt Norms. In *ALT, Singapore, October 8-11, 2005*, 63–77.

Gururangan, S.; Marasovic, A.; Swayamdipta, S.; Lo, K.; Beltagy, I.; Downey, D.; and Smith, N. A. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *ACL, Virtual Event, July 5-10, 2020*, 8342–8360.

He, X.; Deng, K.; Wang, X.; Li, Y.; Zhang, Y.; and Wang, M. 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation. In *SIGIR, Virtual Event, July 25-30, 2020*, 639–648.

Hua, W.; Xu, S.; Ge, Y.; and Zhang, Y. 2023. How to Index Item IDs for Recommendation Foundation Models. In *SIGIR-AP, Beijing, China, November 26-28, 2023*, 195–204.

Jacobs, R. A.; Jordan, M. I.; Nowlan, S. J.; and Hinton, G. E. 1991. Adaptive Mixtures of Local Experts. *Neural Computation*, 3(1): 79–87.

Jain, Y.; Behl, H. S.; Kira, Z.; and Vineet, V. 2023. DAMEX: Dataset-aware Mixture-of-Experts for visual understanding of mixture-of-datasets. In *NeurIPS, New Orleans, LA, USA, December 10 - 16, 2023*.

Jiang, Y.; Li, Q.; Zhu, H.; Yu, J.; Li, J.; Xu, Z.; Dong, H.; and Zheng, B. 2022. Adaptive Domain Interest Network for Multi-Domain Recommendation. In *CIKM, Atlanta, GA, USA, October 17-21, 2022*, 3212–3221.

Kang, W.-C.; and McAuley, J. 2018. Self-Attentive Sequential Recommendation. In *ICDM, Singapore, November 17-20, 2018*, 197–206.

Lepikhin, D.; Lee, H.; Xu, Y.; Chen, D.; Firat, O.; Huang, Y.; Krikun, M.; Shazeer, N.; and Chen, Z. 2021. GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding. In *ICLR, Virtual Event, Austria, May 3-7, 2021*.

Li, Y.; Pogodin, R.; Sutherland, D. J.; and Gretton, A. 2021. Self-Supervised Learning with Kernel Dependence Maximization. In *NeurIPS, December 6-14, 2021, Virtual Event*, 15543–15556.

Ma, R.; Tan, Y.; Zhou, X.; Chen, X.; Liang, D.; Wang, S.; Wu, W.; and Gui, T. 2022. Searching for Optimal Subword Tokenization in Cross-domain NER. In *IJCAI, Vienna, Austria, 23-29 July 2022*, 4289–4295.

Ning, W.; Yan, X.; Liu, W.; Cheng, R.; Zhang, R.; and Tang, B. 2023. Multi-Domain Recommendation with Embedding Disentangling and Domain Alignment. In *CIKM, Birmingham, United Kingdom, October 21-25, 2023*, 1917–1927.

Oord, A. v. d.; Li, Y.; and Vinyals, O. 2018. Representation Learning with Contrastive Predictive Coding. *arXiv preprint arXiv:1807.03748*.

Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 21(140): 1–67.

Rajput, S.; Mehta, N.; Singh, A.; Keshavan, R. H.; Vu, T.; Heldt, L.; Hong, L.; Tay, Y.; Tran, V. Q.; Samost, J.; Kula, M.; Chi, E. H.; and Sathiamoorthy, M. 2023. Recommender Systems with Generative Retrieval. In *NeurIPS, New Orleans, LA, USA, December 10 - 16, 2023*.

Rendle, S.; Freudenthaler, C.; Gantner, Z.; and Schmidt-Thieme, L. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *UAI, Montreal, QC, Canada, June 18-21, 2009*, 452–461.

Shazeer, N.; Mirhoseini, A.; Maziarz, K.; Davis, A.; Le, Q. V.; Hinton, G. E.; and Dean, J. 2017. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer. In *ICLR, Toulon, France, April 24-26, 2017*.

Sun, F.; Liu, J.; Wu, J.; Pei, C.; Lin, X.; Ou, W.; and Jiang, P. 2019. BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer. In *CIKM, Beijing, China, November 3-7, 2019*, 1441–1450.

Ullah, I.; Carrión-Ojeda, D.; Escalera, S.; Guyon, I.; Huisman, M.; Mohr, F.; van Rijn, J. N.; Sun, H.; Vanschoren, J.; and Vu, P. A. 2022. Meta-Album: Multi-domain Meta-Dataset for Few-Shot Image Classification. In *NeurIPS, New Orleans, LA, USA, November 28 - December 9, 2022*.

Wang, W.; Bao, H.; Lin, X.; Zhang, J.; Li, Y.; Feng, F.; Ng, S.-K.; and Chua, T.-S. 2024. Learnable Item Tokenization for Generative Recommendation. In *CIKM, Boise, ID, USA, October 21-25, 2024*, 2400–2409.

Zhang, Y.; DING, H.; Shui, Z.; Ma, Y.; Zou, J.; Deoras, A.; and Wang, H. 2021. Language Models as Recommender Systems: Evaluations and Limitations. In *I (Still) Can't Believe It's Not Better! NeurIPS 2021 Workshop*.

Zheng, B.; Hou, Y.; Lu, H.; Chen, Y.; Zhao, W. X.; Chen, M.; and Wen, J.-R. 2024. Adapting Large Language Models by Integrating Collaborative Semantics for Recommendation. In *ICDE, Utrecht, The Netherlands, May 13-16, 2024*, 1435–1448.

# *Appendix* for "Tokenize Once, Recommend Anywhere: Unified Item Tokenization for Multi-domain LLM-based Recommendation"

## List of contents in Appendix

- Summary of **notations** and their meanings used in the main manuscript

- **Related work** on LLM for generative recommendation and item tokenization: **Appendix A**

- Detailed **algorithm** of training procedure for UniTok: **Appendix B**

- Formal **proofs** of the theorems in the main manuscript: **Appendix C**

- **Additional experimental evaluations**:

  - ✓ Description and statistics of all **datasets** used: **Appendix D.1**
  - ✓ Overview of **competitors** for comparison: **Appendix D.2**
  - ✓ Formal definitions of **evaluation metrics**, including Recall@$M$ and NDCG@$M$: **Appendix D.3**
  - ✓ **Implementation details** including hyperparameters and training setup: **Appendix D.4**
  - ✓ **Full performance results** across all evaluation metrics: **Appendix D.5**
  - ✓ **Ablation study** for evaluation of contributions from individual components: **Appendix D.6**
  - ✓ **Sensitivity analysis** on robustness under different hyperparameter settings: **Appendix D.7**
  - ✓ **Empirical validation for theorems**: **Appendix D.8**

## Table of notations

| Notation | Description |
|---|---|
| $\mathcal{D}_k$ | The $k$-th dataset |
| $K$ | The number of datasets |
| $\mathcal{I}_k$ | Item set of $\mathcal{D}_k$ |
| $\mathbf{X}^k$ | Semantic embeddings for domain $k$ |
| $d$ | Embedding dimensionality |
| $\mathbf{x}_i^k$ | Semantic embedding of the $i$-th item in $\mathcal{D}_k$ |
| $f_\theta$ | Shared encoder |
| $g_\phi$ | Shared decoder |
| $\mathbf{Z}^k$ | Hidden embeddings of $\mathbf{X}^k$ |
| $\mathbf{z}_i^k$ | Hidden embedding of $\mathbf{x}_i^k$ |
| $G(\cdot)$ | Router function |
| $G_k$ | Routing probability assigned to the $k$-th expert |
| $N$ | Number of selected top experts |
| $E_k$ | The $k$-th expert module |
| $E_{\text{share}}$ | The shared expert module |
| $L$ | Codebook size |
| $T$ | Number of code vector in each codebook |
| $C_\ell$ | The $\ell$-th codebook |
| $\mathbf{c}_i^k$ | Discrete token of $\mathbf{x}_i^k$ |
| $\widehat{\text{HSIC}}(\mathbf{X}^k, \mathbf{Z}^k)$ | Empirical estimate of HSIC between $\mathbf{X}^k$ and $\mathbf{Z}^k$ |
| $\mathbf{U}, \mathbf{V}$ | Gaussian kernel matrices |
| $\widehat{I}^{(k)}$ | $\widehat{\text{HSIC}}(\mathbf{X}^k, \mathbf{Z}^k)$ |
| $\text{Var}\left[\widehat{I}^{(k)}\right]$ | Variance of MI estimates across domains |
| $\mathbb{E}\left[\widehat{I}^{(k)}\right]$ | Expectation of MI estimates across domains |
| $\mathcal{L}_{\text{Total}}$ | Total loss |
| $\mathcal{L}_{\text{Rec}}$ | Reconstruction loss |
| $\mathcal{L}_{\text{RQ}}$ | Resdual quantization loss |
| $\mathcal{L}_{\text{MI}}$ | Mutual information calibration loss |

# A. Related Work

Our proposed method is related to two broader areas of research, namely LLM for generative recommendation and item tokenization.

## 1. LLM for Generative Recommendation

LLMs have been widely adopted in recommender systems for their powerful generative capabilities (Li et al. 2024; Gong et al. 2023; Lin et al. 2024d). They have been utilized not only to enhance traditional recommendation methods but also to directly generate recommendations. A growing body of studies leveraged LLMs for data augmentation by synthesizing user—item interactions (Wei et al. 2024), incorporating external knowledge (Xi et al. 2024), and enriching user and item representations at the representation level (Liu et al. 2024; Qiu et al. 2021; Ren et al. 2024; Wu et al. 2022). LLM-based recommender systems employed a projection layer that aligns user and item embeddings within a shared latent space, enabling the computation of similarity scores for ranking tasks (Li et al. 2023a; Lin et al. 2024a; Zhang et al. 2025). Other methods focused on click-through rate prediction by feeding user profiles and item descriptions into LLMs to model their interactions, predicting the likelihood of user engagement (Bao et al. 2023; Lin et al. 2024b; Prakash et al. 2023). However, since LLMs originated from language processing, a gap exists between the language space and the item space when applying them to recommendation tasks.

## 2. Item Tokenization

In recommender systems powered by LLMs, a fundamental challenge arises from the gap between the language space and the item space. To address this issue, recent research has explored item tokenization strategies that effectively bridge this gap. Existing approaches can be broadly categorized into three types based on the nature of item identifiers. First, ID-based identifiers represent items as discrete tokens using unique IDs (Geng et al. 2022; Hua et al. 2023; Wang et al. 2024). While simple, these methods lack semantic information and generalization capability. Second, textual identifiers leverage item metadata, such as titles or descriptions, enabling LLMs to utilize pretrained linguistic knowledge for improved semantic understanding and generalization (Bao et al. 2025; Dai et al. 2023; Li et al. 2023b; Liao et al. 2023; Zhang et al. 2023, 2021; Lin et al. 2024c). Lastly, codebook-based identifiers employ learned discrete representations through vector quantization (Rajput et al. 2023; Wang et al. 2024; Zheng et al. 2024; Wang et al. 2024; Zheng et al. 2024; Zhai et al. 2025). These methods combine the interpretability of discrete tokens with the semantic richness of embeddings, providing structured and generalizable representations of items. However, common limitations of these item tokenization techniques hinder the generalization ability of LLMs. Specifically, most approaches rely on training separate models for each dataset (Rajput et al. 2023; Wang et al. 2024; Zheng et al. 2024; Wang et al. 2024; Zheng et al. 2024), while a more recent approach requires additional costly pre-training to transfer knowledge across domains for adaptation (Zhai et al. 2025).

# B. Algorithm of UniTok

Algorithm 1 outlines the training procedure of the proposed UniTok framework, which unifies item tokenization across multiple domains. Initially (Line 1), all core components, including the encoder, decoder, router, expert modules, and codebooks, are initialized. Each item $i_i^k$ from domain $\mathcal{D}_k$ is embedded into a semantic embedding $\mathbf{x}_i^k$ using a pretrained embedding model (Lines 2–6). During training (Lines 7–13), mini-batches are sampled from mixed domains to encourage multi-domain generalization. Each input $\mathbf{x}_i^k$ is first encoded into $\mathbf{z}_i^k$ (Line 10), then passed through the TokenMoE module to obtain a quantized embedding $\hat{\mathbf{z}}_i^k$ and discrete token $\mathbf{c}_i^k$ via expert routing and residual quantization (Line 11), and finally reconstructed $\hat{\mathbf{x}}_i^k$ via the decoder (Line 12). The total training objective combines a reconstruction loss (Line 14), a residual quantization loss (Line 15), and a mutual information calibration loss (Line 16) to enforce informativeness and semantic consistency. All parameters are updated jointly (Line 17), and the final output consists of discrete tokens for all items (Line 19).

# C. Theorems and Proofs

**Theorem 1.** *The token space induced by UniTok exhibits strictly higher entropy than that of standard codebook-based methods:*

$$H(\mathcal{C}_{\mathsf{UniTok}}) > H(\mathcal{C}_{\mathrm{standard}}),$$

*where $\mathcal{C}_{UniTok}$ and $\mathcal{C}_{standard}$ denote the discrete token distributions generated by UniTok and standard codebook-based methods, respectively.*

**Proof.** Consider the latent space representation of both models:

**Entropy of standard codebook-based methods:** In standard codebook-based methods, a single codebook with $T$ code vectors is used at each of the $L$ levels. The total number of unique tokens in the token space is

$$|\mathcal{C}_{\mathrm{standard}}| = T^L. \tag{15}$$

Assuming a uniform probability distribution over the quantized representations:

$$P(c) = \frac{1}{T^L}, \tag{16}$$

---

Algorithm 1: **UniTok** Training Procedure

---

**Input**: Multi-domain item datasets $\{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_K\}$ with text metadata *e.g.*, title, category, description)
**Parameters**: Encoder $f_\theta$, decoder $g_\phi$, router $h(\cdot)$, expert modules $\{E_1, ..., E_K, E_{\text{share}}\}$ with associated codebooks
**Output**: Discrete tokens $\mathbf{c}_i^k \in \mathcal{C}$ for all items

 1: Initialize all modules: $f_\theta$, $g_\phi$, router $h(\cdot)$, experts $E_k$ (including $E_{\text{share}}$), and their codebooks
 2: **for** $k = 1$ to $K$ **do**
 3:    **for** item $i_i^k \in \mathcal{D}_k$ **do**
 4:       $\mathbf{x}_i^k \leftarrow \text{SemanticEmbedding}(i_i^k)$
 5:    **end for**
 6: **end for**
 7: **for** epoch $= 1$ to $E$ **do**
 8:    Sample mini-batch $\mathcal{B} = \{\mathbf{x}_i^k\}$ from mixed domains
 9:    **for** each item $\mathbf{x}_i^k$ in batch $\mathcal{B}$ **do**
10:       $\mathbf{z}_i^k \leftarrow f_\theta(\mathbf{x}_i^k)$                             // Encode input
11:       $\hat{\mathbf{z}}_i^k, \mathbf{c}_i^k \leftarrow \text{TokenMoE}(\mathbf{z}_i^k)$               // Discretize via TokenMoE
12:       $\hat{\mathbf{x}}_i^k \leftarrow g_\phi(\hat{\mathbf{z}}_i^k)$                         // Decode reconstruction
13:    **end for**
14:    Compute reconstruction loss $\mathcal{L}_{\text{Rec}}$ in the main manuscript of Eq. (3)
15:    Compute RQ loss $\mathcal{L}_{\text{RQ}}$ in the main manuscript of Eq. (8)
16:    Compute MI loss $\mathcal{L}_{\text{MI}}$ in the main manuscript of Eq. (10)
17:    Update $f_\theta$, $g_\phi$, $h(\cdot)$, and expert parameters
18: **end for**
19: **return** token assignments $\mathbf{c}_i^k$ for all items

---

The entropy of the token space in standard codebook-based methods is

$$H(\mathcal{C}_{\text{standard}}) = -\sum_{c \in \mathcal{C}} P(c) \log P(c). \tag{17}$$

Substituting $P(c)$, we have

$$H(\mathcal{C}_{\text{standard}}) = L \log T. \tag{18}$$

**Entropy of UniTok:** In our UniTok, there are $K$ domain-specific experts, each with its own independent codebook of size $T^L$. Given an item embedding $\mathbf{z}$, the router function assigns a probability distribution over experts:

$$G_k = \frac{\exp(h(\mathbf{z}^{(k)}))}{\sum_{j=1}^{K} \exp(h(\mathbf{z}^{(j)}))}, \tag{19}$$

where the $h(\cdot)$ is the linear transformation and $\mathbf{z}^{(k)}$ denotes the logit corresponding to the $k$-th domain-specific expert.

Let $\mathcal{G}$ denote the random variable representing expert selection. An item token $c$ in $\mathcal{C}_{\text{UniTok}}$ is then generated by

- Sampling an expert index $k \sim \mathcal{G}$,
- Sampling a token $c \sim \mathcal{C}_k$.

Here, $\mathcal{C}_k$ is the discrete token distributions in the $k$-th expert.

By the chain rule of entropy, it follows that

$$H(\mathcal{C}_{\text{UniTok}}) = H(\mathcal{G}) + H(\mathcal{C}_k \mid \mathcal{G}). \tag{20}$$

The conditional entropy term $H(\mathcal{C}_k \mid \mathcal{G})$ is the expectation of the entropy of each expert's latent space, weighted by the distribution from the router, and is written as

$$H(\mathcal{C}_k \mid \mathcal{G}) = \sum_{k=1}^{K} G_k \cdot H(\mathcal{C}_k) = \mathbb{E}_{k \sim \mathcal{G}}[H(\mathcal{C}_k)]. \tag{21}$$

Then, we have

$$H(\mathcal{C}_{\text{UniTok}}) = H(\mathcal{G}) + \mathbb{E}_{k \sim \mathcal{G}}[H(\mathcal{C}_k)]. \tag{22}$$

Since each expert has the same quantization entropy as standard codebook-based methods, it follows that

$$H(\mathcal{C}_k) = L \log T \quad \text{for all } k, \tag{23}$$

We then obtain

$$\mathbb{E}_{k \sim \mathcal{G}}[H(\mathcal{C}_k)] = L \log T. \tag{24}$$

The router entropy $H(\mathcal{G})$ is

$$H(\mathcal{G}) = -\sum_{k=1}^{K} G_k \log G_k, \tag{25}$$

which is always positive when $K > 1$, i.e., $H(\mathcal{G}) > 0$.

Thus, the entropy of $\mathcal{C}_{\mathsf{UniTok}}$ is

$$H(\mathcal{C}_{\mathsf{UniTok}}) = H(\mathcal{G}) + L \log T > H(\mathcal{C}_{\mathrm{standard}}). \tag{26}$$

This completes the proof of Theorem 1. ∎

Before proving Theorem 2, we begin with Lemma 1 below.

**Lemma 1.** *The function $f(\mathbf{x}) = \|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^n$, is convex.*

**Proof.** To show that $f(\mathbf{x})$ is convex, we check whether its Hessian matrix is positive semidefinite.

**Step 1: Compute the gradient of $f(\mathbf{x})$ as follows:**

$$\nabla f(\mathbf{x}) = \nabla(\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}. \tag{27}$$

**Step 2: Compute the Hessian of $f(\mathbf{x})$ as follows:**

$$\nabla^2 f(\mathbf{x}) = \nabla(2\mathbf{x}) = 2\mathbf{I}, \tag{28}$$

where $\mathbf{I}$ is the $n \times n$ identity matrix.

Therefore, the Hessian is positive definite, which implies that $f(\mathbf{x})$ is a strictly convex function. This completes the proof of Lemma 1. ∎

**Theorem 2.** *Let $\mathbb{E}[\mathcal{L}_{UniTok}]$ and $\mathbb{E}[\mathcal{L}_{standard}]$ denote the expected quantization error of UniTok and standard codebook-based methods using a single shared codebook, respectively. Then, the following inequality holds:*

$$\mathbb{E}[\mathcal{L}_{\mathsf{UniTok}}] \leq \mathbb{E}[\mathcal{L}_{\mathrm{standard}}].$$

**Proof.** Let $\mathbf{z} \in \mathbb{R}^d$ be an item latent embedding drawn from a distribution $p(z)$. Let $\hat{\mathbf{z}} = E(\mathbf{z})$ denote the output of the quantization function, which generates a combination of code vectors to approximate the item latent embedding $\mathbf{z}$.

For standard codebook-based methods using a single set of codebooks, the expected quantization error is given by

$$\mathbb{E}[\mathcal{L}_{\mathrm{standard}}] = \mathbb{E}_{\mathbf{z} \sim p(z)} \left[ \|\mathbf{z} - E(\mathbf{z})\|^2 \right]. \tag{29}$$

Now, consider UniTok with $K$ domain-specific experts, each with its own set of codebooks. Given the embedding $\mathbf{z}$, the router function that assigns probabilities to each expert:

$$G_k = \frac{\exp(h(\mathbf{z}^{(k)}))}{\sum_{j=1}^{K} \exp(h(\mathbf{z}^{(j)}))}. \tag{30}$$

The approximated latent embedding is generated by $\hat{\mathbf{z}} = \sum_{k=1}^{K} G_k \cdot E_k(\mathbf{z})$. Thus, the quantization error is

$$\mathcal{L}_{\mathsf{UniTok}}(\mathbf{z}) = \left\| \mathbf{z} - \sum_{k=1}^{K} G_k \cdot E_k(\mathbf{z}) \right\|^2. \tag{31}$$

where $E_k(\mathbf{z})$ is the quantization function of the $k$-th expert module. As $\sum_{k=1}^{K} G_k = 1$, we have

$$\mathcal{L}_{\mathsf{UniTok}}(\mathbf{z}) = \left\| \sum_{k=1}^{K} G_k \cdot (\mathbf{z} - E_k(\mathbf{z})) \right\|^2. \tag{32}$$

We define the approximation error for each $z$ as

$$\epsilon_k(\mathbf{z}) = \mathbf{z} - E_k(\mathbf{z}). \tag{33}$$

Replacing with $\epsilon_k(\mathbf{z})$, we have

$$\mathcal{L}_{\mathsf{UniTok}}(\mathbf{z}) = \left\| \sum_{k=1}^{K} G_k \cdot \epsilon_k(\mathbf{z}) \right\|^2. \tag{34}$$

Using Jensen's inequality, we apply the convexity of the squared norm function based on Lemma 1:

$$\mathcal{L}_{\text{UniTok}}(\mathbf{z}) = \left\| \sum_{k=1}^{K} G_k \cdot \epsilon_k(\mathbf{z}) \right\|^2 \leq \sum_{k=1}^{K} G_k \cdot \|\epsilon_k(\mathbf{z})\|^2. \tag{35}$$

Taking the expectation over the data distribution $p(z)$ on both sides of Eq. (21), we obtain

$$\mathbb{E}_{\mathbf{z} \sim p(z)} \left[ \left\| \sum_{k=1}^{K} G_k \cdot \epsilon_k(\mathbf{z}) \right\|^2 \right] \leq \mathbb{E}_{\mathbf{z} \sim p(z)} \left[ \sum_{k=1}^{K} G_k \cdot \|\epsilon_k(\mathbf{z})\|^2 \right]. \tag{36}$$

Since the single quantization function $E(\mathbf{z})$ can be viewed as a degenerate case where $G_k = 1$ for a single expert and $0$ for all others, the UniTok always provides a better or equal approximation:

$$\mathbb{E}_{\mathbf{z} \sim p(z)} \left[ \mathcal{L}_{\text{UniTok}}(\mathbf{z}) \right] \leq \mathbb{E}_{\mathbf{z} \sim p(z)} [\mathcal{L}_{\text{standard}}(\mathbf{z})]. \tag{37}$$

This completes the proof of Theorem 2. ∎

**Theorem 3.** *Suppose that the loss $\mathcal{L}^{(k)}$ on the $k$-th domain is Lipschitz-continuous with respect to the informativeness of representations. Then, the performance variability across domains is upper-bounded by the variance of MI:*

$$|\mathcal{L}^{(i)} - \mathcal{L}^{(j)}| \leq C \sqrt{\text{Var}\left[\widehat{I}^{(k)}\right]}, \forall i, j$$

*where $\text{Var}\left[\widehat{I}^{(k)}\right]$ is the variance of MI estimates across domains and $C$ is a constant.*

**Proof.** Let $\widehat{I}^{(k)} = \widehat{\text{HSIC}}(\mathbf{X}^k, \mathbf{Z}^k)$ be the MI between the input $\mathbf{X}^k$ and its representation $\mathbf{Z}^k$ for the $k$-th domain. The loss $\mathcal{L}^{(k)}$ on the $k$-th domain is Lipschitz-continuous with respect to the informativeness of representation, *i.e.*,

$$\left| \mathcal{L}^{(i)} - \mathcal{L}^{(j)} \right| \leq L \left| \widehat{\text{HSIC}}(\mathbf{X}^i; \mathbf{Z}^i) - \widehat{\text{HSIC}}\left(\mathbf{X}^j; \mathbf{Z}^j\right) \right|, \tag{38}$$

where $L$ is the Lipschitz constant that quantifies the sensitivity of the downstream loss to changes in representation informativeness (i.e., MI).

This assumption ensures that variations in representation informativeness (measured by MI) translate into bounded variations in the downstream task loss, effectively connecting the input–output MI stability (**?**). As a result, controlling the variance of MI across domains allows us to stabilize performance fluctuations.

Now, we take the maximum gap in MI across all domain pairs and relate it to the variance of MI:

$$
\begin{aligned}
& \left| \widehat{I}^{(i)} - \widehat{I}^{(j)} \right| \\
\leq & \max_{i,j} \left| \widehat{I}^{(i)} - \widehat{I}^{(j)} \right| \\
= & \max_{i,j} \left| \widehat{I}^{(i)} - \mu - (\widehat{I}^{(j)} - \mu) \right| \\
\leq & \max_{i,j} \left( \left| \widehat{I}^{(i)} - \mu \right| + \left| \left( \widehat{I}^{(j)} - \mu \right) \right| \right) \\
\leq & 2 \max_{k} \left| \widehat{I}^{(k)} - \mu \right| \\
\leq & 2 \sqrt{\sum_{k=1}^{K} \left( \widehat{I}^{(k)} - \mu \right)^2} \\
\leq & 2 \sqrt{K Var\left[ \widehat{I}^{(k)} \right]}.
\end{aligned}
\tag{39}
$$

Therefore, we obtain

$$\left| \mathcal{L}^{(i)} - \mathcal{L}^{(j)} \right| \leq L \left| \widehat{I}^{(i)} - \widehat{I}^{(j)} \right| \leq C \sqrt{Var\left[ \widehat{I}^{(k)} \right]},$$

where $C = 2L\sqrt{K}$.

This completes the proof of Theorem 3. ∎

## D. Additional Experimental Evaluations

### 1. Datasets

We conduct our experiments on ten real-world datasets widely adopted for evaluating the performance of recommendations, which include Beauty, Cellphones, Grocery, Instruments, Office, Pet Supplies, Tools, Toys, Games[6], and Yelp[7]. Each item in

---

[6]https://nijianmo.github.io/amazon/index.html.

[7]https://www.yelp.com/dataset.

these datasets contains rich content information, including title, category, and description. Additionally, we use three datasets, Clothing, Health, and Sports[8], for zero-shot evaluation. Table 6 provides a summary of the statistics for all datasets.

| Dataset | # of users | # of items | # of interactions |
|---|---|---|---|
| **Beauty** | 22,363 | 12,101 | 198,502 |
| **Cellphones** | 27,879 | 10,429 | 194,439 |
| **Grocery** | 14,681 | 8,713 | 151,254 |
| **Instruments** | 24,772 | 9,922 | 206,153 |
| **Office Products** | 4,905 | 2,420 | 53,258 |
| **Pet Supplies** | 19,856 | 8,510 | 157,836 |
| **Tools** | 16,638 | 10,217 | 134,476 |
| **Toys** | 19,412 | 11,924 | 167,597 |
| **Games** | 24,303 | 10,672 | 231,780 |
| **Yelp** | 30,431 | 20,033 | 316,354 |
| **Clothing** | 39,387 | 23,033 | 278,677 |
| **Health** | 38609 | 18533 | 346,355 |
| **Sports** | 35,598 | 18,357 | 296,337 |

Table 6: The statistics of datasets from ten domains suited to recommendations.

## 2. Competitors

To comprehensively demonstrate the superiority of UniTok, we present nine benchmark recommendation methods, including four widely-used collaborative filtering methods, MF (Rendle et al. 2009), SASRec (Kang and McAuley 2018), LightGCN (He et al. 2020), Bert4Rec (Sun et al. 2019), and five item tokenization-aided recommendations, P5-TID (Hua et al. 2023), P5-SemID (Hua et al. 2023), TIGER (Rajput et al. 2023), LC-Rec (Zheng et al. 2024), and LETTER (Wang et al. 2024). We implemented all these methods using the parameter settings described in their original articles.

- **MF (Rendle et al. 2009)**. This models user-item interactions by decomposing them into latent user and item embeddings through matrix factorization.
- **SASRec (Kang and McAuley 2018)**. This method uses self-attention to model user sequences, effectively capturing both short- and long-term dependencies with high efficiency.
- **LightGCN (He et al. 2020)**. This model simplifies graph convolutional networks (GCNs) for recommendation by focusing solely on neighborhood aggregation, eliminating feature transformation and nonlinear activation to enhance performance.
- **Bert4Rec (Sun et al. 2019)**. This method leverages bidirectional self-attention with a Cloze-style objective to model user behavior sequences, enabling richer context modeling than unidirectional sequential methods.
- **P5-TID (Hua et al. 2023)**. This method uses item titles as textual identifiers to enable LLM-based generative recommendation.
- **P5-SemID (Hua et al. 2023)**. This method constructs item identifiers from metadata (e.g., attributes) for LLM-based generative recommendation.
- **TIGER (Rajput et al. 2023)**. This method trains a sequence-to-sequence model to autoregressively generate item identifiers composed of semantic codeword tuples for next-item prediction.
- **LC-Rec (Zheng et al. 2024)**. This method integrates language and collaborative semantics by learning meaningful item indices via vector quantization and tuning LLMs through alignment tasks for direct item generation.
- **LETTER (Wang et al. 2024)**. This method is a learnable tokenizer for LLM-based generative recommendation that combines residual quantization, contrastive alignment, and diversity loss to encode hierarchical semantics, capture collaborative signals, and mitigate code assignment bias.

## 3. Performance Metrics

We follow the **full-ranking evaluation protocol** (He et al. 2020), where all non-interacted items are ranked for each user. To evaluate performance, we adopt two widely used ranking metrics: **Recall@$M$ (R@$M$)** and **NDCG@$M$ (N@$M$)**, where $M \in \{5, 10\}$. We formally define the ranking metrics used in our experiments.

**Recall@$M$.** Recall@$M$ measures the proportion of relevant items successfully retrieved in the top-$M$ recommendations:

$$\text{Recall}@M = \frac{|\text{Top-}M \text{ recommended items} \cap \text{Ground truth ttems}|}{|\text{Ground truth items}|}.$$

---

[8]https://nijianmo.github.io/amazon/index.html.

**NDCG@$M$ (Normalized Discounted Cumulative Gain).** NDCG@$M$ considers both the relevance and the ranking positions of the recommended items:

$$\text{NDCG@}M = \frac{1}{\text{IDCG@}M} \sum_{i=1}^{M} \frac{\mathbb{1}\{\text{item}_i \in \text{Ground truth}\}}{\log_2(i+1)},$$

where IDCG@$M$ denotes the ideal DCG, i.e., the maximum possible DCG for the given ground truth, ensuring normalization between 0 and 1.

## 4. Implementation Details

To evaluate recommendation performance, we adopt TIGER (Rajput et al. 2023), a representative LLM-based generative recommender system that uses T5 (Raffel et al. 2020) as its backbone. We follow this design and retain T5 in our experiments due to its unified text-to-text framework, strong performance across a wide range of generative tasks, and efficiency in conditional sequence generation—making it well-suited to the generative recommendation paradigm adopted by TIGER.

| Method | Beauty | | Cellphones | | Grocery | | Insturments | | Office Products | |
|---|---|---|---|---|---|---|---|---|---|---|
| | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 |
| MF | 0.0614 | 0.0369 | 0.0604 | 0.0267 | 0.0418 | 0.0216 | 0.0930 | 0.0710 | 0.0569 | 0.0255 |
| LightGCN | 0.0639 | 0.0285 | 0.0668 | 0.0456 | 0.0697 | 0.0357 | 0.1008 | 0.0781 | 0.0587 | 0.0301 |
| SASRec | 0.0646 | 0.0314 | 0.0651 | 0.0446 | 0.0701 | 0.0376 | 0.0905 | 0.0609 | 0.0574 | 0.0285 |
| Bert4Rec | 0.0372 | 0.0194 | 0.0507 | 0.0268 | 0.0448 | 0.0237 | 0.0791 | 0.0573 | 0.0563 | 0.0274 |
| P5-TID | 0.0532 | 0.0255 | 0.0648 | 0.0357 | 0.0617 | 0.0316 | 0.0928 | 0.0721 | 0.0557 | 0.0239 |
| P5-SemID | 0.0584 | 0.0304 | 0.0737 | 0.0406 | 0.0641 | 0.0351 | 0.0964 | 0.0730 | 0.0592 | 0.0283 |
| TIGER | 0.0624 | 0.0324 | 0.0838 | 0.0446 | 0.0706 | 0.0375 | 0.1047 | 0.0788 | 0.0594 | 0.0295 |
| LC-Rec | <u>0.0684</u> | <u>0.0381</u> | 0.0859 | 0.0458 | 0.0722 | 0.0369 | 0.1066 | 0.0802 | 0.0637 | 0.0311 |
| LETTER | 0.0672 | 0.0364 | <u>0.0876</u> | <u>0.0473</u> | <u>0.0731</u> | <u>0.0392</u> | <u>0.1122</u> | <u>0.0831</u> | <u>0.0649</u> | <u>0.0326</u> |
| **UniTok** | **0.0934** | **0.0478** | **0.1251** | **0.0647** | **0.1061** | **0.0533** | **0.1361** | **0.0884** | **0.0897** | **0.0432** |
| Improve | 36.55% | 25.46% | 42.81% | 36.78% | 45.14% | 35.97% | 21.30% | 6.38% | 38.21% | 32.52% |

| Method | Pet Supplies | | Tools | | Toys | | Games | | Yelp | |
|---|---|---|---|---|---|---|---|---|---|---|
| | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 |
| MF | 0.0503 | 0.0268 | 0.0356 | 0.0169 | 0.0405 | 0.0192 | 0.0359 | 0.0366 | 0.0304 | 0.0144 |
| LightGCN | 0.0529 | 0.0289 | 0.0482 | 0.0257 | 0.0495 | 0.0287 | 0.0407 | 0.0417 | 0.0368 | 0.0195 |
| SASRec | 0.0538 | 0.0301 | 0.0475 | 0.0234 | 0.0473 | 0.0239 | 0.0401 | 0.0412 | 0.0354 | 0.0183 |
| Bert4Rec | 0.0313 | 0.0161 | 0.0184 | 0.0092 | 0.0333 | 0.0177 | 0.0363 | 0.0379 | 0.0272 | 0.0131 |
| P5-TID | 0.0485 | 0.0243 | 0.0377 | 0.0198 | 0.0419 | 0.0202 | 0.0372 | 0.0388 | 0.0316 | 0.0154 |
| P5-SemID | 0.0569 | 0.0282 | 0.0405 | 0.0237 | 0.0445 | 0.0231 | 0.0398 | 0.0432 | 0.0324 | 0.0188 |
| TIGER | 0.0546 | 0.0279 | 0.0507 | 0.0284 | 0.0486 | 0.0268 | 0.0438 | 0.0427 | 0.0394 | 0.0208 |
| LC-Rec | <u>0.0648</u> | <u>0.0335</u> | <u>0.0561</u> | <u>0.0307</u> | 0.0538 | 0.0279 | 0.0442 | 0.0451 | 0.0418 | 0.0215 |
| LETTER | 0.0596 | 0.0307 | 0.0556 | 0.0298 | <u>0.0546</u> | <u>0.0291</u> | <u>0.0559</u> | <u>0.0469</u> | <u>0.0426</u> | <u>0.0231</u> |
| **UniTok** | **0.0955** | **0.0496** | **0.0852** | **0.0439** | **0.0902** | **0.0442** | **0.0565** | **0.0476** | **0.0684** | **0.0321** |
| Improve | 47.38% | 48.06% | 51.87% | 42.99% | 65.20% | 51.89% | 1.07% | 1.49% | 60.56% | 38.96% |

Table 7: Performance comparison among UniTok and recommendation competitors for the ten benchmark datasets. Here, the best and second-best performers are highlighted by bold and underline, respectively. The improvements are statistically significant on average ($p = 0.0219 < 0.05$) based on paired t-tests over five runs across all datasets.

While other LLMs (*e.g.*, GPT, BART, and LLaMA) are viable options, our focus is to ensure a fair and controlled evaluation of our proposed tokenization framework, UniTok. Using T5, consistent with other codebook based methods like TIGER (Rajput et al. 2023) and LETTER (Wang et al. 2024), allows us to attribute performance improvements to our method rather than differences in model architecture. Additionally, T5 offers a favorable balance between capability and computational efficiency, which is important for reproducible experimentation. We leave the exploration of other LLMs as promising future work.

For item tokenization, we use 4-level codebook-based identifiers, where each codebook comprises 256 code vectors with an embedding dimension of 32. To obtain item semantic embeddings, we follow (Rajput et al. 2023; Wang et al. 2024) and use a pre-trained language model to encode item content information. The resulting semantic embedding has a dimensionality of 2048, while the hidden embedding dimension is 32. The shared expert is initialized with the average item embedding across all domains, capturing general semantic patterns over domains, whereas the $K = 10$ domain experts—corresponding to ten domains—are initialized with their respective domain-specific average embeddings, naturally inducing domain-aware

assignment without supervision. Following (Rajput et al. 2023; Wang et al. 2024), we set $\alpha$ as 0.25 and $\beta$ as 1. And $\lambda_{\text{RQ}}$ is typically set to 1 and $\lambda_{\text{MI}}$ in the range of $\{0.01, 0.03, 0.05, 0.07, 0.1\}$. We train UniTok for 10k epochs via AdamW [27] optimizer with a learning rate of $1e - 3$ and a batch size of $1,024$. TIGER is fine-tuned for convergence based on the validation performance, with a learning rate in 1e-3, 5e-4 and 1e-4, 2e-4, 3e-4.

All experiments are carried out with Intel (R) 12-Core (TM) E5-1650 v4 CPUs @ 3.60GHz and NVIDIA GeForce RTX 3090 GPUs.

## 5. Complete Set of Experimental Results

To evaluate the recommendation accuracy of UniTok, we compare it with baseline tokenization methods across multiple benchmark datasets, measuring performance using both Recall@10 and NDCG@10 (note that we have shown the results only with respect to NDCG@10 in the main manuscript due to space limit). Unlike existing approaches that train a separate tokenizer per dataset, UniTok is trained once across all domains. As shown in Table 7, UniTok consistently outperforms competitors, achieving up to 65.20% improvement in Recall@10 on the Toys dataset, as confirmed by paired t-tests showing statistically significant improvements ($p < 0.05$). The results demonstrate the effectiveness of our unified tokenization framework; its shared tokenization captures item semantics across diverse domains without domain-specific customization, highlighting the model's scalability and generality.

LLM-based recommender systems using item tokenization, such as TIGER, LC-Rec, LETTER, and UniTok, consistently outperform traditional collaborative filtering methods (*e.g.*, MF, LightGCN, SASRec, and Bert4Rec), benefiting from LLMs' semantic reasoning capabilities. Moreover, incorporating learned tokenization yields further gains over metadata-based approaches (*e.g.*, P5-TID and P5-SemID), as tokenization bridges the item and language domains, enabling discrete, semantically rich item representations that enhance generalization.

## 6. Ablation Study

To assess the contribution of each component in UniTok, we perform an ablation study by progressively removing or modifying its core modules: the TokenMoE module, the shared expert, and the MI calibration part. UniTok-1 removes the TokenMoE and MI calibration; UniTok-2 keeps TokenMoE, but removes the shared expert and MI calibration; and UniTok-3 only removes the MI calibration (note that we have shown the results only on three datasets in the main manuscript due to space limit). UniTok includes all components. As shown in Table 8, removing any module leads to a noticeable drop in recommendation accuracy, which confirms that the combination of modules is crucial to UniTok's effectiveness. In particular, comparing UniTok-1 and UniTok-2 highlights the importance of the TokenMoE module, which significantly improves performance across all datasets by capturing dataset-specific token semantics. Additionally, in comparison with UniTok-3, introducing MI calibration further enhances performance by encouraging the learned representations to retain essential semantic information from the original input space. For clarity and focus, we report results on five representative domains that exhibit trends consistent with those observed across the full range of evaluated datasets.

| Method | Beauty | | Cellphones | | Grocery | | Instruments | | Yelp | |
|---|---|---|---|---|---|---|---|---|---|---|
| | R@10 | N@10 | R@10 | N@10 | R@10 | N@10 | N@10 | R@10 | N@10 | R@10 |
| **UniTok-1** | 0.0558 | 0.0304 | 0.0702 | 0.0371 | 0.0633 | 0.0342 | 0.0926 | 0.0742 | 0.0345 | 0.0177 |
| **UniTok-2** | 0.0896 | 0.0436 | 0.1194 | 0.0606 | 0.0989 | 0.0497 | 0.1273 | 0.0851 | 0.0624 | 0.0281 |
| **UniTok-3** | 0.0915 | 0.0457 | 0.1225 | 0.0622 | 0.1044 | 0.0515 | 0.1327 | 0.0868 | 0.0657 | 0.0303 |
| **UniTok** | **0.0934** | **0.0478** | **0.1251** | **0.0647** | **0.1061** | **0.0533** | **0.1361** | **0.0884** | **0.0684** | **0.0321** |

Table 8: Ablation study results on the Beauty, Cellphones, Grocery, Instrument, and Yelp datasets.

## 7. How Sensitive is UniTok to Key Parameters?

We analyze the sensitivity of UniTok to key hyperparameters, including the number of quantization levels $L$, codebook size $T$, and the loss weights $\lambda_{\text{RQ}}$ and $\lambda_{\text{MI}}$ in Eq. (11) of the main manuscript. The results for the Beauty, Cellphones, and Grocery datasets are shown in Figures 4, 5, and 6, respectively. For clarity and focus, we present results on three representative domains that reflect trends consistent with those observed across the full range of evaluated datasets.

We first vary the number of quantization levels $L$ from 2 to 8 and report NDCG@10 in Figures 4a, 5a, and 6a. Performance improves as $L$ increases from 2 to 4, likely because longer sequences provide better capacity to capture fine-grained semantic information. However, further increasing $L$ beyond 4 leads to performance degradation, which can be attributed to error accumulation in longer autoregressive sequences during recommendations.

Next, we evaluate codebook sizes $T \in \{64, 128, 256, 512\}$, with the results shown in Figures 4b, 4b, 6b. Performance generally improves with larger codebooks, as they provide greater flexibility and token diversity distinguishing items. However, excessively enlarging the codebooks may lead to performance degradation. This may be due to the increased sensitivity to noise in item semantics, which can cause the model to overfit to spurious or less meaningful patterns.
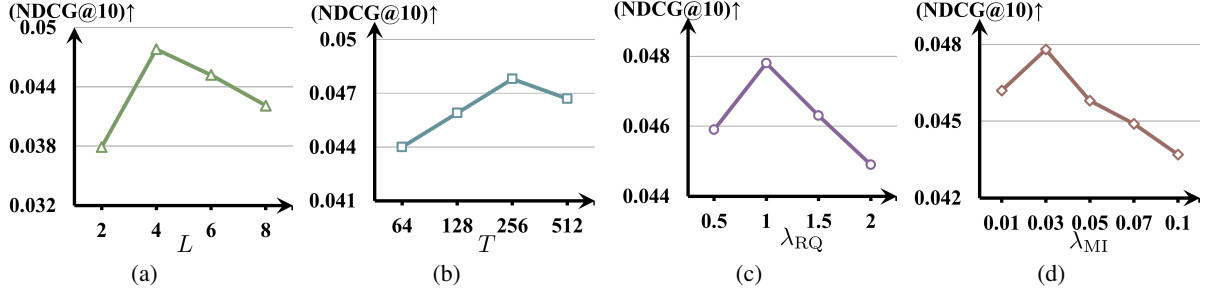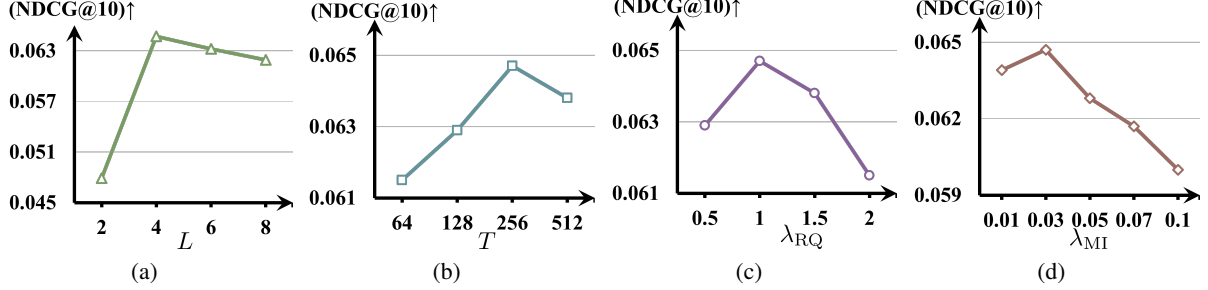
Figure 4: Sensitivity analysis on Beauty.



Figure 5: Sensitivity analysis on Cellphones.

We then examine the impact of the hyperparameter $\lambda_{RQ}$ in Eq.(11) of the main manuscript in terms of NDCG@10. As shown in Figure 4c, 4c, 6c, performance peaks at $\lambda_{RQ} = 1$ across datasets. The results suggest that setting $\lambda_{RQ}$ too high degrades performance, as it overemphasizes the residual quantizations that may be less relevant. Conversely, overly small $\lambda_{RQ}$ may underutilize codebook supervision. Notably, these findings highlight the importance of carefully tuning $\lambda_{RQ}$ to balance codebook-based identifiers and optimize performance.

Finally, we analyze the effect of the hyperparameter $\lambda_{MI}$ in Eq.(11) of the main manuscript in terms of NDCG@10. As shown in Figures 4d, 5d, and 6d, the highest NDCG@10 is achieved at $\lambda_{MI} = 0.03$. The results indicate that setting $\lambda_{MI}$ too high may degrade performance by overemphasizing mutual information calibration, potentially amplifying irrelevant variations. On the other hand, setting $\lambda_{MI}$ too low weakens the influence of mutual information preservation, limiting the semantic alignment of token representations. These findings highlight the importance of properly tuning $\lambda_{MI}$ to balance semantic preservation and generalization for optimal performance.

## 8. Empirical Validation of Theoretical Claims

We empirically validate the technical correctness and practical relevance of **Theorems 1–3**. Entropy analysis (Theorem 1) shows improved token entropy; quantization error comparison (Theorem 2) confirms a lower reconstruction error; and MI variance analysis (Theorem 3) demonstrates enhanced performance stability across domains.

**Entropy analysis of token space (Supporting Theorem 1).** We compare the entropy of the token distributions produced by UniTok and codebook-based methods. Specifically, we analyze how much entropy gain is contributed by the router module in UniTok based on Eq. (12) of the main manuscript. The entropy is calculated based on the frequency of full token combinations across all items.

| Method | Token Space Entropy |
|---|---|
| Codebook-based methods | 9.63 |
| UniTok without the router | 9.63 |
| UniTok (full) | **10.42** |

Table 9: Token space entropy comparison.

As shown in Table 9, we observe that the router alone contributes an additional 0.79 Hartleys of entropy (measured using base-10 logarithm), expanding the capacity of the token space compared to standard codebook-based methods.
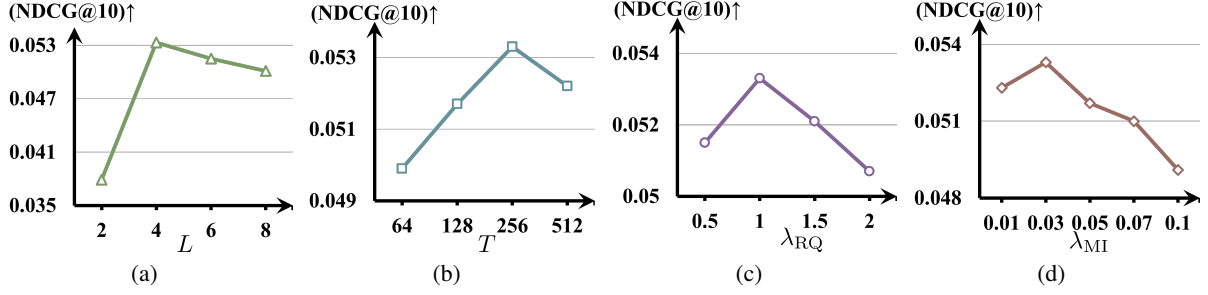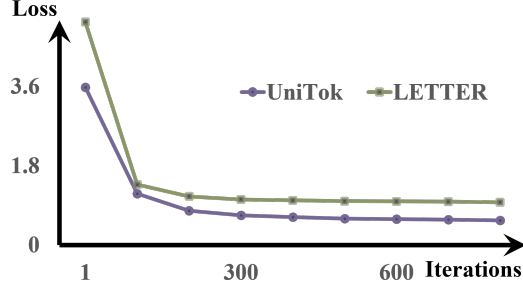
Figure 6: Sensitivity analysis on Grocery.



Figure 7: Comparison of residual quantization loss over training epochs between UniTok and LETTER.

**Quantization error comparison (Supporting Theorem 2).** To empirically support Theorem 2, which states that UniTok yields a lower expected quantization error than that of standard codebook-based methods, we compare their quantization losses on multi-domain setting. As shown in Figure 7, UniTok consistently achieves a lower quantization error than the baseline case with a single set of codebooks. This confirms the theoretical insight that the TokenMoE architecture—by leveraging expert specialization—reduces a representation error and provides more precise item encoding. Notably, the expert-wise partitioning allows UniTok to partially compensate for domain-specific tokenization inaccuracies, leading to improved modeling fidelity in heterogeneous environments.

**MI variance and performance stability (Supporting Theorem 3).** To validate Theorem 3, which posits that the variability in downstream performance across domains is upper-bounded by the variance of MI, we empirically examine the relationship between the variance of MI and the maximum difference in task loss across domains. Specifically, we compute the variance of MI estimates $\widehat{I}^{(k)}$ across domains, denoted as $\text{Var}\left[\widehat{I}^{(k)}\right]$, and compare it against the observed performance variability $\max_{i,j}\left|\mathcal{L}^{(i)} - \mathcal{L}^{(j)}\right|$ (Refer to Eq. (14) of the main manuscript). As shown in Figure 8, we observe a strong positive correlation between $\sqrt{\text{Var}\left[\widehat{I}^{(k)}\right]}$ and performance variability, consistent with the Lipschitz continuity assumption. This trend indicates that lower MI variance leads to more stable performance across domains, supporting our theoretical claim. These results suggest that MI serves as a reliable indicator of multi-domain representation consistency and can guide the design of more robust multi-domain recommendation systems.

# References

Bao, K.; Zhang, J.; Wang, W.; Zhang, Y.; Yang, Z.; Luo, Y.; Chen, C.; Feng, F.; and Tian, Q. 2025. A Bi-step Grounding Paradigm for Large Language Models in Recommendation Systems. *ACM Transactions on Recommender Systems*, 3(4): 1–27.

Bao, K.; Zhang, J.; Zhang, Y.; Wang, W.; Feng, F.; and He, X. 2023. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, 1007–1014.

Dai, S.; Shao, N.; Zhao, H.; Yu, W.; Si, Z.; Xu, C.; Sun, Z.; Zhang, X.; and Xu, J. 2023. Uncovering ChatGPT's Capabilities in Recommender Systems. In *Proceedings of the 17th ACM Conference on Recommender Systems*, 1126–1132.

Geng, S.; Liu, S.; Fu, Z.; Ge, Y.; and Zhang, Y. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *Proceedings of the 16th ACM conference on recommender systems*, 299–315.

Gong, Y.; Ding, X.; Su, Y.; Shen, K.; Liu, Z.; and Zhang, G. 2023. An unified search and recommendation foundation model for
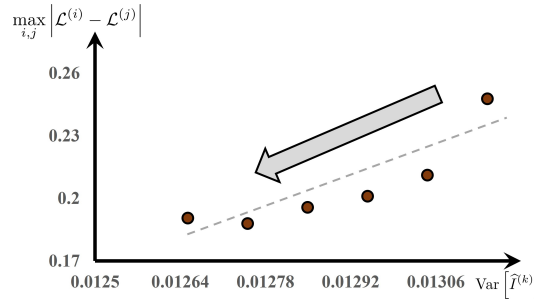
Figure 8: Relationship between MI variance and performance gap.

cold-start scenario. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 4595–4601.

Li, X.; Chen, C.; Zhao, X.; Zhang, Y.; and Xing, C. 2023a. E4srec: An elegant effective efficient extensible solution of large language models for sequential recommendation. *arXiv preprint arXiv:2312.02443*.

Li, Y.; Lin, X.; Wang, W.; Feng, F.; Pang, L.; Li, W.; Nie, L.; He, X.; and Chua, T.-S. 2024. A survey of generative search and recommendation in the era of large language models. *arXiv preprint arXiv:2404.16924*.

Li, Y.; Yang, N.; Wang, L.; Wei, F.; and Li, W. 2023b. Generative Retrieval for Conversational Question Answering. *Information Processing & Management*, 60(5): 103475.

Liao, J.; Li, S.; Yang, Z.; Wu, J.; Yuan, Y.; and Wang, X. 2023. LLaRA: Aligning Large Language Models with Sequential Recommenders. *CoRR*.

Lin, J.; Chen, B.; Wang, H.; Xi, Y.; Qu, Y.; Dai, X.; Zhang, K.; Tang, R.; Yu, Y.; and Zhang, W. 2024a. Clickprompt: CTR models are strong prompt generators for adapting language models to CTR prediction. In *Proceedings of the ACM Web Conference 2024*, 3319–3330.

Lin, J.; Shan, R.; Zhu, C.; Du, K.; Chen, B.; Quan, S.; Tang, R.; Yu, Y.; and Zhang, W. 2024b. Rella: Retrieval-enhanced large language models for lifelong sequential behavior comprehension in recommendation. In *Proceedings of the ACM Web Conference 2024*, 3497–3508.

Lin, X.; Wang, W.; Li, Y.; Feng, F.; Ng, S.-K.; and Chua, T.-S. 2024c. Bridging items and language: A transition paradigm for large language model-based recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 1816–1826.

Lin, X.; Wang, W.; Li, Y.; Yang, S.; Feng, F.; Wei, Y.; and Chua, T.-S. 2024d. Data-efficient Fine-tuning for LLM-based Recommendation. In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*, 365–374.

Liu, Q.; Chen, N.; Sakai, T.; and Wu, X.-M. 2024. Once: Boosting content-based recommendation with both open-and closed-source large language models. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, 452–461.

Prakash, T.; Jalan, R.; Singh, B.; and Onoe, N. 2023. Cr-sorec: Bert driven consistency regularization for social recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, 883–889.

Qiu, Z.; Wu, X.; Gao, J.; and Fan, W. 2021. U-BERT: Pre-training user representations for improved recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 4320–4327.

Ren, X.; Wei, W.; Xia, L.; Su, L.; Cheng, S.; Wang, J.; Yin, D.; and Huang, C. 2024. Representation learning with large language models for recommendation. In *Proceedings of the ACM Web Conference 2024*, 3464–3475.

Wang, Y.; Ren, Z.; Sun, W.; Yang, J.; Liang, Z.; Chen, X.; Xie, R.; Yan, S.; Zhang, X.; Ren, P.; et al. 2024. Enhanced Generative Recommendation via Content and Collaboration Integration. *arXiv e-prints*, arXiv–2403.

Wei, W.; Ren, X.; Tang, J.; Wang, Q.; Su, L.; Cheng, S.; Wang, J.; Yin, D.; and Huang, C. 2024. Llmrec: Large language models with graph augmentation for recommendation. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, 806–815.

Wu, C.; Wu, F.; Qi, T.; and Huang, Y. 2022. Userbert: Pre-training user model with contrastive self-supervision. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2087–2092.

Xi, Y.; Liu, W.; Lin, J.; Cai, X.; Zhu, H.; Zhu, J.; Chen, B.; Tang, R.; Zhang, W.; and Yu, Y. 2024. Towards open-world recommendation with knowledge augmentation from large language models. In *Proceedings of the 18th ACM Conference on Recommender Systems*, 12–22.

Zhai, J.; Mai, Z.; Wang, C.; Yang, F.; Zheng, X.; Li, H.; and Tian, Y. 2025. Multimodal Quantitative Language for Generative Recommendation. In *ICLR, Singapore, April 24-28, 2025*.

Zhang, J.; Xie, R.; Hou, Y.; Zhao, X.; Lin, L.; and Wen, J.-R. 2023. Recommendation as Instruction Following: A Large Language Model Empowered Recommendation Approach. *ACM Transactions on Information Systems*.

Zhang, Y.; Feng, F.; Zhang, J.; Bao, K.; Wang, Q.; and He, X. 2025. Collm: Integrating collaborative embeddings into large language models for recommendation. *IEEE Transactions on Knowledge and Data Engineering*.