# IDGenRec: LLM-RecSys Alignment with Textual ID Learning

Juntao Tan
Rutgers University
New Brunswick, USA
juntao.tan@rutgers.edu

Shuyuan Xu
Rutgers University
New Brunswick, USA
shuyuan.xu@rutgers.edu

Wenyue Hua
Rutgers University
New Brunswick, USA
wenyue.hua@rutgers.edu

Yingqiang Ge
Rutgers University
New Brunswick, USA
yingqiang.ge@rutgers.edu

Zelong Li
Rutgers University
New Brunswick, USA
zelong.li@rutgers.edu

Yongfeng Zhang
Rutgers University
New Brunswick, USA
yongfeng.zhang@rutgers.edu

## ABSTRACT

Generative recommendation based on LLMs have transformed the traditional ranking-based recommendation style into a text-to-text generation paradigm. This approach has attracted significant attention. However, in contrast to standard NLP tasks that inherently operate on human vocabulary, current research in generative recommendations struggles to effectively encode recommendation items within the text-to-text framework using concise yet meaningful ID representations. Due to this unresolved issue, the potential of LLM-based generative recommendation systems remains largely unexplored.

To better align LLMs with recommendation needs, we propose **IDGenRec**, representing each item as a unique, concise, semantically rich, platform-agnostic textual ID using human language tokens. This is achieved by training a textual ID generator alongside the LLM-based recommender, enabling seamless integration of personalized recommendations into natural language generation. Notably, as user history is expressed in natural language and decoupled from the original dataset, our approach suggests the potential for a foundational generative recommendation model. Experiments show that our framework consistently surpasses existing models in sequential recommendation under standard experimental setting. Then, we explore the possibility of training a foundation recommendation model with the proposed method on data collected from 19 different datasets and tested its recommendation performance on 6 unseen datasets across different platforms under a completely zero-shot setting. The results show that the zero-shot performance of the pre-trained foundation model is comparable to or even better than some traditional recommendation models based on supervised training, showing the potential of the IDGenRec paradigm serving as the foundation model for generative recommendation. Code and data are open-sourced at https://github.com/agiresearch/IDGenRec.

## CCS CONCEPTS

• **Information systems → Recommender systems**; • **Computing methodologies → Natural language generation**.

## KEYWORDS

Recommender System, Natural Language Processing

## 1 INTRODUCTION

Generative models with LLMs pre-trained on extensive amounts of information [1, 2, 28, 32] are continually revolutionizing the field of machine learning. Due to their successful in understanding complex instructions and generate creative, contextually relevant predictions, their usage has quickly extended beyond NLP tasks to provide a foundation for applications across diverse research areas. One such example is generative recommendation.

While traditional methods treat recommendation as a retrieval (candidate selection) and ranking process, generative recommendation interprets it as a direct text-to-text generation task: a user's history is expressed as a textual prompt, and the target recommendation is generated in natural language form. However, unlike NLP tasks that are solely reading and generating human language tokens, items in recommendation platforms are individual entities in an ever-growing universe. Therefore, how to encode items as language tokens (i.e., Item IDs) that can be easily integrated into the text-to-text paradigm is a unique and crucial problem within generative recommendation research.

A few attempts have been made to tackle this problem. P5 [8], one of the first works in generative recommendation, proposes allocating out-of-vocabulary (OOV) tokens to items within the recommendation platform. These assigned IDs are fixed-length numerical tokens roughly created based on their sequential appearance in the dataset (e.g., 1001 for the first item, 1002 for the second item). Later, [14] evaluates how recommendation performance can be further improved by using different strategies to initialize the numerical IDs.

However, while these pioneering works achieve considerable performance in standard recommendation settings, the true capability of LLMs in recommendation remains largely unexplored. First, these methods overlook the wealth of semantic information contained in textual descriptions of items, which undermines one of the primary motivations for using LLMs—harnessing the semantic

knowledge gained during their pre-training phase. Second, the numbers assigned to items are meaningless tokens that lack any real contextual meaning. Training such models with a recommendation objective does not lead them to learn the general characteristics of the items. Instead, they merely learn the co-occurrence patterns of these IDs within each dataset. Therefore, although these models present themselves as text-to-text, in essence, they do nothing more than learning representations of each item in a traditional key-value dictionary style, which imposes a ceiling on the quality of the generated recommendations. Simultaneously, since the learned ID representations lack general meanings, the knowledge they acquire is non-transferable across datasets. This means that pre-trained recommendation models lack any zero-shot recommendation ability on unseen data. Consequently, a foundational recommendation model, which has long been pursued in the recommendation community, cannot be achieved by any of the aforementioned methods.

We propose that the limitations mentioned above primarily arise from inadequate item encoding. Consider tasks in human language, such as question answering and machine translation, where every piece of knowledge is represented within a finite set of tokens. This allows a foundational model to easily learn universally applicable knowledge from training on large text corpus and to adapt effortlessly to any downstream task. If, however, items in recommendation systems were also fully represented using human vocabulary, with each item described by a specific set of natural language tokens, then the capabilities of LLMs could more closely align with the requirements of recommendation systems. In this way, by training on recommendation-specific corpora, LLMs would be able to learn genuine recommendation-related knowledge, which could significantly improve the models' accuracy and generalizability in recommendation tasks.

Hence, we suggest that the ideal IDs in generative recommendation should possess the following properties: 1) They should be textual IDs composed of tokens originally processed by the pre-trained LLMs; 2) They should be meaningful, informative, and suitable for recommendation purposes; 3) The generated IDs should be short yet unique, effectively identifying the recommendation items. However, IDs that meet such stringent requirements are clearly not available in existing item information. Therefore, in this paper, we propose training an ID generator that automatically learns a textual ID for each item that fulfills the above criteria. The new framework, named as **IDGenRec**, treat ID generation as another text-to-text process. As shown in Figure 1, the ID generator, which is also a language model, takes an item's metadata (i.e., all available textual information about the item) and produces qualified textual IDs. Consequently, the user's history and the target item for recommendation can be represented in natural language, without any "uncontextualized" tokens, thus making it suitable for training an LLM-based generative recommender. This overall process is illustrated in Figure 2. Notably, by considering all items' text in the user's history, the same ID generator can produce another textual ID, serving as the user's ID that represents a "high-level profile" of the user's preferences. The creation of user ID is optional, and we will provide ablation study results in the experiments.

Many challenges lie in this work, and we propose related strategies to address each of them in the paper, including:

**Table 1: Comparison of LLM-based recommendation models: P5 and its variant versions are generative models but not foundation models due to the use of OOV tokens. UniSRec and Recformer have encoder-only structures and are, therefore, not generative models. Additionally, Recformer employs a rigorously defined item text template that is specifically designed for the Amazon dataset only, and is thus classified as partly a foundation model.**
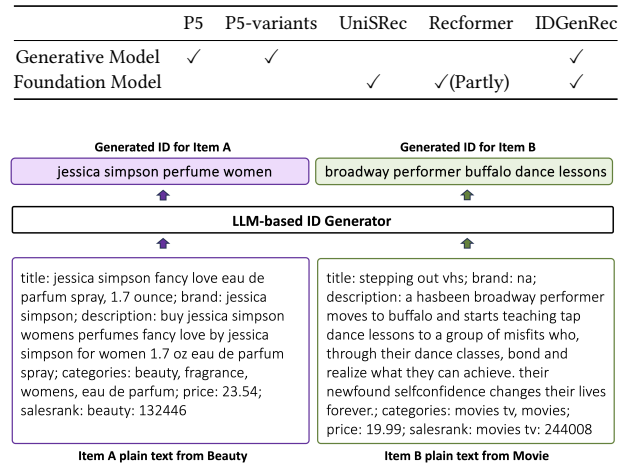
| | P5 | P5-variants | UniSRec | Recformer | IDGenRec |
|---|---|---|---|---|---|
| Generative Model | ✓ | ✓ | | | ✓ |
| Foundation Model | | | ✓ | ✓(Partly) | ✓ |



**Figure 1: The ID generator takes plain text from each item's meta textual information and generates abstractive textual IDs for the item's representation.**

(1) The ID generator should understand lengthy metadata that may include unnecessary information, and should generate tokens that cover the crucial details of the item which are important for recommendations. For this purpose, we have selected a T5 model originally trained for article tag generation and fine-tuned it with recommendation objectives.

(2) The generated IDs should be short yet unique, suitable for identifying the recommendation items. However, the automatically generated IDs may not always satisfy the uniqueness criterion, especially as the number of items increases. Therefore, we propose a diverse ID generation algorithm to always ensure each item has a unique ID allocated.

(3) Since the framework relies on collaboration between two LLMs—the ID generator and the base recommender—a meticulously designed training strategy is required to enable seamless collaboration between them. We propose an alternate training strategy that trains the LLM-based ID generator and the base recommender asynchronously, ensuring that their learned knowledge is well-aligned.

We note that some recent LLM-based recommendation models employ an encoder-only structure, similar to BERT, as their representation function. Some of these models [12, 18] also possess (or partially possess) characteristics of a foundational recommendation model, though they do not function as generative models. The distinctions among these models are depicted in Table 1. However, generative models present several advantages over discriminative methods. These include transforming the retrieval and ranking processes into a more streamlined generative process, eliminating

the need of one-by-one item score calculation, and leveraging the extensive knowledge embedded in pre-trained generative LLMs. Nonetheless, these encoder-only methods remain valuable as baselines for zero-shot evaluation.

We conduct two types of experiments to demonstrate the effectiveness of our proposed method. First, we evaluate the method under a standard sequential recommendation setting. Experiments on 4 widely-used public datasets show its significant improvements compared to sequential recommendation baselines, including both traditional and generative models. Then, to explore the possibility of training a foundational generative model that learns general recommendation knowledge, inspired by the training paradigm of LLMs in standard NLP tasks, we compile user histories from 19 datasets from the Amazon Review Datasets, encompassing a diverse range of recommendation domains, to build a massive recommendation training corpus. After training the model on this extensive dataset, we directly apply the foundational model to 6 unseen datasets either intra- or inter-platform and evaluate the recommendation performance in a completely zero-shot setup. The results show very promising recommendation performance, even surpassing many traditional recommendation models that are based on supervised training.

## 2 APPROACH

We first introduce the generation process, including how the prompts are constructed and how the generated IDs are integrated into the text-to-text format, as discussed in Section 2.1. Then, in Sections 2.2 and 2.3, we describe how the IDs are generated by utilizing the metadata of items, employing a diverse ID generation algorithm to ensure that the IDs are unique for each item. With these generated IDs, the base recommender system is presented in Section 2.4. Finally, in Section 2.5, we demonstrate how the ID generator and recommender are trained alternately with respect to the recommendation objective, ensuring that they work collaboratively and effectively.

### 2.1 Generative Process

In this study, we introduce the foundational generative model within the context of sequential recommendation systems. This approach is particularly apt as it aligns naturally with the sequential representation of user history, serving as an input prompt. Building upon existing work in generative recommendation systems [8, 14, 38, 41], our model requires predefined prompt templates for the generation process. For instance, a typical template could be: "User [user_ID] has purchased items [item_ID], [item_ID], …, [item_ID]; predict the next possible item to be bought by the user." In this template, all item IDs from the user's history are interpolated at the placeholders, preserving their sequential order. Besides, a user ID can also be produced by the same ID generator, taking the meta information of all items in the user history. Generating a user ID is optional and proved to be beníficial to the recommendation performance in experiments. We have developed 10 such templates, each with minor differences from the others, randomly selecting one of them for each training instance to ensure the model focuses more on recommendation-relevant information such as the item sequence and less on the exact format of the prompt. An example

of a completed prompt is illustrated in Figure 2. The decoder then generates the target item "[target_item_ID]" token by token, which is able to uniquely identify the target recommendation. Later, we will detail the two primary LLM components of our generative process: the ID generator and the base recommender.

### 2.2 ID Generator

The ID generator is a generative model that produces item IDs using the item's meta-information. This meta-information encompasses all textual data related to the item, including both relevant and irrelevant aspects for recommendation purposes. Potential elements of this information may comprise the item's title, category, price, general description, creation time, popularity, location, etc. The specific content largely depends on the platform and dataset. Although the meta-information is typically presented in a key-value dictionary format, we convert it into plain text during processing, such as "name: zeppelin; categories: cocktail bars, restaurants; stars: 4.0; …" This allows the ID generator to freely learn which pieces of information should be prioritized when generating IDs.

Consider an item whose plain description is a lengthy sequence of tokens $w = [w_1, w_2, \ldots, w_m]$. Token embeddings will be generated from $w$ by the model's parameters and combined with position embeddings before being fed into the language model. The combination with position embeddings will be omitted in the rest of the formulations since this process is common in large language models. The output of the ID generator will be a concise set of ID tokens $d = [d_1, d_2, \ldots, d_n]$, where $n \ll m$. When generating each token of the item ID, the model attends to both the item's entire plain description and the previously generated ID units. Thus, the probability of a generated ID is denoted as:

$$p(d_1, \cdots, d_n) = \prod_{i=1}^{n} p_\theta(d_i | d_{<i}, w) \tag{1}$$

Where $d_{<i} = [d_1, \cdots, d_{i-1}]$, and $\theta$ are the parameters of the ID generator. This process is illustrated in Figure 1.

### 2.3 Diverse ID Generation

The two primary properties of the generated IDs are: 1) the IDs should have a reasonable length, and 2) the IDs should be unique. However, these two properties are somewhat contradictory: a set of IDs constrained to a shorter length is more likely to result in duplicates when generated by the ID generator. Therefore, we propose an algorithm to ensure that the generated IDs are both short and unique. The core concept of the algorithm is fundamentally based on diverse beam search (DBS) [33], a method commonly used in sentence generation. It is a variation of the standard beam search that is designed to generate a more diverse set of sequences. DBS partitions the beams into groups and introduces a diversity penalty, denoted as $\lambda$, as a hyperparameter to discourage the selection of similar sequences within the same group. A higher value of $\lambda$ promotes more diversity, while a lower value of $\lambda$ gives more weight to the model's probabilities, potentially leading to less diverse outputs.

The ID generator employs DBS in generating IDs. To ensure the uniqueness of the generated IDs, the algorithm generates $k$ groups of IDs each time and compares each generated ID against a set of already existing IDs. If a duplicate is detected, the algorithm
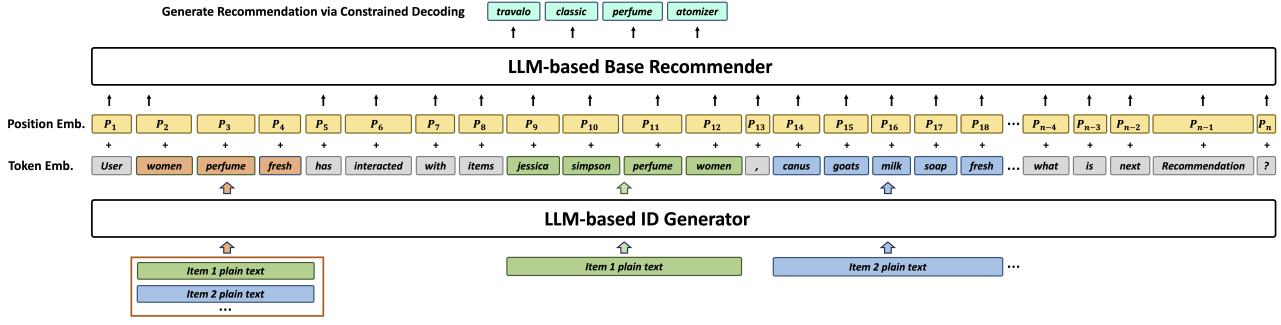
**Figure 2: A real example showing the generative recommendation workflow. The ID generator generates item IDs for items from the user history by taking their plain text. Then, the generated IDs are interpolated into the template. Addtionally, the user ID is generated by using all items' text in the user's history, showing a "high-level profile" of the user's preference. The position embeddings are subsequently combined with token embeddings to capture the sequence of interactions. Finally, the base recommender generates the ID of the recommended item based on constrained decoding.**

increases the diversity penalty in the beam search. This increase in penalty continues until a unique ID is produced, or until the diversity penalty reaches a pre-set maximum threshold (e.g., 10 in this paper). In cases where the maximum penalty is insufficient to generate a unique ID, the algorithm extends the permissible ID length and repeats this process with the initial diversity penalty. Algorithm 1 elaborates on this process.

---

**Algorithm 1** Diverse ID Generation Algorithm

---

1: Initialize set $\mathcal{U}$ to store unique IDs
2: Initialize diversity penalty $\lambda$ to 1
3: Initialize ID length limit $L$ to 10
4: **for** each item in the dataset **do**
5:    Initialize $found$ as False
6:    **while** not $found$ **do**
7:       Generate $k$ IDs using ID Generator with current $L$ and $\lambda$
8:       **for** each generated ID $id$ **do**
9:          **if** $id$ not in $\mathcal{U}$ **then**
10:            Add $id$ to $\mathcal{U}$ and save the item-ID pair
11:            Set $found$ to True
12:            **break**
13:          **end if**
14:       **end for**
15:       **if** not $found$ **then**
16:          $\lambda \leftarrow \lambda + 1$
17:          **if** $\lambda$ exceeds predefined limit **then**
18:            Increase $L$ and reset $\lambda$ to 1
19:          **end if**
20:       **end if**
21:    **end while**
22: **end for**

---

## 2.4 Base Recommender

After generating item IDs and incorporating them into the prompt template, the system tokenizes the completed prompt and feeds it into the LLM-based recommender. The recommender then generates the tokens of the target recommended item ID in an autoregressive manner.

To ensure that the decoded ID corresponds to an actual existing item, we adopt a constrained sequence decoding strategy [5]. More specifically, a prefix tree is used to store all generated candidate IDs. Each newly generated token is constrained by the previously generated tokens, ensuring that the generation process only considers tokens that can potentially form an existing candidate ID in the dataset. Suppose the completed input prompt is denoted as a sequence of tokens $\boldsymbol{x} = [x_1, x_2, \cdots, x_n]$. In this case, the base recommender model aims to generate $\boldsymbol{y} = [y_1, y_2, \cdots, y_n]$, where $\boldsymbol{y}$ is an ID for a real item in the dataset. We define $\mathcal{V}(y_{<i})$ as the subset of valid tokens in the vocabulary, constrained by the prefix tree with previously generated tokens as nodes. The generation of each token by the decoder is defined as:

$$p(y_i|y_{<i}, \boldsymbol{x}) = \begin{cases} p_\phi(y_i|y_{<i}, \boldsymbol{x}) & \text{if } y_i \in \mathcal{V}(y_{<i}), \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Therefore, the probability of the recommendation of a target item $[y_1, \cdots, y_n]$ is:

$$p(y_1, \cdots, y_n) = \prod_{i=1}^{n} p_\phi\big(y_i|y_{<i}, \boldsymbol{x}, \mathcal{V}(y_{<i})\big) \quad (3)$$

## 2.5 Alternate Training

Training the ID generator and the base recommender are two separate but interdependent tasks. The ID generator is trained to produce optimal IDs that the base recommender can easily interpret. Concurrently, the base recommender adjusts its parameters to enhance the correct recommendation for items, each represented by its currently generated ID. Training both components simultaneously may result in an unstable training process. Hence, we propose alternating the training sessions of the ID generator and the base recommender, proceeding through a specified number of iterations. This approach involves asynchronously updating the IDs between the two training phases of the base recommender for better integration and performance.

Both the ID generator and the base recommender are trained by minimizing the negative log-likelihood of the final prediction made by the recommendation pipeline compared to the ground-truth target item ID.

*2.5.1  Training Base Recommender.* At each round of training the base recommender, we pre-compute item IDs for all items using the ID generator at that point. For each user in the training data, the current IDs of the items in the user's history are filled into the sampled template to complete the input prompt $x$. Then, the base recommender $\omega$ is trained with a common teacher forcing strategy [34], i.e., the loss of the next token is computed under the ground-truth value of the previous token:

$$\mathcal{L}_{\text{rec}} = -\sum_{i=1}^{|\boldsymbol{y}|} \log P_\omega(y_i|y_{<i}, \boldsymbol{x}) \tag{4}$$

*2.5.2  Training ID Generator.* During this training process, all parameters in the base recommender are fixed, and only the ID generator is updated. The goal is for the ID generator to produce IDs that are suitable for the base recommendation model.

Since the output of the ID generator is a set of discrete tokens (IDs), it is inherently non-differentiable. This poses a challenge for training the model using gradient-based optimization techniques, as gradients cannot flow back through these discrete outputs. To circumvent this, for each item in the user history, we calculate the output logits of each token by the ID generator across all vocabulary, denoted as $\text{Logits}_\phi(\mathcal{V})$, where $\phi$ is the ID generator model and $\mathcal{V}$ is the vocabulary. We then compute the average embedding for each token of the ID through the parameters of the base recommendation model, denoted as $\text{Emb}_\omega(\text{Logits}_\phi(\mathcal{V}))$, where $\omega$ is the base recommender model. This creates a continuous, differentiable representation of the generated IDs. These ID embeddings are then directly interpolated into the prompt template at the related positions at the embedding level. We use $\text{Emb}_{\text{interp}}$ to represent this completed input embedding. In this way, the ID generator $\phi$ can be trained, guided by the loss computed from the recommendation output. This process is formulated as:

$$\mathcal{L}_{\text{id}} = -\sum_{i=1}^{|\boldsymbol{y}|} \log P_\omega\left(y_i \mid y_{<i}, \text{Emb}_{\text{interp}}\right)$$

where $\text{Emb}_{\text{interp}} = \text{Insert}\left(\text{Emb}_\omega\left(\text{prompt}\right), \text{Emb}_\omega\left(\text{Logits}_\phi(\mathcal{V})\right)\right)$ (5)

The parameters of the base recommender model (i.e., $\omega$) are fixed, and the loss is only backpropagated to the ID generator (i.e., $\phi$), ensuring that the IDs generated capture the essential characteristics of each item, as determined by their meta-information, in a format that the base recommendation model can effectively interpret.

## 2.6  Model Initialization

We choose the T5 model [28] as the backbone for both the ID generator and the base recommender for two main reasons: 1) To maintain the model's simplicity, as this paper does not aim to conduct extensive empirical studies of LLM structures, but rather focuses on the core concept of ID generation; 2) To ensure a fair comparison with previous generative recommendation works, which are also based on T5, thereby demonstrating that the improvement

in recommendation ability comes solely from a more elegant ID selection.

For the base recommender, the standard pretrained T5 checkpoint is chosen as the backbone to incorporate pre-learned knowledge into the recommendation task. For the ID generator, given that generating IDs from lengthy texts is non-trivial and highly task-specific, a more dedicated starting point is preferable. Consequently, we select a T5 small model fine-tuned on the article tag generation task[1], as the initial configuration for the ID generator. This model, trained on 190k Medium articles[2], is adept at generating concise tags from article textual content. This selection is driven by the significant similarity between tag generation for news articles and the summarization of items in a few words.

## 3  EXPERIMENT

Our experiments comprise two components: the first is an evaluation of standard sequential recommendation to compare the basic supervised learning capabilities of our model against widely-used baselines. The second component is a zero-shot evaluation designed to assess the model's potential as the backbone for a foundational generative recommender system. We will begin by introducing the datasets used in the experiments and detailing our model's training process. Subsequently, we will discuss the experimental results for each of the two experimental settings.

## 3.1  Datasets

For the standard evaluation of sequential recommendation, we selected four widely-used datasets. Three of them, namely Sports, Beauty, and Toys, are from the Amazon review dataset [10, 24], along with another dataset from Yelp[3]. These datasets are also used in previous papers [8, 14, 44], and we follow the exact data processing steps to filter out users and items with fewer than 5 interactions, thereby allowing for a fair comparison with all the baselines.

For the zero-shot experiments for foundational model evaluation, we have two groups of datasets: pre-training datasets and testing datasets. The pre-training datasets are all from the Amazon review dataset, containing various domains, and the testing data are selected from both Amazon review datasets (intra-platform) and the Yelp dataset (inter-platform). We propose a detailed data selection rule for deciding which data are used as pre-training datasets and which are used for testing, as follows:

First, we split all 24 Amazon review datasets into different groups according to their densities, as shown in Table 2. Guided by their density range, we include Sports, Beauty, Toys, Music, and Instruments in the test datasets. These datasets cover all the density categories across Amazon review datasets. Besides, including Sports, Beauty, and Toys in the test datasets provides a better view of the foundational model's ability compared to traditional models, since the three datasets are previously used in standard evaluation. Along with Yelp, these form all the test datasets for zero-shot evaluation.

We have 19 Amazon datasets left, which are used to create the massive recommendation corpus for training the foundational

---

[1]https://huggingface.co/nandakishormpai/t5-small-machine-articles-tag-generation
[2]https://www.kaggle.com/datasets/fabiochiusano/medium-articlesdataset
[3]https://www.yelp.com/dataset

**Table 2: Amazon review datasets categorized by density**

| Density Range | Amazon Datasets |
|---|---|
| Den. ≥ 0.5 | Instruments, Patio |
| 0.1 ≤ Den. < 0.5 | Automotive, Instant, Office, Music, Grocery, Baby |
| 0.05 ≤ Den. < 0.1 | Tools, Pet, Toys, Phones, Beauty, Games, Apps |
| Den. ≤ 0.05 | Clothing, Sports, Health, Home, Kindle, CDs, Electronics, Movies, Books |

**Table 3: Dataset Statistics**

| Category | Datasets | # Users | # Items | # Interactions | Density |
|---|---|---|---|---|---|
| | Sports | 35,598 | 18,357 | 296,337 | 0.0453% |
| Std. Eval. | Beauty | 22,363 | 12,101 | 198,502 | 0.0734% |
| | Toys | 19,412 | 11,924 | 167,597 | 0.0724% |
| | Yelp | 30,431 | 20,033 | 316,354 | 0.0519% |
| Pre-training | Fusion | 183,918 | 233,323 | 2,875,446 | 0.0067% |
| | Sports | 35,598 | 18,357 | 296,337 | 0.0453% |
| | Beauty | 22,363 | 12,101 | 198,502 | 0.0734% |
| Zero-shot | Toys | 19,412 | 11,924 | 167,597 | 0.0724% |
| | Music | 5,541 | 3,568 | 64,706 | 0.3273% |
| | Instruments | 1,429 | 900 | 10,261 | 0.7978% |
| | Yelp (Cross Platform) | 30,431 | 20,033 | 316,354 | 0.0519% |

model. Since their sizes vary extremely (e.g., Books contains 603,668 users and Automotive only has 2,928 users), we randomly down-sample the large datasets to only include $30,000$ users, which is around the median number of users of the Amazon datasets. All the selected recommendation records together form a "Fusion" dataset for training the foundation recommendation model. Complete and detailed data statistics can be seen in Table 3.

## 3.2 Implementation Details

We use SentencePiece [17] with a vocabulary size of $32,128$ as the tokenizer. The predefined templates for sequential recommendation are largely adopted from P5 [8], with the only difference being that the "dataset" information is removed from the templates, as our model is more generalized and possesses cross-dataset capabilities.

In the diverse ID generation algorithm, we set $k = 10$ as the number of groups for DBS and start with $\lambda = 1$ as the initial diversity penalty. This penalty is increased by 1 in each iteration until the algorithm successfully generates a unique item ID. If $\lambda$ reaches 10 without producing a valid ID, we then extend the length limit for the item ID. This iterative process of adjusting the diversity penalty and, if necessary, the length of the ID, ensures that the Diverse ID Generation algorithm can successfully produce IDs that are both succinct and distinct. Initially, the token length is set within the range of $[1, 10]$, but if the algorithm is unable to generate a unique ID once the diversity penalty has reached its maximum, we increase the token length limit to the range of $[10, 20]$. In our experiments, even for item IDs in the largest ID set, specifically the Fusion dataset, only 11.20% of the items required a second attempt at ID generation with an increased diversity penalty, and merely 3.72% of the items necessitated an extension in ID length.

In the standard recommendation experiments, where the model is trained on a single dataset, we begin by training the ID generator for 1 epoch, followed by training the base recommender for 10 epochs, for a total of 3 iterations. The learning rate is set to $1e - 3$ for training the base recommender and $1e - 8$ for training the ID generator. This approach is applied to all the datasets in the standard recommendation setting.

Regarding the training of the foundational model for zero-shot recommendation, we find that training the model for one epoch achieved the best performance when tested on the test datasets. Further training negatively affected the zero-shot performance. Interestingly, a similar pattern was also observed in NLP research [26, 32], where pre-training LLMs for more than one epoch led to overfitting on the training data. This observe suggests that the proposed framework narrows the gap between recommendation and language generation.

## 3.3 Evaluation Metrics

In all experiments, we evaluate the ranking performance of the recommendation models using Normalized Discounted Cumulative Gain (NDCG) at 5 and 10, as well as Hit Ratio (HR) at 5 and 10. For fair comparison with baselines, we adopt a leave-one-out strategy for testing. Meanwhile, negative sampling is not used in the evaluation. Instead, we rank over all items for evaluation.

## 3.4 Exp1: Standard Evaluation

*3.4.1* ***Baselines***. The baselines for standard evaluation include two types of recommendation methods:

- Traditional sequential recommendation methods, including GRU4Rec [11], Caser [31], HGN [23], SASRec [16], Bert4Rec [29], FDSA [43], and S3Rec [44]. These models are popular and cover various model structures, with GRU4Rec based on RNN, Caser on CNN, and SASRec, HGN, Bert4Rec, FDSA, and S3Rec on transformers.
- Generative recommendation methods, including P5 [8, 38] and its variations [14] with different ID generation strategies. P5-SID generates numerical item IDs with respect to their sequential appearance order in the dataset. P5-CID creates item IDs guided by collaborative filtering. P5-SemID uses item category as semantic information to construct item IDs. More details can be seen in [14].

*3.4.2* ***Results***. The standard evaluation results are shown in Table 4. Generally speaking, our proposed method significantly outperforms all other baselines. When compared to the second best baseline on each dataset (highlighted with an underline), our method shows improvements of 39.44%, 23.55%, 42.37%, and 36.76% on the Sports, Beauty, Toys, and Yelp datasets, respectively. These improvements are calculated by averaging the increase across four metrics relative to the corresponding best performance from baselines, i.e., the bold value against the underline value of each row for each dataset. Such results underscore the robustness and significance of the improvements brought by our generative recommendation method. Besides, according to comparisons with P5 and its variants, representing items with learned textual IDs indeed better utilizes the semantic understanding abilities of LLMs.

*3.4.3* ***Ablation Studies***. We conduct ablation studies regarding two components of the model design.

**Table 4: Standard evaluation for single dataset recommendation. All improvements are significant at $p < 0.05$ compared to the best baseline under the student's t-test.**

| Dataset | Metric | GRU4Rec | Caser | HGN | SASRec | Bert4Rec | FDSA | S3Rec | P5-SID | P5-CID | P5-SemID | **IDGenRec** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | HR@5 | 0.0129 | 0.0116 | 0.0189 | 0.0233 | 0.0115 | 0.0182 | 0.0251 | 0.0264 | 0.0313 | 0.0274 | **0.0429** |
| Sports | NDCG@5 | 0.0086 | 0.0072 | 0.0120 | 0.0154 | 0.0075 | 0.0122 | 0.0161 | 0.0186 | 0.0224 | 0.0193 | **0.0326** |
| | HR@10 | 0.0204 | 0.0194 | 0.0313 | 0.0350 | 0.0191 | 0.0288 | 0.0385 | 0.0358 | 0.0431 | 0.0406 | **0.0574** |
| | NDCG@10 | 0.0110 | 0.0097 | 0.0159 | 0.0192 | 0.0099 | 0.0156 | 0.0204 | 0.0216 | 0.0262 | 0.0235 | **0.0372** |
| | HR@5 | 0.0164 | 0.0205 | 0.0325 | 0.0387 | 0.0203 | 0.0267 | 0.0387 | 0.0430 | 0.0489 | 0.0433 | **0.0618** |
| Beauty | NDCG@5 | 0.0099 | 0.0131 | 0.0206 | 0.0249 | 0.0124 | 0.0163 | 0.0244 | 0.0288 | 0.0477 | 0.0299 | **0.0486** |
| | HR@10 | 0.0283 | 0.0347 | 0.0512 | 0.0605 | 0.0347 | 0.0407 | 0.0647 | 0.0602 | 0.0680 | 0.0652 | **0.0814** |
| | NDCG@10 | 0.0137 | 0.0176 | 0.0266 | 0.0318 | 0.0170 | 0.0208 | 0.0327 | 0.0368 | 0.0357 | 0.0370 | **0.0541** |
| | HR@5 | 0.0097 | 0.0166 | 0.0321 | 0.0463 | 0.0116 | 0.0228 | 0.0443 | 0.0231 | 0.0215 | 0.0247 | **0.0655** |
| Toys | NDCG@5 | 0.0059 | 0.0107 | 0.0221 | 0.0306 | 0.0071 | 0.0140 | 0.0294 | 0.0159 | 0.0133 | 0.0167 | **0.0481** |
| | HR@10 | 0.0176 | 0.0270 | 0.0497 | 0.0675 | 0.0203 | 0.0381 | 0.0700 | 0.0304 | 0.0327 | 0.0376 | **0.0870** |
| | NDCG@10 | 0.0084 | 0.0141 | 0.0277 | 0.0374 | 0.0099 | 0.0189 | 0.0376 | 0.0183 | 0.0170 | 0.0209 | **0.0551** |
| | HR@5 | 0.0176 | 0.0150 | 0.0186 | 0.0170 | 0.0051 | 0.0158 | 0.0201 | 0.0346 | 0.0261 | 0.0202 | **0.0468** |
| Yelp | NDCG@5 | 0.0110 | 0.0099 | 0.0115 | 0.0110 | 0.0033 | 0.0098 | 0.0123 | 0.0242 | 0.0171 | 0.0131 | **0.0368** |
| | HR@10 | 0.0285 | 0.0263 | 0.0326 | 0.0284 | 0.0090 | 0.0276 | 0.0341 | 0.0486 | 0.0428 | 0.0324 | **0.0578** |
| | NDCG@10 | 0.0145 | 0.0134 | 0.0159 | 0.0147 | 0.0090 | 0.0136 | 0.0168 | 0.0287 | 0.0225 | 0.0170 | **0.0404** |

First, we assess how critical the alternate training strategy is for the model, compared to 1) training only the ID generator while the base recommender remains fixed with the default pre-trained T5 parameters; 2) training only the base recommender with the item IDs directly generated by the T5 model trained on article tag generation. The total number of training epochs is the same as when each is trained in the alternate training setup, i.e., $1 \times 3$ epochs for the ID generator and $10 \times 3$ epochs for the base recommender. The results, shown in Table 5, indicate that alternate training significantly boosts the model's performance by enabling better collaboration between the two components. Moreover, training only the ID generator does not yield good results. However, training only the base recommender with initially generated IDs still outperforms all baselines.

Second, we evaluate whether generating optional user IDs enhances recommendation performance, as shown in Table 6. In summary, while using only item IDs already allows the model to make excellent recommendations in the standard setting, including user IDs does prove beneficial. Nonetheless, relying solely on user IDs does not provide enough information for making recommendations.

*3.4.4* ***Case Studies***. As generating item IDs using human vocabulary is the core concept of the proposed method, we conduct an extensive quality study with real examples of the generated item IDs from their plain text information, as shown in Figure 3. Specifically, we select two examples from each dataset: one with lengthy plain text data and one with shorter plain text data. For each example, we present the ID generated by the initial ID generator and the ID generated by the fine-tuned ID generator at the end of alternate training. Initially, the ID generator tends to select the first few words from the item's meta information, regardless of their relevance. Excitingly, after training, the ID generator is capable of selecting words that more accurately represent the items. The learned IDs are generally more informative and representative, containing less extraneous information such as numbers.

**Table 5: Comparison of recommendation accuracy for different training strategies: 1) ID-only: Training only the ID generator. 2) Rec-only: Training only the base recommender. 3) Alternate: Training ID generator and base recommender alternately for 3 iterations.**

| | | ID-only | Rec-only | Alternate |
|---|---|---|---|---|
| | HR@5 | 0.0102 | 0.0350 | **0.0429** |
| Sports | NDCG@5 | 0.0070 | 0.0271 | **0.0326** |
| | HR@10 | 0.0155 | 0.0461 | **0.0574** |
| | NDCG@10 | 0.0087 | 0.0307 | **0.0372** |
| | HR@5 | 0.0111 | 0.0601 | **0.0618** |
| Beauty | NDCG@5 | 0.0067 | 0.0442 | **0.0486** |
| | HR@10 | 0.0192 | 0.0797 | **0.0814** |
| | NDCG@10 | 0.0093 | 0.0505 | **0.0541** |

## 3.5 Exp2: Zero-shot Evaluation

We test the potential of our model to serve as a foundation recommendation model. The experiment is conducted by first training the model on the fusion dataset, i.e., an extensive recommendation corpus collected from 19 recommendation datasets from Amazon Review, then directly test its recommendation performance on each test dataset in a completely zero-shot setting (all the users and items are not seen in the training dataset).

*3.5.1* ***Baselines***. UniSRec [12] is the main comparable baseline in the zero-shot setting, which uses an encoder-only structure to learn item representations from their meta text information. We note that UniSRec incorporates some human-inductive knowledge in data preprocessing, e.g., the "title," "category," and "brand" information are carefully selected for Amazon datasets. For the Amazon dataset, we still use the pre-selected categories as in the original paper to showcase its best performance, therefore introducing a slight advantage for UniSRec. For Yelp, we use all the plain text information. Since the patterns have not been seen in both models, this is a fair comparison.

**Dataset: Yelp**

| Plain Text: | name: zeppelin; categories: cocktail bars, tapassmall plates, restaurants, nightlife, bars; stars: 4.0; review_count: 146; address: 235 w tremont ave; city: charlotte; state: nc; false; 2; goodforkids: false; casual; happyhour: true; outdoorseating: true; goodformeal: dessert: false, latenight: false, lunch: false, dinner: true, brunch: false, breakfast: false; goodfordancing: false; wifi: free; music: dj: false, false, no_music: false, jukebox: false, live: false, video: false, karaoke: false; garage: false, street: true, validated: false, lot: true, valet: false; ambience: touristy: false, hipster: false, romantic: false, divey: false, intimate: false, trendy: false, upscale: false, classy: true, casual: false; true; true; bestnights: monday: false … | name: richards window tinting; categories: home services, auto glass services, home window tinting, automotive, car window tinting; stars: 5.0; review_count: 58; address: 3863 s valley view blvd, ste 9; city: las vegas; state: nv; true; true; wifi: ufree |

| Initial ID: | zeppelin zeppel | richards window tinting categories home |
| Fine-tuned ID: | zeppelin cocktail bars tap | richards window tinting auto glass services |

**Dataset: Beauty**

| Plain Text: | title: farouk chi pro gf 1505 1300 watt ceramic anion infared low emf professional hair dryer; categories: beauty, hair care, styling tools, hair dryers; brand: na; description: the chi pro hair dryer is the latest in advanced american technology. this lightweight quiet dryer utilizes ceramic technology, negative ions, far infrared and low emf. chi pro hair dryer includes the chi diffuser. all chi products that are purchased direct from amazon.com are covered by a 30 day amazon warranty. this item is not covered by the manufacture warranty. applies only to products sold by amazon.com. does not apply to products sold on our site by thirdparty merchants or through thirdparty areas such as amazon.com … | title: truth by calvin klein for women, eau de parfum spray, 3.4 ounce; categories: beauty, fragrance, womens, eau de parfum; brand: na; description: launched by the design house of calvin klein in 2000, truth perfume is classified as a refreshing, oriental, woody fragrance. this feminine scent possesses a blend of fresh, natural rich woods that react to individual skin chemistry. it is recommended for casual wear.; price: 31.69; salesrank: beauty: 71873 |

| Initial ID: | farouk chi pro gf 15 | truth by calvin klein eau de |
| Fine-tuned ID: | chi pro hair dryer chi diffuser | truth perfume calvin klein oriental |

**Dataset: Sports**

| Plain Text: | title: kor one bpa free 750ml water bottle; categories: sports outdoors, accessories, sports water bottles; brand: na; description: the average u.s. consumer discarded 167 singleuse water bottles last year. of these 167, only 20 were recycled; the rest of the oilbased bottles went into landfills or wound up as litter. now you can do your part to promote a greener lifestyleand enjoy a strikingly beautiful design in the processwith the kor one hydration vessel. offering everything you want in a multiuse water bottle, the kor one boasts an elliptical shape thats inspired by the organic beauty of blown glass, with a sparkling body that stands out in a crowd. more importantly, the bottle eagerly quenches your thirst thanks to … | title: nathan speed 4 waist pack with four 10ounce nutrition flasks; categories: sports outdoors, exercise fitness, running, waist packs; brand: na; description: four 10 oz. nutrition flasks plus a versatile rear pocket make the speed 4 the ultimate mobile aid station.; price: 68.51; salesrank: sports amp; outdoors: 267505 |

| Initial ID: | kor one bpa free 750 | speed 4 waist pack nathan speed |
| Fine-tuned ID: | kor one hydration vessel polycarbon | speed 4 waist pack exercise fitness running |

**Dataset: Toys**

| Plain Text: | title: star wars r2d2 interactive astromech droid; categories: toys games, action figures statues, action figures; brand: na; description: collectors young and old will appreciate the details of this star wars interactive electronic r2d2 astromech droid. complete with movieaccurate messages and flashing lights, this droid responds to voice commands and has a special arm designed to keep your beverage handy. this friendly robot is designed to be a fun companion for kids aged eight and up..caption fontfamily: verdana, helvetica neue, arial, serif; fontsize: 10px; fontweight: bold; fontstyle: italic; ul.indent liststyle: inside disc; textindent: 15px; table.callout fontfamily: verdana; fontsize: 11px; lineheight: 1.3em … | title: chixos design a luxury loft home set; categories: toys games, dolls accessories, playsets; brand: na; description: the chixos house is the ultimate party palace. with 4 rooms to explore, theres so much to do.; price: 11.99; salesrank: toys games: 196539 |

| Initial ID: | star wars r2d2 interactive | chixos design a luxury loft |
| Fine-tuned ID: | star wars interactive electronic r2d2 | luxury loft home set toys games dolls |

**Figure 3: Quality study of the generated item IDs: two examples for each dataset, one with lengthy plain text data and one with shorter plain text data. The blue IDs are generated by the initial ID generator pre-trained on article tag generation, while the green IDs are generated by the ID generator after alternate training.**

*3.5.2* **Results**. The results are presented in Table 7. In this zero-shot evaluation, our model generally outperforms UniSRec, with the exception of HR@10 on the Music dataset, where UniSRec demonstrates better performance. Notably, for the Yelp dataset, which is cross-platform from the training data, our model significantly surpasses the baseline with a 353.46% improvement. This demonstrates the superior generalizability of our proposed model, making it more suitable for serving as a foundational model backbone. This

**Table 6: Recommendation accuracy for using 1) only user IDs, 2) only item IDs, and 3) both user and item IDs.**

| | | User ID | Item ID | User & Item ID |
|---|---|---|---|---|
| Sports | HR@5 | 0.0177 | 0.0404 | **0.0429** |
| | NDCG@5 | 0.0118 | 0.0308 | **0.0326** |
| | HR@10 | 0.0300 | 0.0528 | **0.0574** |
| | NDCG@10 | 0.0141 | 0.0348 | **0.0372** |
| Beauty | HR@5 | 0.0202 | 0.0577 | **0.0618** |
| | NDCG@5 | 0.0325 | 0.0441 | **0.0486** |
| | HR@10 | 0.0138 | 0.0778 | **0.0814** |
| | NDCG@10 | 0.0177 | 0.0506 | **0.0541** |

**Table 7: Zero-shot Evaluation. Intra-platform datasets are from the same platform of the training data (i.e., Amazon) and inter-platform datasets are from an unseen recommendation platform (i.e., Yelp).**

| Scenario | Dataset | Metric | UniSRec | **IDGenRec** |
|---|---|---|---|---|
| Intra-platform | Sports | HR@5 | 0.0060 | **0.0156** |
| | | NDCG@5 | 0.0034 | **0.0134** |
| | | HR@10 | 0.0098 | **0.0218** |
| | | NDCG@10 | 0.0046 | **0.0154** |
| | Beauty | HR@5 | 0.0118 | **0.0174** |
| | | NDCG@5 | 0.0068 | **0.0135** |
| | | HR@10 | 0.0206 | **0.0310** |
| | | NDCG@10 | 0.0096 | **0.0177** |
| | Toys | HR@5 | 0.0097 | **0.0103** |
| | | NDCG@5 | 0.0055 | **0.0079** |
| | | HR@10 | 0.0175 | **0.0215** |
| | | NDCG@10 | 0.0080 | **0.0114** |
| | Music | HR@5 | 0.0152 | **0.0184** |
| | | NDCG@5 | 0.0087 | **0.0148** |
| | | HR@10 | **0.0294** | 0.0238 |
| | | NDCG@10 | 0.0133 | **0.0165** |
| | Instruments | HR@5 | 0.0154 | **0.0203** |
| | | NDCG@5 | 0.0084 | **0.0139** |
| | | HR@10 | 0.0280 | **0.0440** |
| | | NDCG@10 | 0.0125 | **0.0215** |
| Inter-platform | Yelp | HR@5 | 0.0064 | **0.0300** |
| | | NDCG@5 | 0.0051 | **0.0248** |
| | | HR@10 | 0.0081 | **0.0329** |
| | | NDCG@10 | 0.0057 | **0.0258** |

is indeed expected, as our method can automatically create representative IDs for the items without requiring any special data processing.

Moreover, by comparing with the average HR and NDCG scores across the shared datasets presented in the standard evaluations, the foundational model's zero-shot recommendation capability is surprisingly comparable to, or even surpasses, that of some traditional recommendation models based on supervised training. For example, the zero-shot performance surpasses that of GRU4Rec on all four datasets, Bert4Rec on three out of four datasets, and Caser on two out of four datasets. Impressively, on the Yelp dataset, our model's zero-shot performance outperforms all traditional methods based on supervised training, only falling short of the P5 model. As the model's zero-shot recommendation capability may further improve with larger and more carefully collected training data, this suggests great potential for its future use as a foundational model.

## 4 RELATED WORKS

Recent LLM-based recommender systems can be broadly classified into two categories: discriminative and generative [36].

In LLM-based discriminative recommendation, LLMs are primarily used to learn better representations of users and items by leveraging contextual information in the recommendation process [12, 18, 20, 25, 27, 35, 37, 39, 40, 42]. Among them, BERT [6] is commonly used as the backbone. Compared to recommender systems that learn user/item embeddings primarily from user-item associations, these models capture rich textual information in conjunction with collaborative filtering. The core idea of these models is the learning of embeddings with LLMs and integrating the embeddings into a ranking score calculation based paradigm. Since this types of works are not the main focus of this paper, readers may refer to [36] for a more comprehensive study.

Our work belongs to the second category: LLM-based generative recommendation [3, 8, 9, 14]. This novel approach transitions from the traditional ranking-based recommendation to a pure text-to-text method. In comparison with the original score computation and re-ranking framework, this method can directly generate the target item from the complete pool of items [19], thus gathering significant attention. Geng et al. [8] introduced one of the first text-to-text generative recommendations, leveraging a pre-trained T5 as the recommendation backbone. The model was fine-tuned on 5 diverse recommendation tasks with predefined prompts and demonstrated promising results on the downstream tasks. Following [8], [9] proposed a multimodal version with a similar architecture, considering item visual features during recommendation generation. In these works, numerical IDs were assigned to each item, enabling the recommended items to be generated token-by-token. Subsequently, [14] posited that indexing item IDs is the most pivotal aspect of generative recommendation. They conducted an extensive study on indexing methods and assessed their effects on sequential recommendation performance. Although these methods have shown potential, a limitation exists: The rich textual information in recommender platforms is not fully harnessed by using numerical IDs. Furthermore, the employment of OOV tokens for item representation limits the generalizability of these methods, confining the trained model to a single dataset. Alternative approaches, such as using item titles [15] or categories [14] as semantic IDs, may initially seem to offer a more meaningful representation. However, these approaches may encounter duplicated ID issues as the number of items increases, and the IDs may have unintended overlaps. For example, similar titles may represent completely different items. Our method distinguishes itself by proposing an ID-generator that derives semantic item IDs from rich contextual information, enabling both the utilization of this information and facilitating cross-platform, zero-shot recommendations.

Another branch of generative recommendation research has concentrated on probing the potential of LLMs to directly generate recommendations without the need for training or fine-tuning [4, 7, 13, 21, 22, 30]. The primary focus in these works is the design of prompts. The emergence of this line of research was inspired by the exceptional zero-shot capabilities of ChatGPT, and they aim to discover how ChatGPT performs in a recommendation scenario.

However, as pointed out by [21], ChatGPT cannot generate recommendations with accuracy competitive to that of traditional recommendation methods.

## 5 CONCLUSION

In this paper, we address the item encoding problem in generative recommendation systems by introducing a novel framework that incorporates textual ID learning. Specifically, we employ an ID generator to produce unique, concise, and semantically rich textual IDs that are platform-agnostic and are based on human vocabulary. We also propose a diverse ID generation algorithm and an alternative training strategy to better align the LLM-based ID generator and the base recommender. This model has been proven to outperform existing recommendation baselines in the standard sequential recommendation setting. Furthermore, by training our model on some datasets while testing on other unseen datasets, our model shows strong performance under zero-shot recommendation scenario. Our research offers a new perspective to better align large language models and recommender systems by bridging the two through meticulously learned textual IDs, which may serve as a solid basis for training foundational recommendation models in the future.

## REFERENCES

[1] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311* (2022).

[2] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416* (2022).

[3] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. M6-rec: Generative pretrained language models are open-ended recommender systems. *arXiv preprint arXiv:2205.08084* (2022).

[4] Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. 2023. Uncovering ChatGPT's Capabilities in Recommender Systems. *arXiv preprint arXiv:2305.02182* (2023).

[5] Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2020. Autoregressive entity retrieval. *arXiv preprint arXiv:2010.00904* (2020).

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[7] Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. 2023. Chat-rec: Towards interactive and explainable llms-augmented recommender system. *arXiv preprint arXiv:2303.14524* (2023).

[8] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 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.

[9] Shijie Geng, Juntao Tan, Shuchang Liu, Zuohui Fu, and Yongfeng Zhang. 2023. VIP5: Towards Multimodal Foundation Models for Recommendation. *arXiv preprint arXiv:2305.14302* (2023).

[10] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*. 507–517.

[11] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).

[12] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. 2022. Towards universal sequence representation learning for recommender systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 585–593.

[13] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2023. Large language models are zero-shot rankers for recommender systems. *arXiv preprint arXiv:2305.08845* (2023).

[14] Wenyue Hua, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. 2023. How to Index Item IDs for Recommendation Foundation Models. *SIGIR-AP* (2023).

[15] Jianchao Ji, Zelong Li, Shuyuan Xu, Wenyue Hua, Yingqiang Ge, Juntao Tan, and Yongfeng Zhang. 2024. Genrec: Large language model for generative recommendation. In *European Conference on Information Retrieval*. Springer, 494–502.

[16] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.

[17] Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226* (2018).

[18] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian McAuley. 2023. Text Is All You Need: Learning Language Representations for Sequential Recommendation. *arXiv preprint arXiv:2305.13731* (2023).

[19] Lei Li, Yongfeng Zhang, Dugang Liu, and Li Chen. 2023. Large Language Models for Generative Recommendation: A Survey and Visionary Discussions. *arXiv preprint arXiv:2309.01157* (2023).

[20] Ruyu Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. 2023. Exploring the upper limits of text-based collaborative filtering using large language models: Discoveries and insights. *arXiv preprint arXiv:2305.11700* (2023).

[21] Junling Liu, Chao Liu, Renjie Lv, Kang Zhou, and Yan Zhang. 2023. Is chatgpt a good recommender? a preliminary study. *arXiv preprint arXiv:2304.10149* (2023).

[22] Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. 2023. A First Look at LLM-Powered Generative News Recommendation. *arXiv preprint arXiv:2305.06566* (2023).

[23] Chen Ma, Peng Kang, and Xue Liu. 2019. Hierarchical gating networks for sequential recommendation. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 825–833.

[24] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*. 43–52.

[25] Aashiq Muhamed, Iman Keivanloo, Sujan Perera, James Mracek, Yi Xu, Qingjun Cui, Santosh Rajagopalan, Belinda Zeng, and Trishul Chilimbi. 2021. CTR-BERT: Cost-effective knowledge distillation for billion-parameter teacher models. In *NeurIPS Efficient Natural Language and Speech Processing Workshop*.

[26] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

[27] Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. 2021. U-BERT: Pre-training user representations for improved recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 4320–4327.

[28] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.

[29] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1441–1450.

[30] Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agent. *arXiv preprint arXiv:2304.09542* (2023).

[31] Jiaxi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the eleventh ACM international conference on web search and data mining*. 565–573.

[32] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[33] Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424* (2016).

[34] Ronald J Williams and David Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural computation* 1, 2 (1989), 270–280.

[35] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2021. Empowering news recommendation with pre-trained language models. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1652–1656.

[36] Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2023. A Survey on Large Language Models for Recommendation. *arXiv preprint arXiv:2305.19860* (2023).

[37] Shitao Xiao, Zheng Liu, Yingxia Shao, Tao Di, Bhuvan Middha, Fangzhao Wu, and Xing Xie. 2022. Training large-scale news recommenders with pretrained language models in the loop. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4215–4225.

[38] Shuyuan Xu, Wenyue Hua, and Yongfeng Zhang. 2024. OpenP5: An Open-Source Platform for Developing, Training, and Evaluating LLM-based Recommender Systems. *SIGIR* (2024).

[39] Shaowei Yao, Jiwei Tan, Xi Chen, Juhao Zhang, Xiaoyi Zeng, and Keping Yang. 2022. ReprBERT: Distilling BERT to an Efficient Representation-Based Relevance Model for E-Commerce. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4363–4371.

[40] Zheng Yuan, Fajie Yuan, Yu Song, Youhua Li, Junchen Fu, Fei Yang, Yunzhu Pan, and Yongxin Ni. 2023. Where to go next for recommender systems? id-vs. modality-based recommender models revisited. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2639–2649.

[41] Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. 2023. Recommendation as instruction following: A large language model empowered recommendation approach. *arXiv preprint arXiv:2305.07001* (2023).

[42] Song Zhang, Nan Zheng, and Danli Wang. 2022. GBERT: Pre-training User representations for Ephemeral Group Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2631–2639.

[43] Tingting Zhang, Pengpeng Zhao, Yanchi Liu, Victor S Sheng, Jiajie Xu, Deqing Wang, Guanfeng Liu, Xiaofang Zhou, et al. 2019. Feature-level Deeper Self-Attention Network for Sequential Recommendation.. In *IJCAI*. 4320–4326.

[44] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 1893–1902.