# LEADRE: Multi-Faceted Knowledge Enhanced LLM Empowered Display Advertisement Recommender System

Fengxin Li
Renmin University of China
lifengxin@ruc.edu.cn

Yi Li
Yue Liu
Chao Zhou
Yuan Wang
sincereli@tencent.com
herculesliu@tencent.com
derekczhou@tencent.com
leoyuanwang@tencent.com
Tencent Inc.

Xiaoxiang Deng
Wei Xue
Dapeng Liu
reesedeng@tencent.com
weixue@tencent.com
rocliu@tencent.com
Tencent Inc.

Lei Xiao
Haijie Gu
Jie Jiang
shawnxiao@tencent.com
jerrickgu@tencent.com
zeus@tencent.com
Tencent Inc.

Hongyan Liu
Tsinghua University
hyliu@tsinghua.edu.cn

Biao Qin
Jun He
qinbiao@ruc.edu.cn
hejun@ruc.edu.cn
Renmin University of China

## Abstract

Display advertising provides significant value to advertisers, publishers, and users. Traditional display advertising systems utilize a multi-stage architecture consisting of retrieval, coarse ranking, and final ranking. However, conventional retrieval methods rely on ID-based learning to rank mechanisms and fail to adequately utilize the content information of ads, which hampers their ability to provide diverse recommendation lists.

To address this limitation, we propose leveraging the extensive world knowledge of LLMs. However, three key challenges arise when attempting to maximize the effectiveness of LLMs: "*How to capture user interests*", "*How to bridge the knowledge gap between LLMs and advertising system*", and "*How to efficiently deploy LLMs*". To overcome these challenges, we introduce a novel LLM-based framework called **LLM E**mpowered Display **AD**vertisement **RE**commender system (LEADRE). LEADRE consists of three core modules: (1) The **Intent-Aware Prompt Engineering** introduces multi-faceted knowledge and designs intent-aware *<Prompt, Response>* pairs that fine-tune LLMs to generate ads tailored to users' personal interests. (2) The **Advertising-Specific Knowledge Alignment** incorporates auxiliary fine-tuning tasks and Direct Preference Optimization (DPO) to align LLMs with ad semantic and business value. (3) The **Efficient System Deployment** deploys LEADRE in an online environment by integrating both latency-tolerant and latency-sensitive service. Extensive offline experiments demonstrate the effectiveness of LEADRE and validate the contributions of individual modules. Online A/B test shows that LEADRE leads to a 1.57% and 1.17% GMV lift for serviced users on WeChat Channels and Moments separately. LEADRE has been deployed on both platforms, serving tens of billions of requests each day.

## CCS Concepts

• **Information systems** → **Display advertising**; *Personalization*.

## Keywords

Display Advertising, LLM-based Recommendation, Generative Retrieval

## 1 Introduction

Online display advertising provides significant value to advertisers, publishers, and users by facilitating targeted content delivery and meeting users' personal interests [36, 37]. Traditionally, display advertising systems utilize a multi-stage architecture, including retrieval, coarse ranking, ranking, re-ranking, etc. The retrieval stage, serving as the entry of the process, is crucial for identifying user interests and overcoming the "information cocoon" by presenting a diverse set of ad options[5, 24].

Traditional retrieval solutions typically employ ID-based learning-to-rank mechanisms to learn collaborative semantics for efficient ad filtering. However, these methods fail to adequately utilize the content information of ads, bearing limitations in generating diverse recommendations—especially when dealing with sparse user behaviors and long-tail ads [10, 43].

Recently, large language models (LLMs) have demonstrated exceptional capabilities in understanding, generalization, and reasoning by leveraging vast amounts of general knowledge [25, 27, 53]. To address the shortcomings of conventional retrieval approaches, researchers have begun incorporating LLMs into recommender systems [33, 42]. However, these methods often struggle to find applications in industrial display advertising systems, which introduces several critical challenges:

**(1) How to capture user interests in scenarios with implicit user intents and sparse user behaviors:** While LLMs have demonstrated strong capabilities in understanding user intent, display advertising often operates in an environment with implicit user intents, where explicit queries are absent. Additionally, user behavior in the ad domain tends to be sparse, making it crucial to effectively utilize the limited available data and supplement it appropriately to capture commercial intent. **(2) How to bridge the knowledge gap within LLMs and advertising system:** Although LLMs excel at generating natural language responses based on general knowledge, display advertising systems require the generation of specific ads from a predefined set of ads. Bridging the gap between the general capabilities of LLMs and the specific needs of the advertising domain is essential. Moreover, the generated ads must align with business objectives, necessitating that LLMs are tuned to meet the specific goals of the advertising system. **(3) How to efficiently deploy LLMs in industrial display advertising system which serves billions of users requests per day:** Deploying LLMs within an online advertising system introduces significant computational overhead, which can conflict with the need to maintain cost efficiency. Therefore, it is crucial to deploy LLM-based methods in a way that balances computational costs with advertising performance, ensuring scalability and efficiency within the system.

To address these challenges, we propose a novel multi-faceted knowledge enhanced **LLM** **E**mpowered display **AD**vertisement **RE**commender system (LEADRE), termed **LEADRE**. The framework consists of three key modules:

**(1) Intent-Aware Prompt Engineering:** This module designs intent-aware *<Prompt, Response>* pairs that fine-tune LLMs to generate ads tailored to users' personal interests. To overcome the scarcity of user behaviors in the ad domain, this module introduces user behaviors from content domains such as micro-videos, and news. Additionally, it models user commercial intent by combining **long-term interests** (derived from user profiles and historical behaviors) and **short-term interests** (derived from recent behaviors) within the prompt. **(2) Advertising-Specific Knowledge Alignment:** This module incorporates auxiliary fine-tuning tasks and Direct Preference Optimization (DPO). The auxiliary fine-tuning tasks are specifically designed to bridge the semantic gap between language modality and advertising system. Additionally, to balance user intent and business value, we apply **DPO**, which encourages the model to generate ads with higher business value. **(3) Efficient**

**System Deployment:** This module integrates both latency tolerant service and latency sensitive service into the LLM deployment architecture. To further improve computational efficiency, we optimize the deployment process using **TensorRT LLM Acceleration**.

We implement LEADRE using Hunyuan with 1B parameters and conduct extensive offline and online experiments within Tencent's display advertising system. The offline results validate the contributions of the individual modules, demonstrating their effectiveness in the overall framework. In the online A/B test, LEADRE achieved a 1.57% increase in Gross Merchandise Value (GMV) on Tencent WeChat Channels for serviced users, along with a 1.17% increase in GMV on Tencent WeChat Moments for serviced users. Currently, LEADRE has been successfully deployed on both platforms, serving billions of users and processing tens of billions of user requests each day. The retrieved ads are further incorporated into the ranking phase as new features. Specifically, on the user side, the retrieved ads are used to extended user interests, while on the item side, the match scores between the retrieved ads and the target ad are introduced as new item-level features. These new features contribute to an extra 1.43% improvement in GMV on Tencent WeChat Channels.

The contributions of this work can be summarized as follows:

(1) To the best of our knowledge, this is the first work to deploy LLMs in an online display advertising system. We propose a novel LLM-based generative retrieval framework, LEADRE, for display advertising, and deploy it using a combination of latency tolerant service system and latency sensitive service system to optimize performance and resource utilization.

(2) LEADRE integrates an Intent-Aware Prompt Engineering and an Advertising-Specific Knowledge Alignment to ensure the generated ads are accurate, diverse, aligned with high business value, and meet user interests.

(3) Extensive experiments, both offline and online, demonstrate the effectiveness of the proposed method. We observe significant improvements in offline metrics such as HitRatio, as well as in online metrics like GMV.

## 2 Related Work

### 2.1 Large Language Models-based Recommender

Large Language Models (LLMs) have demonstrated remarkable capabilities in both understanding and generation tasks[18, 49], and they have been widely applied to a variety of domains, including document summarization[22, 55], conversational agents[15, 40], code completion[13, 54], and others[51, 57]. Recently, researchers have begun exploring the application of LLMs in recommender systems[46]. These efforts can be broadly categorized into two roles for LLMs: *feature encoder* and *ranker*. As feature encoders, LLMs are leveraged primarily for their ability to understand and represent textual information. In these works, user or item features are often assembled into textual descriptions via pre-designed templates, which are then encoded with LLMs to generate feature embeddings [1, 35, 42]. These embeddings can subsequently be utilized in downstream tasks, such as ranking items in a recommender system. On the other hand, as rankers, LLMs are employed for their ability to generate. In these works, user behavior sequences are transformed into text sequences, and LLMs are tasked with

generating the next likely item (or user action) based on the input text [14, 52]. To this end, LLMs typically require fine-tuning on recommendation-specific tasks to acquire the domain knowledge and capture user preference in recommender systems [2, 14, 26].

However, these works fail to find applications in industrial display advertising systems due to implicit user intents, high computational costs, etc. In this work, we present the first application of LLM-based generative retrieval in an industrial display advertising system. We employ intent-aware prompt engineering and ad-specific knowledge alignment to enable LLMs to generate ads that are diverse, aligned with business objectives, and tailored to user interests. Additionally, we address computational cost challenges by deploying LLMs through a hybrid architecture that integrates both latency-tolerant and latency-sensitive services.

## 2.2 Retrieval in Advertising Systems

Display advertising systems typically follow a multi-stage architecture composed of multiple stages, including retrieval, corse-ranking, ranking, re-ranking, and others[8, 12]. The retrieval stage, positioned at the top of the funnel, plays a critical role in identifying user interests and generating a pool of candidate ads for further processing [21, 47]. Most advertising systems adopt ID-based learning to ranking mechanisms to recall all ads that users may find appealing. These methods often rely on a two-tower model, where one tower encodes user features and the other encodes item (ad) features[3, 19].

Despite their efficiency, ID-based retrieval methods can suffer from issues such as the "information cocoon" effect, limited diversity in the retrieved list, and the inability to effectively leverage the ads' content information[20, 44]. In this work, we propose incorporating LLMs to address these issues. LLMs, with their powerful generative and contextual understanding capabilities, can enhance the diversity of the retrieved list by expanding the range of ads presented to users, ultimately improving the overall effectiveness of the advertising system.

## 3 Preliminary

### 3.1 Problem Definition

Let $\mathcal{U}$ and $\mathcal{A}$ denote the sets of users and advertisements(ads), respectively, within the target advertising system. Due to the limited user interactions with the ads, we introduce a content domain to enrich user behavior data. Let $C$ denote the set of content items in the system. The sequence of user behaviors is represented as $\mathcal{S}_u = [i_u^1, i_u^2, \cdots, i_u^L]$ in chronological order, where $u \in \mathcal{U}$, $i_u \in \mathcal{A} \cup C$, and $L = |\mathcal{S}_u|$ is the length of the user's behavior sequence. Given this user behavior sequence, the goal of generative retrieval is to generate the next relevant ad, $a_u^{L+1} \in \mathcal{A}$.

### 3.2 Ads Indexing

Traditional advertising systems typically index ads incrementally and learn dense embeddings for each ad [4, 45]. However, the index built in this way usually lacks semantic information, creating a significant gap between the advertising system and human understanding about ads. To bridge this gap, we adopt the strategy of learning Semantic IDs (S-IDs) based on ad features (illustrated in Figure 1), inspired by previous works on item indexing [11, 33, 56].
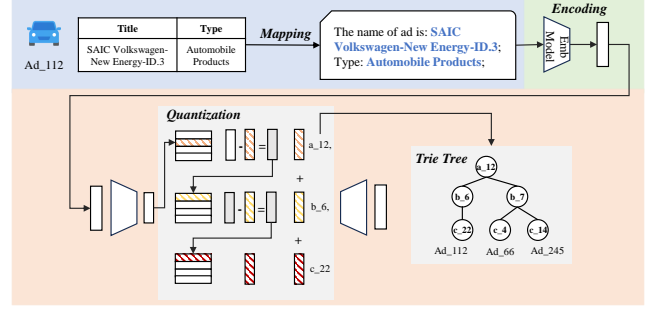


**Figure 1: Overall Framework of Ads indexing in LEADRE.**

Specifically, we first define a **feature mapping rule** in the form of a template string to generate text description of each ad, as follows:

> *The name of the ad is <ads_name>; The product type is <ads_type>; The first-level category is <first_cat>; The second-level category is <second_cate>; The attributes include: <basic_att>.*

For each ad $a \in \mathcal{A}$, all its features are filled in above template string to generate a textual description $t_a$. For example:

> *The name of the ad is **SAIC Volkswagen-New Energy-ID.3**; The product type is **Automobile Products**; The first-level category is **Automobile**; The second-level category is **SAIC Volkswagen·New Energy**; The attributes include: **automobile brand_Volkswagen, automobile series_Volkswagen ID.3**.*

After obtaining the textual description $t_a$, we employ a pre-trained language model (e.g. Hunyuan [38], E5[41]) to **encode** $t_a$, bridging the gap between the advertising system and human understandable text. This step can be formulated as $\mathbf{x}_a = f_{PLM}(t_a; \theta_{PLM})$, where $f_{PLM}$ represents the pre-trained language model parameterized by $\theta_{PLM}$. The resulting $\mathbf{x}_a \in \mathbb{R}^{d_{PLM}}$ is the textual embedding of ad $a$, and $d_{PLM}$ denotes the output dimension of the pre-trained language model.

After that, we employ a Residual-Quantized Variational AutoEncoder (RQ-VAE) [23] to generate concise semantic IDs for the ads. The RQ-VAE takes the ad's embedding $\mathbf{x}_a$ as input and uses an **encoder-quantization-decoder** mechanism to generate a list of discrete semantic tokens. For details of RQ-VAE training please refer to Appendix A.1

### 3.3 Trie-Tree Construction

In the process of ad indexing, each ad $a$ is represented by a list of semantic tokens denoted as $\hat{\mathbf{c}}_a = [c_a^1, c_a^2, \ldots, c_a^M]$. To handle potential collisions where different ads may map to the same S-IDs, we introduce an additional code to distinguish them [33]. As a result, the S-IDs for each ad are updated to $\mathbf{c}_a = [c_a^1, c_a^2, \ldots, c_a^L, c_a^{(M+1)}]$, where $c_a^{(M+1)}$ represents the additional code used to disambiguate ads with identical list of S-IDs. To distinguish different level semantic tokens, we use <a_ , b_, c_, $\cdots$ > as prefix of different S-IDs level.
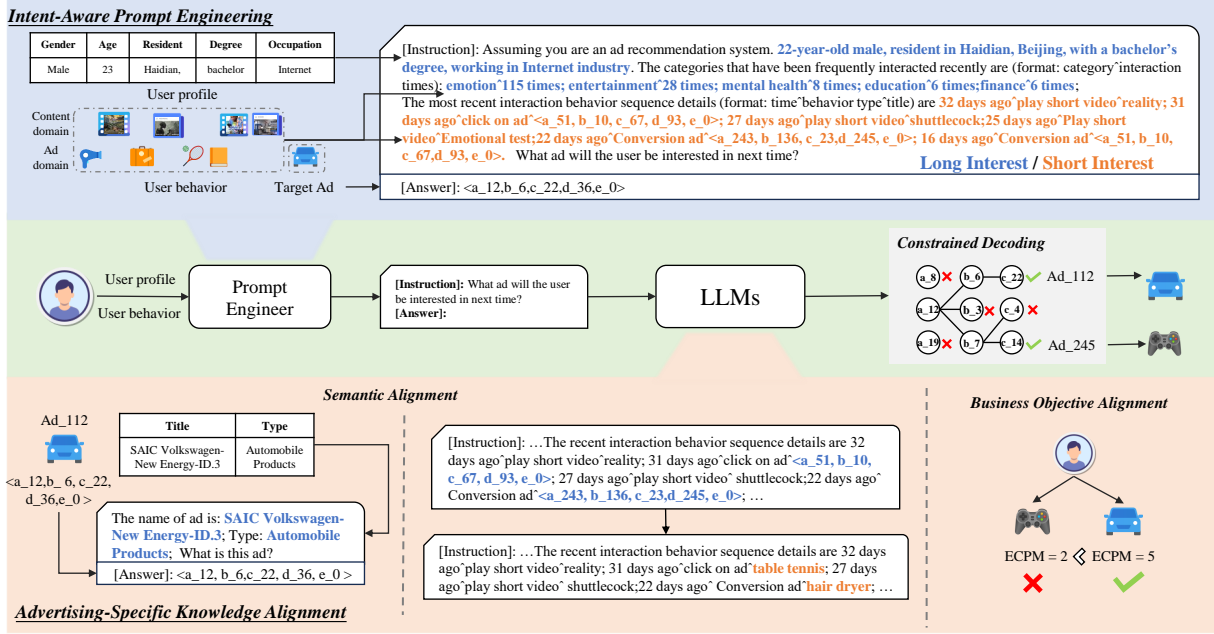
**Figure 2: Overall Framework of Ads Constrained Generation Module and LLM Fine-tuning Module.**

Next, we represent all ads within a trie-tree, where each path from the root node to a leaf node corresponds to the S-IDs of a specific ad. The **Trie-Tree Construction Algorithm** (detailed in Appendix A) iteratively adds each semantic token of an ad's S-IDs as a child node, creates new nodes when necessary and marks the end of each complete ad's S-IDs. This hierarchical structure enables efficient organization of ads, allowing for quick retrieval and prefix-based searches based on S-IDs.

## 3.4 Constrained Decoding

Different from conventional LLM-based chatting bot, generative retrieval for advertising system requires that the generated ads are valid ads in the predefined ads set, and should be a list of diverse ads instead of only one. To meet these requirements, we employ constrained decoding with the trie-tree [16, 17, 28]. The **Constrained Decoding Algorithm**, detailed in Appendix A, iteratively expands candidate S-IDs lists through layers of the trie-tree, selecting the top $B$ lists at each step based on their scores, until either all layers are processed or the maximum length is reached. This approach ensures that the final generated ads are both high-probability according to the LLMs and valid within the predefined ad set.

## 4 Method

In this section, we present a detailed explanation of multi-faceted knowledge enhanced **L**LM **E**mpowered display **AD**vertisement **RE**commender system (LEADRE), which integrates large language models (LLMs) into our display advertising system.

## 4.1 Overview

LEADRE framework consists of three novel modules: **(1) Intent-Aware Prompt Engineering** designs intent-aware *<Prompt, Response>* pairs which accommodate user profile and behavior sequences, and are used as fine-tuning corpus for LLMs to learn users' interests. **(2) Advertising-Specific Knowledge Alignment** incorporates auxiliary fine-tuning tasks and Direct Preference Optimization (DPO) to bridge the knowledge gap within LLMs and advertising system. **(3) Efficient System Deployment** integrates both latency tolerant service system and latency sensitive service system into LLMs deployment, ensuring scalability and real-time performance.

## 4.2 Intent-Aware Prompt Engineering

To accurately capture user intent and predict the next relevant ad, the LLM should be fine-tuned with proper corpus first. We build textual *<Prompt, Response>* pairs with user behavior sequences and ad descriptions as the corpus. Unlike search advertising [33], where user queries explicitly express intent and could be used as the prompt, display advertising operates in a low-intent environment where explicit queries are absent and user behavior is often sparse. As a result, the prompt must effectively leverage available user data within the advertising system and supplement it with user data outside of the advertising system, such as user behaviors in the content domain, to adequately capture commercial intent. Considering the affordable prompt length is quite limited, user data should be encoded concisely enough. Below, we detail the prompt engineering process.

*4.2.1 Prompt Components.* Given their extensive training on vast amounts of web-crawled corpus, LLMs are capable of understanding

and operating with content from various domains, as long as they are encoded as human understandable text in natural languages. We model user commercial intent by combining both **long-term interests** (derived from the user's profile and historical behaviors) and **short-term interests** (based on recent behaviors). The prompt is assembled with the following key components:

**(1) Task Instruction**: This part describes the task for the LLM to perform. Furthermore, it provides a learnable token for all fine-tuning tasks, similar to Soft Prompt formats [29, 39], helping the LLM recognize instructions related to ad generation and retrieval. The template of the task description is shown below:

> *The following is an instruction describing a task. Please give a response to complete this request appropriately.*
> *[Instruction]:<**Learnable Token**>Assuming you are an ad recommender system, <**Prompt**>, what ad will the user be interested in next time?*

**(2) User Profiles**: This part provides demographic information about the user, such as age, gender, and region. The template is shown below:

> *<**age**><**sex**>, resident in <**residence**>, with a <**education_level**> degree, working in <**occupation**>, with a <**consumption_level**>*

**(3) User Interest Summary**: This part summarizes the user's long-term positive and negative feedbacks on commercial interest categories and standard product units (SPUs) across the ad and content domains. The template is shown below:

> *The categories that have been frequently interacted recently are (format: category^interaction times): <**category_1**>^<**count_1**> times; <**category_2**>^<**count_2**> times; <**category_3**>^<**count_3**> times;*

**(4)User Ad Domain Behavior Sequence**: This part expresses the user's short-term behaviors in the ad domain. To filter the noise signals in the sequence, only positive user actions like clicks and conversions are considered. As discussed in Section 3.1, each ad is represented by a list of semantic tokens denoted as Semantic-IDs (S-IDs). To distinguish different level semantic tokens, we use <a_ , b_, c_, · · · > as prefix of different S-IDs level.

**(5) User Content Domain Behavior Sequence**: This part captures the user's short-term behaviors in the content domain. Each interaction is described in terms of the commercial category and SPU. Similarly, only positive actions like watching full videos and searching are considered. The ad and content domain behaviors are merged and presented in chronological order. The template is shown below:

> *The most recent interaction behavior sequence details (format: time^behavior type^title) are <**time_1**> days ago^<**type_1**>^ <**title_1**>/<**SIDs_1**>; <**time_2**> days ago^<**type_2**>^<**title_2**>/<**SIDs_2**>; <**time_3**> days ago ^ <**type_3**> ^ <**title_3**>/<**SIDs_3**>*

By concatenating the previous components, we obtain a complete prompt. An example prompt is shown below:

> *The following is an instruction describing a task. Please give a response to complete this request appropriately.*

> *[Instruction]:<**Learnable Token**>Assuming you are an ad recommender system. **22-year-old male**, resident in **Haidian, Beijing**, with a **bachelor**'s degree, working in **Internet industry**, with a **medium consumption level**. The categories that have been frequently interacted recently are (format: category^interaction times): **emotion^115 times; entertainment^28 times; mental health^8 times; education^6 times; finance^6 times**; The most recent interaction behavior sequence details (format: time^behavior type^title) are **32 days ago^play short video^reality; 31 days ago^click on ad^<a_51, b_10, c_67, d_93, e_0>; 27 days ago^play short video^shuttlecock; 25 days ago^Play short video^Emotional/psychological age test;22 days ago^Conversion ad^<a_243, b_136, c_23, d_245, e_0>;19 days ago^Click ad^<a_164, b_243, c_38, d_88, e_0>;16 days ago^Conversion ad^<a_51, b_10, c_67, d_93, e_0>**. what ad will the user be interested in next time?*

For the *Response*, the last click or conversion ad in the user behavior sequence is used as a supervisory signal for fine-tuning the LLM. A sample response is shown below:

> *[Response]: <**a_122, b_28, c_35, d_15, e_0**>*

*4.2.2 Prompt Augmentation.* We develop a series of prompt augmentation strategies to enhance the fine-tuning process. These strategies are outlined below:

**Multiple Prompt Templates:** We introduce a variety of prompt templates to map a single user behavior sequence into multiple prompts. By altering the arrangement of sequence components and varying the instruction descriptions, these templates facilitate a broader understanding of user intentions. This diversification allows the model to capture different aspects of user intent, enriching the generation process.

**User Profile Reordering:** To prevent LLMs from simply memorizing specific token sequences or orders, we propose a user profile reordering method. It reshuffles user profile descriptions, encouraging the model to focus on the semantic meaning rather than the fixed order of inputs.

**Ad Positive Interaction Reuse:** Inspired by auto-regressive mechanisms, We augment the original user behavior sequence into multiple prompt samples by reusing positive interaction, enhancing the model's ability to generalize from limited interaction data. For example, Assuming that user behavior sequence is denoted by $S_u = [c_u^1, a_u^2, a_u^3, c_u^4, a_u^5]$, where $c_u^1, c_u^4 \in C$ denote content domain items and $a_u^2, a_u^3, a_u^5 \in \mathcal{A}$ denote ads. We augment the original sequence into prompt samples: $<[c_u^1, a_u^2, a_u^3, c_u^4], a_u^5>$, $<[c_u^1, a_u^2], a_u^3>$, $<[c_u^1], a_u^2>$.

## 4.3 Advertising-Specific Knowledge Alignment

The advertising system poses a great gap with LLMs, which hurts the ad generation ability of LLMs. To bridge this gap, we conduct semantic alignment by auxiliary tuning tasks to align the ad's Semantic-IDs (S-IDs) with the LLM, and business objectiveness alignment by Direct Preference Optimization (DPO) [31, 50] to encourage the generation of ads with higher Effective Cost Per Mille (ECPM).

*4.3.1 Semantic Alignment.* Apart from the main task, we introduce auxiliary tasks to align the ad system with the language model

and enhance the model's advertising knowledge[56]. The auxiliary tasks consist of two components: an explicit alignment task and an implicit alignment task. These tasks not only improve the language model's understanding of ads, but also contribute to better performance in the main task of the next ad generation.

**The explicit alignment task** is designed to predict the S-IDs of an ad based on its detailed description. This task directly aligns the language model with the ad system by making the LLMs understand how textual descriptions map to the corresponding ad S-IDs. A sample *prompt-response* pair for this task is shown below:

---

*[Prompt]: Given the ad's detailed description "**The name of the ad is SAIC Volkswagen-New Energy-ID.3; The product type is Automobile Products; The first-level category is Automobile; The second-level category is SAIC Volkswagen·New Energy; The attributes include: automobile brand_Volkswagen, automobile series_Volkswagen ID.3.**", what is the corresponding ad?*
*[Response]: <a_12, b_22, c_50, d_25, e_0>*

---

In this example, LLM is provided with a detailed textual description of an ad, including the product's name, type, brand, and category. The model's task is to map this description to the corresponding S-IDs, which are internal identifiers used by the ad system.

**The implicit alignment task** builds on the Next Ad Generation mechanism described in the main task but replaces the S-IDs in the user's ad interaction history with the ad descriptions. Instead of predicting the next ad based on a sequence of S-IDs (as in the main task), the model now predicts the next ad using textual descriptions of the ads. This subtle modification enhances the LLMs' ability to associate ad descriptions with user behaviors, further aligning the model with the advertising system in a more implicit manner. By incorporating these auxiliary tasks, the language model's advertising knowledge is enhanced, contributing to better performance in the next ad generation task.

**Task Order:** The tasks are conducted in the following order: explicit alignment task → implicit alignment task → main task. It is essential to prioritize alignment tasks, both explicit and implicit, before addressing the main task. These alignment tasks, particularly the explicit alignment task, serve as crucial preparatory steps that equip the LLMs to understand the mapping between textual ad descriptions and the S-IDs used by the ad system. By completing these tasks first, the model gains a deeper comprehension of the structural and semantic aspects of the ad domain, ensuring that it is properly aligned with the system's internal representations. This alignment process is not merely a supplementary step but a foundational one, directly influencing the model's effectiveness in the primary task. In essence, the explicit and implicit alignment tasks establish the groundwork necessary for the model's success in predicting the next ad, making them indispensable for optimal performance in the main task.

*4.3.2 Business Objectiveness Alignment.* The primary and auxiliary tasks assist the LLM in recognizing user intent and interests, thereby generating ads that users are more likely to click on or convert. However, due to the LLM's inherent lack of business context, the ads it generates may have limited business value, contributing
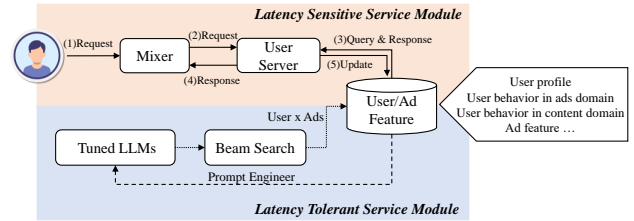


**Figure 3: Framework of Efficient System Deployment.**

minimally to cost and Gross Merchandise Volume (GMV). To address this, we introduce **Direct Preference Optimization (DPO)**, which encourages the generation of ads with higher business value.

For a specific user $u$, the goal of DPO is to increase the likelihood of generating ads with high business value, denoted as $a_h$, while simultaneously reducing the probability of generating ads with lower business value, denoted as $a_l$. To construct the chosen pair $\langle u, a_h \rangle$ (preferred ad) and the reject pair $\langle u, a_l \rangle$ (less preferred ad) for DPO, we leverage the Expected Cost Per Mille (ECPM) metric during the ranking phase.

Specifically, for a given user $u$, we select two potential ads, $a_1$ and $a_2$, that are likely to result in a click or conversion. We then calculate the ECPM for each ad. The ad with the higher ECPM is considered the high-value ad $a_h$, while the ad with the lower ECPM is considered the low-value ad $a_l$. Through DPO, the LLM enhances its ability to recognize business value, thereby motivating the model to balance user intent with business objectives during ad generation.

## 4.4 Efficient System Deployment

In this section, we detail the deployment framework and engineering techniques for LLM inference used to ensure efficient and responsive ad generation.

*4.4.1 Deployment Framework.* Typical LLMs exhibit high latency, which is inadequate for real-time display advertising services. To address this, we implement a hybrid system by integrating latency tolerant service module and latency sensitive service module [5, 24]. The deployment consists of the following components:

**(1) Inference:**

*Latency Sensitive Service Module :* This component is activated when a user makes a request. A Mixer receives the user request and queries the User Server for the ad retrieval list. The User Server retrieves the pre-computed retrieval list from the user feature database and sends it to the Mixer and further responds to the user. After that, the user behaviors update and ad feature update will be stored in the user feature database and ad feature database.

*Latency Tolerant Service Module:* This module adopts nearline computing for LLM inference and is triggered after a user request. This component receives the fine-tuned LLM and the Trie-Tree from the offline stage for inference. It inquiries users' recent behavior to construct prompt samples and generates a list of ads for the users. The generated ad list is stored in the user feature database and made available for online service.

Even for nearline inference, GPU resources are relatively scarce and cannot be allocated to every request. To address this, we introduce the **Adaptive Resource Distribution Strategy**. This strategy categorizes users into 25 user groups according to Average Revenue Per User (ARPU) value, prioritizing GPU inference for high-value users.

**(2) Training:** The training stage is responsible for LLM fine-tuning, and it is updated daily. The process begins with engineering samples from the user feature database to extract ad features to construct *<Prompt, Response>* pairs for LLM fine-tuning. The *<Prompt, Response>* pairs are utilized to fine-tune the LLM, enabling it to generate ads effectively. Once the training is complete, the fine-tuned LLM is deployed to the latency tolerant service module.

*4.4.2 TensorRT LLM Acceleration.* To further enhance the computational efficiency of nearline LLM inference, we employ TensorRT [1] for inference acceleration. This approach leverages several optimization techniques to significantly reduce latency and improve performance, particularly in advertising scenarios:

1) **TensorRT-LLM Kernel Optimization**

*Softmax Kernel Optimization:* The softmax kernel, which accounts for over 80% of decoding time, is optimized by using 'float4' for vectorized memory access, reducing bandwidth usage and doubling kernel performance. This reduces end-to-end time by 5%. Additionally, for the prefill step, we introduced a 'ComputeMode' to calculate softmax only for beam 0, improving prefill performance by 3-4 times.

*Finalize Kernel Optimization:* In beam search, the finalize kernel recursively obtains token IDs. Initially, it used a single thread, underutilizing GPU resources. By assigning a block per beam width, performance improved by 50 times, reducing latency by 10%.

2) **TensorRT-LLM Quantization** LLM inference generally involves two stages: prefill and generate. Prefill is typically the computational bottleneck, while the generate stage is more memory-bound. However, in advertising scenarios with a large beam width and batching functionality supported by the inference engine, the generate stage also becomes a computational bottleneck. Therefore, we adopt quantization schemes that focus on activation values, such as smooth quantization ('w8a8c8') and FP8 quantization ('w8a8c8'), to mitigate these bottlenecks.

*Smooth Quantization ('w8a8c8'):* We apply smooth quantization (int8 precision for weights and activations), doubling the computational power for operations like matrix multiplication. This results in a 50% performance improvement in search advertising.

*FP8 Quantization ('w8a8c8'):* FP8 quantization offers a simpler alternative to smooth quantization and achieves similar performance on H20 GPUs without sacrificing precision.

3) **Proxy Load Balancing for Multi-GPU Utilization:** Previous optimizations, such as proxy global load balancing, mainly focused on improving the performance of individual GPUs. To enhance the utilization of multiple GPUs, we optimized the scheduling strategy to better distribute workloads across the GPU cluster. In our system, multiple proxies interact with multiple TensorRT-LLM instances. Without coordination, multiple proxies may send requests to the same TensorRT-LLM instance, leading to uneven load distribution.

TensorRT-LLM is highly sensitive to queries per second (QPS). If one instance processes even a couple more requests than others, it can become the bottleneck for the entire cluster. To mitigate this, we introduced Redis-based global counters to evenly distribute requests across all TensorRT-LLM instances, ensuring balanced load distribution. As a result, the TensorRT-LLM cluster now achieves over 90% of its theoretical maximum load capacity.

## 5 Experiment

In this section, we verify the effectiveness of LEADRE by conducting extensive offline and online experiments.

### 5.1 Experiment Setup

*5.1.1 Dataset.* We conducted all experiments using industrial-scale datasets on display advertising of Tencent WeChat Channels, as public datasets were deemed unsuitable due to gaps in applicability to our serving system and significant discrepancies with our internal models. The experiment process is suitable for all industrial display advertising systems.

We created user behavior sequences in chronological order and applied the "leave-one-out" strategy for dataset splitting [4, 14, 45]. Specifically, the last ad interaction in each user's sequence was designated as the test set, while the remaining interactions were used for training. The maximum user behavior window was set to 90 days and the maximum prompt token length was set to 2096, covering user behavior in both the ads and content domains.

*5.1.2 Evaluation Metric.* We used the Hit Ratio (HR@K) and Normalized Discounted Cumulative Gain (NDCG@K) as the primary metrics to evaluate the offline predictive performance of the models. Hit Ratio measures whether the model successfully retrieves the clicked or converted ads for users. NDCG accounts for both the relevance of items and their positions in the ranked list.

*5.1.3 Implementation Details.* The LEADRE is implemented on the Hunyuan with 1B parameters [38]. Training was conducted on 16 A100Pro GPUs, while inference was performed on hundreds of L40S GPUs. For the Ads indexing, we employed the Hunyuan to encode ad features. The number of residual quantization steps was set to 4, with each layer containing 1024 code vectors, each having a dimension of 8. The length of the S-ID sequence was set to 5. For LLM fine-tuning, we employed the AdamW optimizer, setting the learning rate to 6e-5, weight decay to 0.9, and a minimum learning rate of 6e-6. By leveraging data parallelism and gradient accumulation, the batch size was set to 64. To avoid overfitting, training was performed for 3 epochs.

### 5.2 Offline Performance

*5.2.1 Effectiveness of Prompt Components.* We perform ablation studies on the different components of the prompt in the main task (Section 4.2) to assess their individual contributions.

*Content Domain Behaviors and Interest Summary:* To evaluate the impact of content domain behaviors and user interest summary, we design two variants: "w.o. content" and "w.o. summary". The "w.o. content" variant removes the content domain behaviors from the prompt, retaining only the ad domain behaviors, while "w.o. summary" removes the user interest summary from the prompt.

---

[1] https://github.com/NVIDIA/TensorRT-LLM/tree/release/0.5.0

**Table 1: Comparison of performances under different components of prompt. "w.o. content" variant removes the content domain behaviors from the prompt, retaining only the ad domain behaviors, and "w.o. summary" removes the user interest summary from the prompt.**

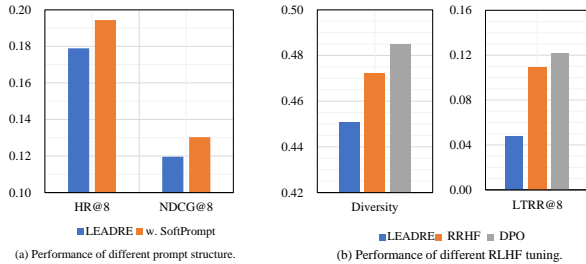| Variants | HR@1 | HR@4 | HR@8 | NDCG@4 | NDCG@8 |
|---|---|---|---|---|---|
| LEADRE | **0.0764** | **0.1567** | **0.2021** | **0.1199** | **0.1360** |
| w.o. content | 0.0660 | 0.1343 | 0.1773 | 0.1029 | 0.1181 |
|  | -13.60% | -14.28% | -12.26% | -14.19% | -13.15% |
| w.o. summary | 0.0760 | 0.1543 | 0.2043 | 0.1181 | 0.1358 |
|  | -0.62% | -1.54% | +1.08% | -1.52% | -0.18% |



(a) Performance of different prompt structure.
(b) Performance of different RLHF tuning.

**Figure 4: Ablation Studies Results on (a) Soft Prompt and (b) DPO tuning.**

As shown in Table 1, both variants lead to a performance drop across all metrics. This demonstrates the effectiveness of including both content domain behaviors and the user interest summary in capturing user preferences and intent.

*Soft Prompt:* To investigate the effectiveness of adding an instruction-based task description, we introduce a learnable token at the beginning of the prompt, referred to as "w. SoftPrompt". The results, depicted in Figure 4(a), show a performance improvement after incorporating the learnable token, confirming the positive impact of the SoftPrompt on the task.

*5.2.2 Effectiveness of Tuning Tasks and Their Order.* To study the effect of various tuning tasks and their order, we design several tuning strategies: "EX → IM → Main", "EX → Main", "Main+IM+EX mix", "Main+EX mix", and "Main only". The "EX → IM → Main" strategy refers to tuning the LLM sequentially on the explicit alignment task, implicit alignment task, and then the main task. The "EX → Main" strategy skips the implicit task and directly tunes on the explicit task followed by the main task. "Main+IM+EX mix" and "Main+EX mix" represent mixed-tuning strategies where pairs are randomly sampled and tuned together. "Main only" refers to tuning solely on the main task, without alignment tasks.

Table 2 presents the tuning results, leading to the following observations: (1) The explicit and implicit alignment tasks enhance the performance of the main task by bridging the gap between the language model and the ad system. (2) The order of tuning matters, as the fixed-order strategies outperform the mixed strategies. This suggests that the alignment tasks provide a foundational understanding of the ad system, which better prepares the LLM for the main task.

*5.2.3 Effectiveness of DPO.* After tuning the LLM on the auxiliary and main tasks, we apply Reinforcement Learning from Human Feedback (RLHF) to align the language model with business objectives, using both RRHF [50] and DPO [31] techniques. To evaluate

**Table 2: Comparison of performances under different tuning tasks and order. "Main" denotes the main task (Next Ad Generation), "EX" denotes the explicit alignment task, "IM" denotes the implicit alignment task, and "mix" denotes mix the tuning $< Prompt, Response >$ pairs and random sample pairs to tune the LLM.**

| Tasks and Order | HR@1 | HR@4 | HR@8 | NDCG@4 | NDCG@8 |
|---|---|---|---|---|---|
| EX → IM → Main | **0.0780** | **0.1610** | **0.2062** | **0.1232** | **0.1391** |
| EX→ Main | 0.0767 | 0.1566 | 0.1998 | 0.1201 | 0.1354 |
|  | -1.71% | -2.74% | -3.09% | -2.50% | -2.67% |
| Main+IM+EX mix | 0.0745 | 0.1523 | 0.2009 | 0.1165 | 0.1337 |
|  | -4.60% | -5.43% | -2.55% | -5.39% | -3.90% |
| Main+EX mix | 0.0766 | 0.1516 | 0.1948 | 0.1174 | 0.1327 |
|  | -1.91% | -5.84% | -5.49% | -4.64% | -4.59% |
| Main only | 0.0707 | 0.1444 | 0.1879 | 0.1107 | 0.1260 |
|  | -9.42% | -10.32% | -8.87% | -10.16% | -9.42% |

the alignment with business values, we use two metrics: Diversity score and LTRR@K.

The Diversity score is calculated based on TopK List Concentration and TopK List Abundance. Concentration is the mean proportion of the most frequent category across all users, while Abundance refers to the mean number of categories present in the top-K list for all users. The Diversity score combines both metrics, where a higher score indicates a more diverse top-K list with lower Concentration and higher Abundance. LTRR@K is computed by Recall@K based on Learning-to-Rank (LTR) labels provided by other retrieval strategies, with higher LTRR@K indicating better alignment with business objectives.

The RLHF results, shown in Figure 4, lead to the following conclusions: (1) RLHF improves the alignment of the language model with business values. (2) DPO outperforms RRHF, likely due to the constraints imposed by the tuned model during the alignment process.

**Table 3: Comparison of different emb models.**

| Emb model | Accuracy | | | |
|---|---|---|---|---|
|  | Top1 | Top10 | Top50 | Top100 |
| Hunyuan Embedding[38] | **98.03%** | 96.46% | 86.61% | 65.67% |
| E5-large-instruct [41] | 97.21% | 96.48% | 87.42% | 66.38% |
| bge-m3[6] | 95.00% | 91.07% | 78.76% | 59.64% |
| Bert-Chinese[9] | 97.24% | 92.65% | 74.12% | 50.20% |
| Electra_Chinese[7] | 94.60% | 85.98% | 63.09% | 43.62% |
| Sentence-Bert/L12[34] | 95.40% | 90.11% | 70.15% | 48.44% |
| Sentence-Bert/L6[34] | 81.97% | 72.65% | 53.61% | 37.97% |
| XLNet Chinese[48] | 80.31% | 63.32% | 38.27% | 26.14% |
| CLIP-32[30] | 80.32% | 67.91% | 45.90% | 29.87% |
| CLIP-16[30] | 81.57% | 67.40% | 45.20% | 30.33% |
| T5[32] | 67.10% | 51.22% | 33.90% | 23.89% |

**Table 4: Performance Comparison of different emb models.**

| S-IDs | HR@4 | HR@8 | LTRR@4 | LTRR@8 |
|---|---|---|---|---|
| Hunyuan S-IDs | 0.1070 | 0.2140 | 0.0294 | 0.0443 |
| E5 S-IDs | 0.1052 | 0.2140 | 0.0256 | 0.0407 |

*5.2.4 Effectiveness of S-IDs.* To assess the effectiveness of the ads' S-IDs, we compare the performance of different embedding models (emb. models). We first sampled 100 ads from each of the 48 ad categories and retrieved the Top-K nearest ads for each sampled ad based on embedding similarity. The accuracy is defined as the proportion of Top-K retrieved ads that belong to the same category
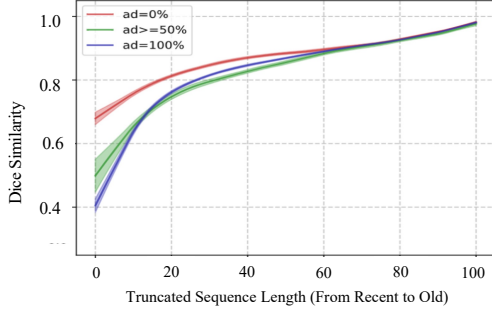
**Figure 5: LLM User Action Response w.r.t Sequence Length.**

as the original ad. The performance comparison is presented in Table 3. Our observations indicate that the E5 model demonstrates the best performance on this task, while the Hunyuan embedding achieves competitive results.

Subsequently, we trained the VQ-VAE using both the E5 and Hunyuan embeddings to infer the ads' S-IDs. We then constructed a trie tree and fine-tuned the LLM using these S-IDs. The performance of the tuned LLM is shown in Table 4. The results indicate that Hunyuan S-IDs exhibit competitive performance in terms of HR@K and outperform in LTRR@K, attributed to their enhanced capability in understanding the ad system.

*5.2.5 Effectiveness of User Behavior.* To assess how user behavior contributes to LEADRE's ad generation process, we conducted a user study to gain a deeper understanding of the generation mechanism.

To investigate the impact of different positions within user behavior sequences, we conducted the following experiment: For a user $u$ with a behavior sequence $\mathcal{S}_{u,1} = [i_u^1, i_u^2, \ldots, i_u^L]$, we sequentially truncated the sequence by removing earlier interactions, forming sequences such as $\mathcal{S}_{u,2} = [i_u^2, i_u^3, \ldots, i_u^L]$ and $\mathcal{S}_{u,3} = [i_u^3, i_u^4, \ldots, i_u^L]$. LEADRE then made predictions using each truncated sequence, generating top-K lists denoted by $\mathcal{R}_{u,1}, \mathcal{R}_{u,2}$, and so on. To quantify the similarity between these lists, we used the Dice similarity coefficient:

$$\text{Dice}(\mathcal{R}_{u,1}, \mathcal{R}_{u,l}) = \frac{2|\mathcal{R}_{u,1} \cap \mathcal{R}_{u,l}|}{|\mathcal{R}_{u,1}| + |\mathcal{R}_{u,l}|} \tag{1}$$

A higher Dice similarity score indicates that the retrieved list generated from the truncated sequence is more similar to that of the complete sequence. This, in turn, suggests that the removed interaction is of lesser importance in influencing the retrieval outcome. We computed the Dice similarity for all users, grouped by the partitioning of their ad behaviors, and plotted the results as a function of truncated sequence length (Figure 5). The following key observations were made: (1) Longer behavior sequences capture more user preference information. (2) LEADRE places more emphasis on the earlier parts of the sequence, particularly the first 20 behaviors, and less emphasis on more recent behaviors. (3) Ad domain behaviors have a stronger influence on the generated outcomes compared to content domain behaviors, and they encourage LEADRE to focus more on the beginning of the sequence.

## 5.3 Online Performance

We conducted a 20% traffic A/B test on both the Tencent WeChat Channels and Tencent WeChat Moments display advertising systems to evaluate the effectiveness of LEADRE. Over several weeks' A/B test on WeChat Channels, LEADRE demonstrated a 1.57% increase in Gross Merchandise Value (GMV) for serviced users. Similarly, several weeks' A/B test on WeChat Moments showed a 1.17% increase in GMV for serviced users. These results highlight LEADRE's capability to generate accurate, diverse, and high-value ads that contribute positively to business outcomes.

Currently, LEADRE has been deployed as a complementary retrieval sub-brunch on both Tencent WeChat Channels and Moments, serving billions of users and processing tens of billions of requests each day, highlighting its scalability and practical impact within a large-scale advertising system.

To effectively leverage the capability of LLMs in breaking the "information cocoon" effect, the retrieved ads are further incorporated into the ranking phase as new features. Specifically, on the user side, the retrieved ads are used as additional features to enhance user profiling, while on the item side, the match scores between the retrieved ads and the target ad are introduced as new item-level features. These new features contribute to an extra 1.43% improvement in GMV on Tencent WeChat Channels.

## 6 Conclusion

In this work, we present the first industrial application of generative retrieval in display advertising. To address the challenges of "*How to capture user commercial intent*", "*How to bridge the gap between LLMs and ad-Specific knowledge*", and "*How to efficiently deploy LLMs*", we introduce a novel LLM-based framework called Multi-Faceted Knowledge Enhanced **L**LM **E**mpowered Display **AD**vertisement **RE**commender system (LEADRE). LEADRE comprises three core modules: the Intent-Aware Prompt Engineering, the Advertising-Specific Knowledge Alignment, and the Efficient System Deployment. To evaluate the effectiveness of LEADRE, we conducted extensive experiments, including both offline and online evaluations. The results demonstrate its ability to generate more accurate and diverse ads. Moreover, online A/B test revealed that LEADRE resulted in a 1.57% increase on Tencent WeChat Channels, as well as a 1.17% increase on Tencent WeChat Moments in GMV for serviced users.

Looking ahead, our future work will explore the following directions: (1) Next-N Generation for LLMs: Next-N ad prediction provides a deeper understanding of user intent and offers a more accurate and diverse retrieval list. However, LEADRE currently supports only the generation of a single ad due to its reliance on a statistical trie-tree and fixed tuning tasks. Future efforts will focus on designing appropriate trie-tree structures and constructing tuning tasks for the next-N generation. (2) Development of More Reasonable S-IDs: Current results indicate that the RQ-VAE tends to allocate most information to the first codebook, while the subsequent codebooks retain limited information. Future research will aim to design an information-controllable quantization model to construct more effective S-IDs.

# References

[1] Keqin Bao, Jizhi Zhang, Xinyu Lin, Yang Zhang, Wenjie Wang, and Fuli Feng. 2024. Large Language Models for Recommendation: Past, Present, and Future. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2993–2996.

[2] Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 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.

[3] Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. Controllable multi-interest framework for recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2942–2951.

[4] Jianxin Chang, Chen Gao, Yu Zheng, Yiqun Hui, Yanan Niu, Yang Song, Depeng Jin, and Yong Li. 2021. Sequential recommendation with graph neural networks. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 378–387.

[5] Jun Chen, Cheng Chen, Huayue Zhang, and Qing Tan. 2022. A Unified Framework for Campaign Performance Forecasting in Online Display Advertising. *arXiv preprint arXiv:2202.11877* (2022).

[6] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation. arXiv:2402.03216

[7] K Clark. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555* (2020).

[8] Zhihua Cui, Xianghua Xu, XUE Fei, Xingjuan Cai, Yang Cao, Wensheng Zhang, and Jinjun Chen. 2020. Personalized recommendation system based on collaborative filtering for IoT scenarios. *IEEE Transactions on Services Computing* 13, 4 (2020), 685–695.

[9] Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[10] Manqing Dong, Feng Yuan, Lina Yao, Xiwei Xu, and Liming Zhu. 2020. Mamo: Memory-augmented meta-optimization for cold-start recommendation. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 688–697.

[11] 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.

[12] Zhabiz Gharibshah, Xingquan Zhu, Arthur Hainline, and Michael Conway. 2020. Deep learning for user interest and response prediction in online display advertising. *Data Science and Engineering* 5, 1 (2020), 12–26.

[13] Daya Guo, Canwen Xu, Nan Duan, Jian Yin, and Julian McAuley. 2023. Longcoder: A long-range pre-trained language model for code completion. In *International Conference on Machine Learning*. PMLR, 12098–12107.

[14] Jesse Harte, Wouter Zorgdrager, Panos Louridas, Asterios Katsifodimos, Dietmar Jannach, and Marios Fragkoulis. 2023. Leveraging large language models for sequential recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*. 1096–1102.

[15] Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian McAuley. 2023. Large language models as zero-shot conversational recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*.

[16] Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. *arXiv preprint arXiv:1704.07138* (2017).

[17] J Edward Hu, Huda Khayrallah, Ryan Culkin, Patrick Xia, Tongfei Chen, Matt Post, and Benjamin Van Durme. 2019. Improved lexically constrained decoding for translation and monolingual rewriting. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 839–850.

[18] Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *arXiv preprint arXiv:2210.11610* (2022).

[19] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-based retrieval in facebook search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2553–2561.

[20] Neil Hurley and Mi Zhang. 2011. Novelty and diversity in top-n recommendation–analysis and evaluation. *ACM TOIT* 10, 4 (2011), 1–30.

[21] Houye Ji, Ye Tang, Zhaoxin Chen, Lixi Deng, Jun Hu, and Lei Su. 2024. Neural Graph Matching for Video Retrieval in Large-Scale Video-driven E-commerce. *arXiv preprint arXiv:2408.00346* (2024).

[22] Huan Yee Koh, Jiaxin Ju, Ming Liu, and Shirui Pan. 2022. An empirical survey on long document summarization: Datasets, models, and metrics. *ACM computing surveys* 55, 8 (2022), 1–35.

[23] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

[24] Jin Li, Jie Liu, Shangzhou Li, Yao Xu, Ran Cao, Qi Li, Biye Jiang, Guan Wang, Han Zhu, Kun Gai, et al. 2021. Truncation-Free Matching System for Display Advertising at Alibaba. *arXiv preprint arXiv:2102.09283* (2021).

[25] Kunze Li and Yu Zhang. 2024. Planning First, Question Second: An LLM-Guided Method for Controllable Question Generation. In *Findings of the Association for Computational Linguistics ACL 2024*. 4715–4729.

[26] Xinyu Lin, Wenjie Wang, Yongqi Li, Fuli Feng, See-Kiong Ng, and Tat-Seng Chua. 2024. 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.

[27] Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. 2024. Using an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.

[28] Matt Post and David Vilar. 2018. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. *arXiv preprint arXiv:1804.06609* (2018).

[29] Guanghui Qin and Jason Eisner. 2021. Learning how to ask: Querying LMs with mixtures of soft prompts. *arXiv preprint arXiv:2104.06599* (2021).

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.

[31] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems* 36 (2024).

[32] 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. *Journal of Machine Learning Research* 21, 140 (2020), 1–67. http://jmlr.org/papers/v21/20-074.html

[33] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan Hulikal Keshavan, Trung Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Tran, Jonah Samost, et al. 2024. Recommender systems with generative retrieval. *Advances in Neural Information Processing Systems* 36 (2024).

[34] N Reimers. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv preprint arXiv:1908.10084* (2019).

[35] Xubin Ren, Wei Wei, Lianghao Xia, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 2024. Representation learning with large language models for recommendation. In *Proceedings of the ACM on Web Conference 2024*.

[36] Vinay Singh, Brijesh Nanavati, Arpan Kumar Kar, and Agam Gupta. 2023. How to maximize clicks for display advertisement in digital marketing? A reinforcement learning approach. *Information Systems Frontiers* 25, 4 (2023), 1621–1638.

[37] Derun Song, Enneng Yang, Guibing Guo, Li Shen, Linying Jiang, and Xingwei Wang. 2024. Multi-scenario and multi-task aware feature interaction for recommendation system. *ACM Transactions on Knowledge Discovery from Data* 18, 6 (2024), 1–20.

[38] Xingwu Sun, Yanfeng Chen, Yiqing Huang, Ruobing Xie, Jiaqi Zhu, Kai Zhang, Shuaipeng Li, Zhen Yang, Jonny Han, Xiaobo Shu, et al. 2024. Hunyuan-Large: An Open-Source MoE Model with 52 Billion Activated Parameters by Tencent. arXiv:2411.02265

[39] Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. 2021. Spot: Better frozen model adaptation through soft prompt transfer. *arXiv preprint arXiv:2110.07904* (2021).

[40] Bryan Wang, Gang Li, and Yang Li. 2023. Enabling conversational interaction with mobile ui using large language models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.

[41] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533* (2022).

[42] Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 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.

[43] Yinwei Wei, Xiang Wang, Qi Li, Liqiang Nie, Yan Li, Xuanping Li, and Tat-Seng Chua. 2021. Contrastive learning for cold-start recommendation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 5382–5390.

[44] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. 2022. A survey on accuracy-oriented neural recommendation: From collaborative filtering to information-rich recommendation. *IEEE Transactions on Knowledge and Data Engineering* 35, 5 (2022), 4425–4445.

[45] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. 2020. SSE-PT: Sequential recommendation via personalized transformer. In *Proceedings of the 14th ACM conference on recommender systems*. 328–337.

[46] Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2024. A survey on large language models for recommendation. *World Wide Web* 27, 5 (2024), 60.

[47] Ruobing Xie, Qi Liu, Liangdong Wang, Shukai Liu, Bo Zhang, and Leyu Lin. 2022. Contrastive cross-domain recommendation in matching. In *Proceedings of the*

*28th ACM SIGKDD conference on knowledge discovery and data mining.*

[48] Zhilin Yang. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. *arXiv preprint arXiv:1906.08237* (2019).

[49] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing* (2024), 100211.

[50] Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. 2024. RRHF: Rank responses to align language models with human feedback. *Advances in Neural Information Processing Systems* 36 (2024).

[51] ChengXiang Zhai. 2024. Large language models and future of information retrieval: Opportunities and challenges. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval.*

[52] Jizhi Zhang, Keqin Bao, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 2023. Is chatgpt fair for recommendation? evaluating fairness in large language model recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems.* 993–999.

[53] Jing Zhang, Hui Gao, Peng Zhang, Boda Feng, Wenmin Deng, and Yuexian Hou. 2024. LA-UCL: LLM-augmented unsupervised contrastive learning framework for few-shot text classification. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024).* 10198–10207.

[54] Mingxuan Zhang, Bo Yuan, Hanzhe Li, and Kangming Xu. 2024. LLM-Cloud Complete: Leveraging cloud computing for efficient large language model-based code completion. *Journal of Artificial Intelligence General science (JAIGS) ISSN: 3006-4023* 5, 1 (2024), 295–326.

[55] Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics* 12 (2024), 39–57.

[56] Bowen Zheng, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, Ming Chen, and Ji-Rong Wen. 2024. Adapting large language models by integrating collaborative semantics for recommendation. In *2024 IEEE 40th International Conference on Data Engineering (ICDE).* IEEE, 1435–1448.

[57] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107* (2023).

## A Preliminary

### A.1 Ads Indexing

The RQ-VAE takes the ad's textual embedding $\mathbf{x}_a$ as input and uses an **encoder-quantization-decoder** mechanism to generate a list of discrete IDs. In the encoding phase, the textual embedding $\mathbf{x}_a$ is encoded by an encoder, $f_{En}(\cdot; \theta_{En})$, to obtain a hidden embedding, formulated as $\hat{\mathbf{z}}_a = f_{En}(\mathbf{x}_a; \theta_{En})$, where $\hat{\mathbf{z}}_a \in \mathbb{R}^{d_{RQ}}$ denotes the hidden embedding, and $d_{RQ}$ is the hidden embedding's dimension. During the quantization phase, we apply $M$ layers of residual vector quantization. For each layer $l$, we maintain a codebook $C^l = \{\mathbf{e}_k^l\}_{k=1}^K$, where $\mathbf{e}_k^l \in \mathbb{R}^{d_{RQ}}$ and $K$ denotes the number of codes in the codebook. The quantization process is defined as follows:

$$c_a^l = \operatorname{argmin}_k ||\mathbf{r}_a^l - \mathbf{e}_k^l||_2^2, \quad \mathbf{r}_a^{l+1} = \mathbf{r}_a^l - \mathbf{e}_{c_a^l}^l, \quad \mathbf{z}_a = \sum_{l=1}^M \mathbf{e}_{c_a^l}^l$$

where $\mathbf{r}_a^l \in \mathbb{R}^{d_{RQ}}$ is the $l$-th layer residual embedding of ad $a$, $\mathbf{r}_a^1 = \hat{\mathbf{z}}_a$, $c_a^l$ is the index from the $l$-th codebook for ad $a$, and $\mathbf{z}_a \in \mathbb{R}^{d_{RQ}}$ is the final quantized embedding of ad $a$. In the decoding phase, the quantized embedding $\mathbf{z}_a$ is used to reconstruct the original textual embedding $\mathbf{x}_a$ through a decoder $f_{De}(\cdot; \theta_{De})$, formulated as $\hat{\mathbf{x}}_a = f_{De}(\mathbf{z}_a; \theta_{De})$

To train the encoder $\theta_{En}$, decoder $\theta_{De}$, and codebooks $\{C^l\}_{l=1}^M$, we introduce two loss components: the reconstruction loss and the quantization loss. The overall loss function for RQ-VAE is given by:

$$\mathcal{L}_{recons} = \sum_a ||\mathbf{x}_a - \hat{\mathbf{x}}_a||_2^2,$$

$$\mathcal{L}_{quant} = \sum_a \sum_{l=1}^M \left( ||\operatorname{sg}[\mathbf{r}_a^l] - \mathbf{e}_{c_a^l}^l||_2^2 + \beta_{quant} ||\mathbf{r}_a^l - \operatorname{sg}[\mathbf{e}_{c_a^l}^l]||_2^2 \right),$$

where $\operatorname{sg}[\cdot]$ is the stop-gradient operator, and $\beta_{quant}$ is a weight factor for the loss term.

This process ensures that ads are represented by discrete S-IDs, which allows the LLMs to generate valid ads from the predefined ads set during constrained decoding.

### A.2 Trie-Tree Construction

We employ a trie-tree to represent all ads, where each path from the root node to a leaf node corresponds to the S-IDs of a specific ad. The construction of the trie-tree follows the **Trie-Tree Construction Algorithm**. The algorithm initializes the trie-tree with a head node and processes each ad $a$ in the set $\mathcal{A}$ by sequentially examining its S-ID list $\mathbf{c}_a$ from the root.

At each step, the algorithm checks whether the current S-ID $c_a^l$ already exists as a child of the current node. If the S-ID is not present, a new node is created and added as a child, labeled with the current S-ID. The algorithm then updates the current node to the child node corresponding to the processed S-ID. After processing all S-IDs of an ad, the current node is marked as the end of the ad, indicating that a complete ad sequence has been stored. This process is repeated for all ads in the set $\mathcal{A}$. Once all ads have been processed, the algorithm returns the constructed trie-tree, which efficiently organizes the ads in a hierarchical structure. This structure allows for quick retrieval of ads based on their S-IDs and supports efficient prefix-based searches.

### A.3 Constrained Decoding

The Constrained Decoding algorithm leverages a trie-tree and beam search to ensure that the ads generated by LLMs conform to a predefined set of valid ads. The trie-tree, encompassing all possible ad S-IDs, constrains the LLM's output to valid S-IDs.

The algorithm initializes a generated ad beam set comprising $B$ lists, each starting with the head node of the trie-tree. For each layer of the trie-tree, corresponding to the position in the S-ID list, the algorithm retrieves valid next S-IDs based on the current node in the trie-tree. The LLM offers a probability distribution over feasible next S-IDs, constrained by the trie-tree structure. For each valid next S-ID, the algorithm appends the S-ID to the current candidate list, computes the score using the LLM's output probabilities, and adds the updated list to a new beam set. After processing all candidate lists in the current beam set, the top $B$ lists are selected based on their scores, and the process continues to the next layer of the trie-tree. This process is repeated until either all layers of the trie have been processed or the maximum S-IDs length ($M+1$) is reached. The algorithm returns the final beam set, which contains the generated ads that are valid according to the trie-tree and have the highest probabilities based on the LLM's predictions.

By combining beam search with trie-tree constraints, the algorithm ensures that the generated ads are both high-probability lists according to the LLM and valid ads from the predefined ads set.