

HLLM-Creator: Hierarchical LLM-based Personalized Creative Generation

Junyi Chen*
chenjunyi.s@bytedance.com
ByteDance

Lu Chi*
chilu@bytedance.com
ByteDance

Siliang Xu*
xusiliang@bytedance.com
ByteDance

Shiwei Ran
ranshiwei@bytedance.com
ByteDance

Bingyue Peng
bingyue.peng@bytedance.com
ByteDance

Zehuan Yuan†
yuanzehuan@bytedance.com
ByteDance

Abstract

AI-generated content (AIGC) technologies are increasingly used for content creation across diverse domains. However, most current AIGC systems depend heavily on the inspiration of content creators, with limited exploration into generating truly personalized content tailored to individual users. In real-world applications such as online advertising, a single product may have multiple selling points, with different users focusing on different features. This underscores the significant value of personalized, user-centric creative generation. Effective personalized content generation faces two main challenges: (1) accurately modeling user interests and integrating them into the content generation process while adhering to factual constraints, and (2) ensuring high efficiency and scalability to handle the massive user base in industrial scenarios. Additionally, the scarcity of personalized creative data in practice complicates model training, making data construction another key hurdle. To tackle these issues, we propose HLLM-Creator, a hierarchical large language model framework designed to efficiently model user interests and generate personalized creative content. During inference, a combination of user clustering and a user-ad-matching-prediction based pruning strategy is employed to significantly enhance generation efficiency and reduce computational overhead, making the approach suitable for large-scale deployment. Moreover, we design a data construction pipeline based on chain-of-thought (CoT) reasoning, which generates high-quality, user-specific creative titles and ensures factual consistency despite limited personalized data. This pipeline serves as a critical foundation for the effectiveness of our model. Extensive experiments on personalized advertising title generation for Douyin Search Ads show the effectiveness of HLLM-Creator. Online A/B test results demonstrate the practical value and scalability of our approach, with Adss increasing by **0.476%**, paving the way for more effective and efficient personalized content generation in real-world industrial applications. Codes for academic dataset are available at <https://github.com/bytedance/HLLM>.

CCS Concepts

• **Information systems** → **Online advertising; Recommender systems.**

Keywords

Personalized Creative Generation, Online Advertising, Recommendation System, Large Language Model

*Equal contribution. † Corresponding author.

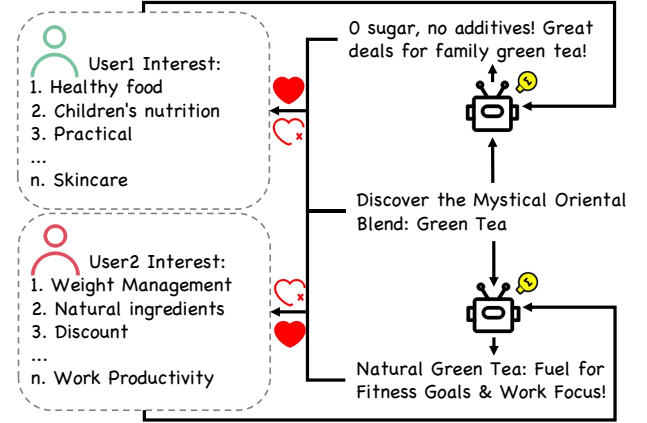


Figure 1: Personalized Creative Generation, which generates creatives more aligned with users’ preferences, achieving “Different strokes for different folks”.

1 Introduction

The landscape of information dissemination and content creation has undergone dramatic changes over the past decades. In the early days of the internet, information was primarily published on static web pages by a limited number of administrators or companies, offering users little opportunity for interaction. The emergence and rapid evolution of recommender systems fundamentally shifted this paradigm: systems began proactively pushing content most likely to interest each user [6, 25, 27], while the barriers to content creation were significantly lowered. Platforms such as TikTok have enabled anyone to become a content creator, resulting in an explosion of user-generated content. More recently, advances in AI-generated content (AIGC) technologies [11, 18, 23, 28] have further democratized content creation, leading to a proliferation of AI-generated materials across the web. Despite these advancements, a critical gap remains: the precise fulfillment of individual users’ personalized needs is still far from being realized. Taking online advertising as an example, a single product often has multiple selling points, and different users may focus on different aspects. An example of a personalized advertising title is illustrated in Figure 1. For advertisers, it is extremely challenging to create finely tailored ads for diverse user groups—both due to the sheer scale of creative work required and the difficulty of accurately identifying user preferences in advance, especially for small and medium-sized advertisers. As a

result, even when a product could meet a user’s needs, suboptimal ad descriptions may cause potential customers to be overlooked. Generating factually accurate, personalized ad descriptions is beneficial for all parties: advertisers can more precisely highlight product features and attract attention, users can quickly determine whether an ad matches their needs, and platforms can deliver content more accurately to relevant audiences. This underscores the importance and value of personalized creative generation.

Effective personalized creative generation must satisfy two key criteria: (1) accurately addressing user needs and pain points, and (2) strictly adhering to factual information, avoiding hallucinations and misleading “clickbait” content. However, most existing works on personalized generation have relatively weak user modeling capabilities, typically relying on simple user attributes (such as age and gender) [39] or keywords from users’ browsing history [1, 29] to represent user interests. Hallucination is unacceptable in the advertising domain. To tackle this challenge, previous works have constructed extensive factually grounded training datasets [21], ensured the completeness of input information to the model [14, 21, 34], and performed additional fine-tuning of generative models [14, 34]. In our work, we largely follow this series of strategies and additionally deploy a strict hallucination detection mechanism in online settings to further mitigate risks.

Scalability and efficiency are also critical for real-world deployment. Industrial applications often serve hundreds of millions of users and manage millions of ads daily. Efficiently generating personalized ad titles for such a massive user base is a significant challenge. Existing approaches tend to either compromise on user modeling (sacrificing effectiveness) [29, 39] or ignore the scale of the user base (sacrificing efficiency) [36, 37].

To address these challenges, we present HLLM-Creator, a novel framework for personalized creative generation. Our approach leverages Hierarchical Large Language Models (HLLM [3]) to accurately extract user interests from historical behavior. HLLM consists of two LLMs: the Item LLM encodes ad titles into item representations, while the User LLM aggregates representations of ads previously clicked by the user to produce a user embedding. This hierarchical modeling strategy not only significantly improves training and inference efficiency compared to directly inputting raw click sequences into an LLM, but also fully exploits the world knowledge embedded in LLMs, achieving state-of-the-art performance in sequential recommendation. Once the user embedding is obtained, it is combined with ad-side constraints (such as the original ad title and product selling points) and fed into a Creative LLM, which generates personalized ad titles in an autoregressive manner. The inclusion of user embeddings enables the Creative LLM to effectively perceive user interests, prioritize relevant selling points, and incorporate expressions likely to attract specific users (e.g., “a must-have for office workers”). At the same time, ad-side constraints guide the model to rewrite only based on the provided information, thereby ensuring factual consistency and minimizing hallucination.

Generating a unique personalized title for every ad-user pair is impractical in industrial settings. To maximize personalization benefits within resource constraints, we design two strategies. First, we cluster user embeddings into groups, with the number of clusters determined by available inference resources. This assumes that

users within a cluster share similar interests, allowing a single generated title to serve many users while still meeting personalization needs. Second, since each ad is typically relevant to only a subset of user groups, we introduce an ad-user matching model to predict relevance scores between each ad and all user cluster centers. We then generate personalized titles only for the top-matched clusters, optimizing resource utilization and maximizing ROI.

Beyond modeling, training data quality is another critical factor. In practice, advertisers often focus on making their titles broadly appealing by stacking multiple attractive selling points, but they rarely create distinct titles for different user groups.

This lack of personalization makes original advertiser titles sub-optimal for training personalized models. Furthermore, the high cost of manual annotation precludes large-scale labeled datasets. Consequently, leveraging synthetic data becomes a necessary alternative.

Some prior works attempt to leverage the generalization ability of generative models to produce personalized creatives by feeding user information into the model, but our experiments show that the quality of generated content is highly sensitive to prompt design. Meanwhile, synthetic data can lead to severe hallucination issues. To address this, we carefully design prompts and employ a chain-of-thought (CoT) approach to construct high-quality, user-personalized training data, substantially raising the performance ceiling. Furthermore, rigorous data cleaning procedures are applied to the constructed dataset to ensure the absence of hallucination issues in the training data.

Finally, our method has been deployed in the Douyin Search Ads Platform, a real-world industrial environment, generating personalized ad titles for hundreds of millions of users and delivering statistically significant improvements in key metrics through A/B testing. In summary, our main contributions are as follows:

- We propose a novel hierarchical LLM framework for personalized creative generation, enabling precise user interest modeling and fine-grained content personalization.
- A Chain-of-Thought-based data construction pipeline is developed to expand personalization space and ensure factual consistency, effectively reducing hallucinations in generated titles.
- A flexible and efficient inference scheme is developed for large-scale industrial deployment, with significant positive results in Douyin search advertising demonstrating its real-world impact.

2 Related Work

With the advancement of AIGC technology, it has been widely adopted in various domains, including text generation [4, 9, 13, 14, 17, 21, 26, 32, 34] and image generation [2, 7, 8, 15]. These technological breakthroughs have revolutionized content production by significantly reducing material creation costs. To align AIGC with user behavior objectives (e.g., click-through rate, CTR), several studies have attempted to integrate AIGC with reinforcement learning (RL) techniques. These efforts utilize real user behaviors as feedback signals to optimize generative models [34] or train reward models based on user actions to guide the refinement of generative processes [5, 12, 16, 40, 43]. While such methods have

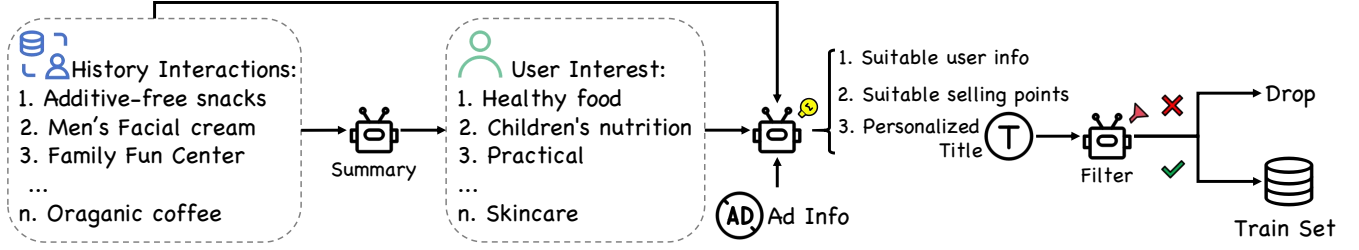


Figure 2: CoT-driven personalized title dataset construction pipeline. The LLM first summarizes user interests based on user historical sequences. Then it takes the user sequences, user interests, and advertising information as input to generate personalized titles. Finally, a strict LLM-based filtering mechanism is employed to remove hallucinated titles.

improved user click rates, they tend to generate high-performing texts or images tailored to the global user base, catering only to the preferences of the majority (i.e., the "head" users) rather than addressing the unique needs of individual users.

To address the limitations of previous methods, several personalized generation approaches have emerged [1, 20, 29, 30, 35–37, 39]. Despite their contributions, these personalized generation methods suffer from two key limitations: inadequate user modeling or inapplicability to large-scale industrial inference scenarios. For example, the personalized prompt model in CG4CTR [39] relies solely on simplistic user attributes (e.g., age, gender), which have been proven insufficient for modeling user interests in recommendation systems [25, 42]. The approach in [1], which summarizes user behaviors into a few keywords, incurs significant information loss and heavily depends on the quality of keyword extraction. While PMG [29] addresses the limitations of keywords by introducing soft preference embeddings—an idea somewhat similar to ours—we observe that direct end-to-end training is challenging to optimize. To mitigate this, we incorporate multiple auxiliary losses and leverage HLLM [3] for more effective modeling of user behavioral features. Regarding industrial deployment, existing works lack practical considerations. PMG [29], for instance, is difficult to deploy in scenarios with hundreds of millions of users. Pigeon [36] and DRC [37] adopt target-aware modeling approaches that result in a user-advertisement ($|\text{user}| \times |\text{ad}|$) scale, making them computationally prohibitive. In contrast, our work explicitly addresses industrial deployment requirements by designing clustering and pruning strategies, ensuring greater flexibility and efficiency in large-scale industrial settings.

3 Method

In this section, we first provide the problem formulation for the personalized creative generation task, followed by a description of the data construction pipeline (Section 3.2), which is one of the key factors for the effectiveness of our approach. Finally, we present the overall architecture of HLLM-Creator (Section 3.3), training objectives (Section 3.4), and inference strategies (Section 3.5) for industrial deployment.

3.1 Problem Formulation

Given a user u and target ad I , we aim to learn a generation model $p_{\Theta}(y|u, I)$ that can generate a personalized creative y that aligns

more closely with user interests (where Θ denotes trainable parameters of the model). More specifically, y refers to the personalized advertising title, and u is modeled with the user behavior sequence $H_u = \{I_1, I_2, \dots, I_n\}$, where n denotes the sequence length. The available information of I includes the original ad title y_{orig} and the set of selling points s .

3.2 CoT-driven Personalized Title Dataset Construction

In real-world scenarios, advertisers rarely have the capacity to create personalized titles for distinct user groups, leading to a scarcity of personalized data for model training. To address this challenge, we leverage LLM (e.g., DeepSeek-R1 [10]) to construct synthetic personalized datasets for training purposes. A naive approach would involve feeding user behavior sequences, the target ad title, and key selling points into the LLM and prompting it to directly generate personalized titles. However, we found that titles generated through this method exhibit limited personalization (experimental results can be found in Section 4.4.2). Inspired by chain-of-thought (CoT) [33] techniques widely adopted in the field of LLM, we decompose the personalized title generation process into several key subtasks and leverage the intermediate results of these subtasks to ultimately generate high-quality personalized titles. Figure 2 illustrates the proposed CoT-driven high-quality dataset construction pipeline. All prompts are provided in the Appendix B.

3.2.1 User Interest Profiling. We first instruct the LLM to extract predefined multi-dimensional user interests (such as long-term interests, short-term preferences, and specific needs) from user behavior sequences H . Then, we prompt the LLM to generate personalized titles based on these user profiles, combined with the original ad title and ad selling points.

3.2.2 Interest-Driven Title Generation. Regarding the forms of user personalization in title generation, we hypothesize they can be categorized into two types: one is directly incorporating descriptions targeting specific user groups (e.g., "a must-have for office workers"), and the other is integrating advertising selling points that users may be more interested in. Therefore, during the generation of personalized titles, we guide the LLM to first extract user information suitable for integration into the current advertisement based on the provided user interests, then identify the selling points that the user is more likely to care about, and finally these

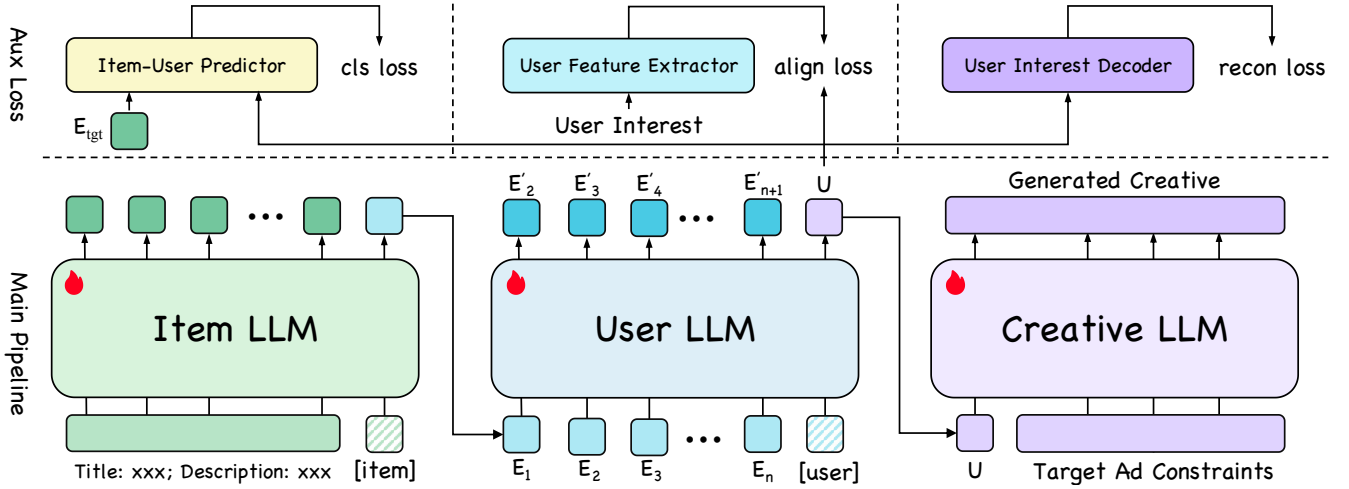


Figure 3: Overview of HLLM-Creator training framework. HLLM-Creator includes three LLMs: Item LLM, User LLM, and Creative LLM. Item LLM and User LLM are used to model user interests, while Creative LLM generates personalized titles for different users. Additionally, three auxiliary losses are added to enhance the extraction of user interests.

elements are combined to generate a personalized ad title. This CoT-based generation method reduces the processing difficulty for the LLM, allows for more thorough consideration of user interests and advertising selling points, and thus generates higher-quality titles.

3.2.3 Hallucination-Free Title Filtering. Synthetic data inevitably suffers from hallucination issues (i.e., fabricating non-existent information), even if we emphasize that the LLM should avoid fabricating information out of thin air during the synthesis process. Although titles with hallucinations may offer superior personalization, such "clickbait-style titles" are strictly prohibited in industrial scenarios. Therefore, we employ the LLM to implement a strict filtering mechanism on the generated personalized titles. It can effectively mitigate this critical issue by ensuring strict adherence to the source material.

3.3 Model Architecture

The model architecture and training process of HLLM-Creator are shown in Figure 3. HLLM-Creator consists of three LLMs in total: Item LLM is used to extract item features, which serve as input to the User LLM; User LLM adopts a late fusion structure as introduced in HLLM [3] for extracting user features; and Creative LLM is used for generating personalized content.

Specifically, for each $I \in H$, we first flatten its corresponding textual attributes into the sentence \mathcal{T} . After passing through the LLM tokenizer, we append a special token [item] at the end, thus the input token sequence for the Item LLM can be formulated as $\{t_1, t_2, \dots, t_m, [\text{item}]\}$ where m represents the length of item text tokens. The hidden state from the last layer corresponding to the special token [item] is considered as the item embedding E .

Then the original user history sequence $H = \{I_1, I_2, \dots, I_n\}$ can be transformed into a historical feature sequence $E = \{E_1, E_2, \dots, E_n\}$ through the Item LLM, where E_i represents the item embedding of I_i . Similar to the Item LLM, we append a special [user] token

after the feature sequence E as input to the User LLM, and take the hidden state corresponding to the [user] token as the user embedding U .

The Creative LLM first converts the original text information of advertisements (such as title, selling points, etc.) into corresponding embeddings $F = \{f_1, f_2, \dots, f_a\}$ through the word embedding layer, where a represents the length of the ad text tokens. These embeddings are then concatenated with the user embedding U to form the input sequence $\{U, f_1, f_2, \dots, f_a\}$ for the Creative LLM, enabling the model to generate personalized titles tailored to the user's interests.

3.4 Training Objectives

3.4.1 Generative Loss. The main training objective of HLLM-Creator is consistent with conventional Supervised Fine-tuning (SFT). It supervises the output of Creative LLM to align with the synthetic high-quality personalized titles introduced in Section 3.2, adopting the paradigm of next token prediction for training. The generative loss function can be formulated as:

$$\mathcal{L}_{gen} = -\frac{1}{L} \sum_{i=1}^L \log(p(r_j | U, f_1, f_2, \dots, f_a, r_1, \dots, r_{i-1})) \quad (1)$$

where r (response) denotes the synthesized training data, and L is the length of the response text tokens.

However, using only the training objective of personalized title generation for end-to-end training of all model parameters yields only suboptimal results. We attribute this to the fact that optimizing the User LLM for extracting user embeddings is relatively challenging. Therefore, three auxiliary losses are introduced in the following sections to explicitly supervise the user embedding extraction.

3.4.2 Recommendation Objective Loss (cls loss). A user feature that adequately models user interests should be sufficiently capable of judging user preferences. Therefore, referring to the late fusion version of HLLM, we introduce a classification loss for whether

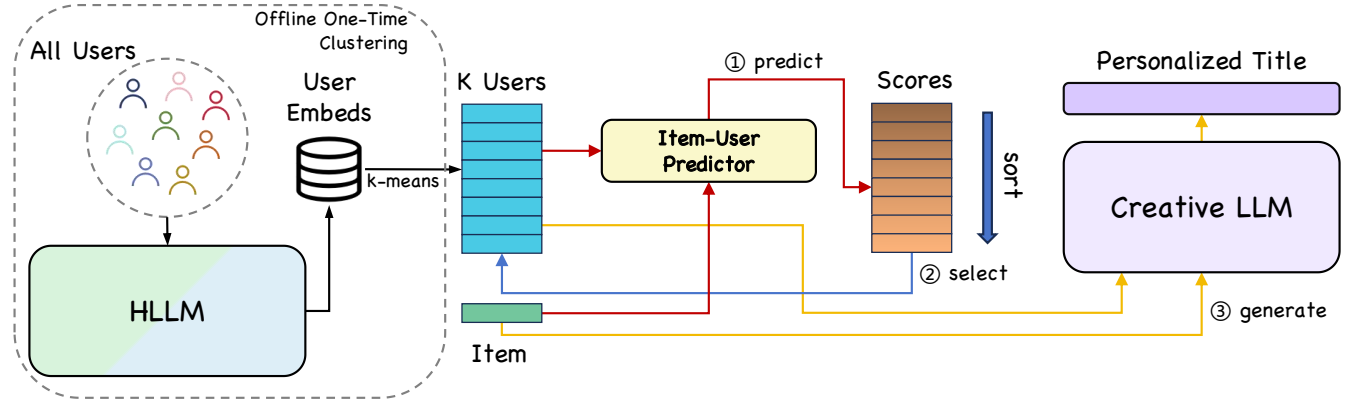


Figure 4: Inference workflow of HLLM-Creator. All users first have their user embeddings extracted via HLLM, followed by offline one-time clustering into K clusters. During actual inference, the top- k user clusters that best match the ad are selected to generate personalized titles.

the user clicks on the candidate item to supervise the training of user features. Here, positive examples are items that were actually clicked by the user, while negative examples are randomly sampled. Specifically, we input the user feature U and the candidate item feature E , which is extracted by the Item LLM, into the Item-User Predictor, implemented by a shallow fully connected neural network, to output the click probability. The cls loss function can be formulated as:

$$\mathcal{L}_{\text{cls}} = -\frac{1}{M} \sum_{i=1}^M [y_i \log(\sigma(f_{\text{dense}}(U_i, E_i))) + (1 - y_i) \log(1 - \sigma(f_{\text{dense}}(U_i, E_i)))] \quad (2)$$

where f_{dense} denotes the Item-User Predictor, U_i and E_i represent user embedding and item embedding respectively, σ is the sigmoid function. y indicates whether the user clicked, and M is the total number of positive and negative samples in a training batch.

3.4.3 Semantic Alignment Loss (align loss). Here, we use the user interests extracted in Section 3.2 to assist in supervising the learning of U . We use an LLM as the User Feature Extractor, and append a special token after the user interests text description to extract interest features V . Contrastive learning is performed between V and the user embeddings U generated by the User LLM, using the InfoNCE [22] loss to pull closer the two types of features from the same user and push apart the two types of features from different users. The align loss function can be formulated as:

$$\mathcal{L}_{\text{align}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(U_i, V_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(U_i, V_j)/\tau)} \quad (3)$$

where sim denotes a cosine similarity function, τ is temperature, and N is the number of samples in a training batch.

3.4.4 Reconstruction Loss (recon loss). We hypothesize that user embeddings should contain compressed user interest information. To achieve this, a decoder (LLM) is connected after the user embeddings, where the supervision target for the decoder output is the user interest description extracted in Section 3.2, and the training

follows the next token prediction paradigm. The recon loss function can be formulated as:

$$\mathcal{L}_{\text{recon}} = -\frac{1}{T} \sum_{t=1}^T \log P(w_t | U, w_1, \dots, w_{t-1}) \quad (4)$$

where w denotes the user interest description, and T is its text length.

The overall training objective of HLLM-Creator is:

$$\mathcal{L} = \mathcal{L}_{\text{gen}} + \lambda_{\text{cls}} \mathcal{L}_{\text{cls}} + \lambda_{\text{align}} \mathcal{L}_{\text{align}} + \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} \quad (5)$$

where each λ represents the weight of the corresponding loss.

3.5 Inference workflow

Generating personalized titles for every ad-user pair is impractical in industrial scenarios. Two strategies are designed during inference to maximize the benefits of personalized creative generation under limited resources.

As shown in Figure 4, we first use HLLM (Item LLM + User LLM) to extract a large number of user embeddings, and then perform k-means clustering based on these user embeddings to categorize users into K groups. During inference, the cluster center embedding U_k is fed into the Creative LLM as the user embedding, while other ad information remains unchanged. While clustering may compromise some granularity of personalization, it enables a single generation to cover more users, thereby enhancing production efficiency. Related experiments can be found in Section 4.4.5.

Furthermore, usually each ad has limited target audience, most ads will not reach all user groups. Therefore, newly created ads do not need to generate personalized titles for all K clusters. Instead, we introduce a predictor to determine the matching score between user groups and advertisements, and only perform personalized generation for user groups with a high matching score. Here we reuse the Item-User Predictor shown in Figure 3. In the inference phase, the K cluster center embeddings are fed into the Item-User Predictor with the current ad embedding E_{tgt} respectively to determine their matching degree with the current ad. Based on these scores, the top- k user clusters are selected for personalized title generation.

4 Experiments

HLLM-Creator is validated on Douyin Search Ads Platform.

4.1 Implementation Details and Evaluation Setup

4.1.1 Dataset Construction. An industrial dataset from Douyin Search Ads is adopted for training and evaluation. First, we selected 2.6 million samples from online user click logs, with each sample consisting of the user’s click sequence, search query, original ad title, and selling points. Then, personalized titles were generated using the CoT-driven data construction method introduced in Section 3.2. After filtering out personalized titles with hallucinations, approximately one quarter of the data remained. The filtered data are split into 650K samples for model training and 500 samples for offline evaluation.

4.1.2 Training Configuration. TinyLlama-1.1B [41] is adopted for both the Item LLM and the User LLM, while the Creative LLM employs Qwen3-8B [38]. The user’s historical click sequence is truncated to a maximum length of 500. Only the title with a maximum of 64 tokens is used as input to the Item LLM. All auxiliary loss weights are set to 1. The user embedding U is aligned to the dimension of the Creative LLM through a single-layer fully connected network. The Item-User Predictor is an 8-layer MLP-Mixer [31] with a hidden dimension of 2048. To save GPU memory, the parameters of the User Feature Extractor, User Interest Decoder, and Creative LLM are shared. All parameters are initialized with the pre-trained LLM weights and trained end-to-end for one epoch with a learning rate of $2e-5$. Unless otherwise specified, all experimental results for HLLM-Creator are based on clustering into 256 groups.

4.1.3 Incorporating Search Query. Our evaluation scenario is primarily the search scenario, where the user’s search query is important information that represents the user’s immediate interests. Therefore, reflecting the user’s search query in the title is very important. Our method is highly flexible in incorporating new constraints, with two main upgrades: First, we introduce the user’s search query during data construction (the specific prompt can be found in the Appendix B); second, we update the input to the Creative LLM by appending the search query after the user embedding U and ad information.

4.1.4 Reproducibility. To facilitate the reproduction and validation of the effectiveness of HLLM-Creator, we conducted experiments on the academic dataset (Amazon Book Reviews [19]). The experimental results are provided in the Appendix A, and the relevant codes are available at <https://github.com/bytedance/HLLM>.

4.1.5 Evaluation Metrics.

Online Evaluation Metrics. We use three key metrics: Adss (Advertiser Score), Advv (Advertiser Value), and RankAdvv (Rank Advertiser Value) for online evaluation. The latter two are defined as:

$$\text{Advv} = \text{cpa_bid} \times \text{conversions} \quad (6)$$

$$\text{RankAdvv} = \text{rank_bid} \times \text{conversions} \quad (7)$$

where cpa_bid denotes the cost-per-action bid, and rank_bid is the bid used by the ads ranking system.

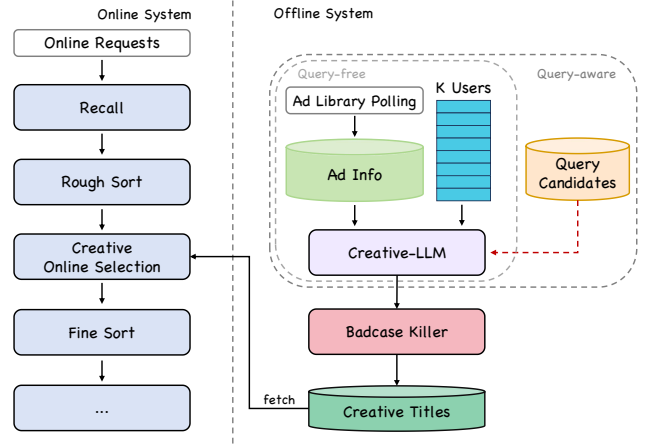


Figure 5: Industry deployment framework. See the main text for details.

Offline Evaluation Metrics. Directly quantifying the degree of user interest in different titles is inherently challenging. We adopt the Good-Same-Bad (GSB) preferences of GPT4.1 [24]. Prompt can be found in Appendix B. To avoid order bias in LLM evaluation, we test both (A, B) and (B, A) prompt orders. Only consistent results are counted; otherwise, the outcome is marked as "Same". We define "advantage" as follows:

$$\text{Advantage} = \frac{N_{\text{Good}} - N_{\text{Bad}}}{N_{\text{Good}} + N_{\text{Same}} + N_{\text{Bad}}} \quad (8)$$

For hallucination, we set up a hallucination detection pass rate metric. Using the same prompt as in the data cleaning process, we perform hallucination detection on the generated titles and calculate the proportion of titles without hallucinations. A higher value for this metric indicates a stronger ability of the model to adhere to factual information.

4.2 Online A/B Test

HLLM-Creator is validated on Douyin Search Ads Platform, a real-world industrial recommendation system. As shown in Figure 5, our online A/B testing deployment consists of two key components: offline generation and online selection.

4.2.1 Offline Generation. Due to the inherent latency of LLM inference, real-time generation of personalized content in an online environment is impractical. To address this limitation, we adopt an offline inference strategy, where personalized titles are pre-generated in batches as described in Section 3.5. Specifically, we cluster user embeddings at the million scale into 256 groups, and then generate personalized titles for each ad campaign. Two generation strategies are adopted for each ad campaign: (1). **Query-Aware:** Personalized titles are generated by incorporating the ad’s potential queries (based on historical statistics or model predictions). These titles fully leverage the relevance among the user, query, and ad. Considering resource constraints, this strategy generates titles only for the top-1 user group predicted by the Item-User Predictor. (2). **Query-Free:** Personalized titles are generated using only user and

Table 1: Douyin Search Ads A/B Test Results.

Adss↑	Advv↑	RankAdvv↑
+0.476%	+0.297%	+0.51%

ad information, targeting the top-5 user groups. This approach enables coverage of a broader range of user groups.

We continuously poll and generate content for all active advertising campaigns across the entire library, with each update cycle taking around 5 hours. Generated titles are then processed by a rule-based filter, referred to as the Badcase Killer, which removes undesirable cases such as fabricated content, inappropriate numbers, locations, brands, and sensitive terms. The qualified titles will be stored in the creative library and used as candidates for material selection during online serving.

4.2.2 Online Selection. The online request pipeline remains largely unchanged. For each search request, relevant ad campaigns are recalled and go through a coarse ranking stage. After that, they enter the creative online selection module, where title selection is performed for each campaign. The personalized titles are added as new candidates to participate in the selection process. The winning title is then used as the title for the corresponding ad campaign in subsequent fine-ranking and other downstream stages.

4.2.3 A/B Test Result. A/B test results can be found in Table 1. Compared to the baseline, introducing personalized titles led to improvements of **0.476%**, **0.297%**, and **0.510%** on the Adss, Advv, and RankAdvv metrics, respectively, which are considered good gains for a highly optimized system. It is worth noting that the baseline already includes various AI-Generated titles, including query-aware titles, which further demonstrate the benefits of introducing personalization.

Additionally, we conducted an in-depth study of the Click-Through Rate (CTR) for personalized titles. During the A/B testing process, we marked ads where personalized titles won in the selection stage. For the corresponding baseline group, personalized titles were skipped, and the top-1 title was selected from the remaining titles for subsequent stages. By directly comparing the CTR performance of these ads between the baseline and experimental groups, the results showed that the experimental group achieved a **1.789%** increase in CTR.

4.3 Comparison with other Methods

Figure 6 shows the GSB results of HLLM-Creator compared with other methods. Pigeon [36] is a recent state-of-the-art work on personalized generation. We reproduced this method on search ads data. Pigeon is much less efficient than HLLM-Creator, mainly because it flattens the historical title sequences into a single sequence and its reference-aware modeling of historical behavior results in an inference complexity of $|\text{user}| \times |\text{ad}|$, whereas ours is only $K \times |\text{ad}|$. In terms of effectiveness, with the same sequence length, Pigeon has a 7.4% advantage. However, when aligning training complexity, HLLM-Creator can extend the sequence length to 500, at which point it achieves a 4.2% advantage.

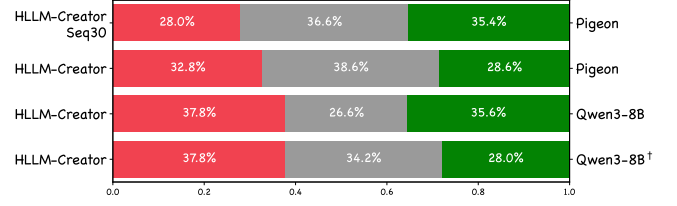


Figure 6: Comparison of HLLM-Creator with other baselines on LLM GSB evaluation. † indicates testing on the subset of titles from both methods that are free of hallucinations.

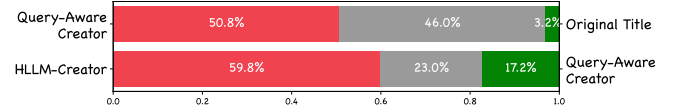


Figure 7: Ablation study on Personalization Modeling.

Table 2: Ablation study on CoT-based personalized title generation in data construction pipeline.

CoT 1	CoT 2	Chat Rounds	Good	Same	Bad	Advantage↑
	✓	1	33.0%	41.4%	25.6%	7.4%
✓	✓	1	35.8%	42.0%	22.2%	13.6%
✓		2	45.0%	39.2%	15.8%	29.2%
✓	✓	2	57.8%	28.0%	14.2%	43.6%

We also compared our method with a general-purpose LLM trained on broad knowledge, by directly inputting user behavior sequences and user interests in text form into the LLM to generate personalized titles. For a fair comparison, we used Qwen3-8B [38]—the same base model as our Creative LLM. The results show that our method has a slight 2% advantage. Additionally, we observed that the LLM without finetuning performs poorly in terms of hallucination: as shown in Table 3, the hallucination detection pass rate for Qwen3-8B is 60%, while our method achieves 75%. When testing on the subset of titles from both methods that are free of hallucinations, the advantage of HLLM-Creator increases to 9.8%.

4.4 Ablation Study

4.4.1 Necessity of Personalization Modeling. We use the model trained on titles that are generated based solely on original ad information and user queries as the baseline, referred to as Query-Aware-Creator, to verify the necessity of further incorporating user modeling. As shown in Figure 7, the Query-Aware Creator already achieves a significant improvement over the original titles. However, our HLLM-Creator is able to achieve further gains on top of this, demonstrating the necessity of personalization modeling.

4.4.2 Effect of the CoT Process in Data Construction. There are two critical CoT steps in our data construction process: CoT1: User Interest Profiling; CoT2: Interest-Driven Title Generation. Ablation studies are shown in Table 2. The baseline here does not involve any

Table 3: Hallucination detection pass rates of different models.

Method	Data Cleaning	Selling Points	Titles	Pass↑
Qwen3-8B [38]	-			60%
HLLM-Creator	✗			29%
HLLM-Creator		✗		53%
HLLM-Creator			✗	9%
HLLM-Creator				75%

Table 4: Ablation study on auxiliary losses.

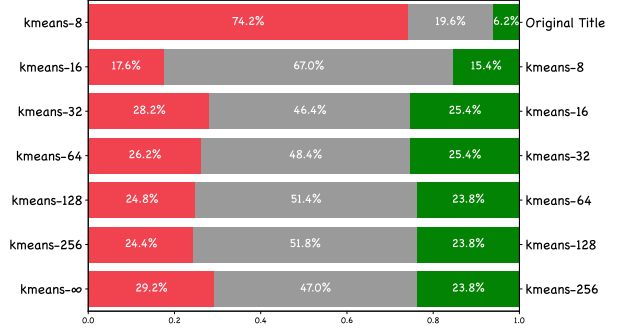
align	cls	recon	Good	Same	Bad	Advantage↑
			52.6%	25.6%	21.8%	30.8%
✓			56.8%	25.4%	17.8%	39.0%
	✓		56.4%	24.6%	19.0%	37.4%
		✓	53.8%	25.4%	20.8%	33.0%
✓	✓		57.6%	24.4%	18.0%	39.6%
✓	✓	✓	59.8%	23.0%	17.2%	42.6%

CoT process; instead, user behavior sequences, ad information, and query are directly input into the LLM to generate personalized titles in a single chat round. As shown, each CoT step brings significant performance improvements. Meanwhile, we attempted to combine the two CoT processes into a single dialogue, that is, requiring the LLM to first extract user interests, then identify key interests and selling points based on the ad information and user interests, and finally generate the title, all within one prompt. We found that, compared to the two-step dialogue, the single-step approach resulted in inferior performance.

4.4.3 Hallucination Issue. To address the hallucination problem, we implemented measures at both the data and model levels. On the data side, we used LLMs to perform hallucination detection on the constructed data and filtered out titles with hallucinations. On the model side, we provided sufficient original ad constraints in the input, including the original title and ad selling points. As shown in Table 3, each step contributes to alleviating the hallucination issue. Ultimately, we achieved a hallucination pass rate of 75%.

4.4.4 Auxiliary Loss. The impact of different auxiliary losses is explored in Table 4. Here, the baseline is Query-Aware-Creator. When trained without any auxiliary losses, the advantage of HLLM-Creator is only 30.8%. Experiments show that regardless of which auxiliary loss is used, it is beneficial for generating personalized titles, with the advantage increasing by 2.2% (30.8% → 33.0%) to 8.2% (30.8% → 39.0%). Moreover, the effects of different auxiliary losses are complementary: when three aux losses are used simultaneously, the advantage of HLLM-Creator reaches **42.6%**.

4.4.5 Clustering. To deploy HLLM-Creator in real-world industrial scenarios, we cluster all users before deployment to reduce inference costs. Clustering inevitably incurs clustering loss and affects the quality of generated personalized titles. Thus, we explore the impact of clustering on the performance in Figure 8. It can be observed from Figure 8 that even when clustered into a very small number of

**Figure 8: Ablation study on the impact of clustering. ∞ indicates no clustering applied.**

clusters (8), there is still a 68% advantage compared to original titles. It can also be observed that as the number of clusters increases, the improvement in advantage becomes less significant. However, there is still considerable room for improvement compared to not using clustering. For example, compared to kmeans-256, the no-clustering setting kmeans-∞ still achieves a 5.4% advantage. This indicates that our method has a higher potential upper bound, and in the future work, we will explore the online ROI of scaling up the number of clusters.

4.5 Case Study

By examining the cases generated by HLLM-Creator, we found that they are generally consistent with our previous assumptions: the personalization of the generated titles is mainly reflected in two aspects—explicitly adding descriptions related to user groups and incorporating personalized selling points. For example, the generated titles may explicitly include descriptions such as "professional women" to highlight user characteristics. Alternatively, for users who are more concerned about quality of life, the model selects selling points like "health" and "smart," while discarding less appealing features such as "cost-effectiveness" or "large capacity" that may not attract the current user as much.

To protect the privacy of users and advertisers, we used LLMs to construct some illustrative data in Appendix C, which is generally consistent with the experimental conclusions.

5 Conclusion

In this paper, we propose HLLM-Creator, an innovative personalized creative generation model based on hierarchical large language models. The model takes user historical behavior and ad-side constraint information as input and outputs personalized creatives. To enable deployment in real-world industrial environments, we designed clustering and pruning strategies based on the matching degree between ads and user interests. To address the lack of personalized data in real scenarios, we developed a CoT-based personalized data construction process and ensured the factual accuracy of generated data through rigorous data cleaning. In the context of Douyin search ads, we added personalized title candidates for each ad campaign. Extensive offline experiments verified the effectiveness of our design, and online A/B testing demonstrated statistically significant gains.

References

- [1] Pengshan Cai, Kaiqiang Song, Sangwoo Cho, Hongwei Wang, Xiaoyang Wang, Hong Yu, Fei Liu, and Dong Yu. 2023. Generating user-engaging news headlines. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 3265–3280.
- [2] Binghui Chen, Chongyang Zhong, Wangmeng Xiang, Yifeng Geng, and Xu-ansong Xie. 2024. Virtualmodel: Generating object-id-retentive human-object interaction image by diffusion model for e-commerce marketing. *arXiv preprint arXiv:2405.09985* (2024).
- [3] Junyi Chen, Lu Chi, Bingyue Peng, and Zehuan Yuan. 2024. Hllm: Enhancing sequential recommendations via hierarchical large language models for item and user modeling. *arXiv preprint arXiv:2409.12740* (2024).
- [4] Jiao Chen, Kehui Yao, Reza Yousefi Maragheh, Kai Zhao, Jianpeng Xu, Jason Cho, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. 2025. CARTS: Collaborative Agents for Recommendation Textual Summarization. *arXiv preprint arXiv:2506.17765* (2025).
- [5] Xingye Chen, Wei Feng, Zhenbang Du, Weizhen Wang, Yanyin Chen, Haohan Wang, Linkai Liu, Yaoyu Li, Jinyuan Zhao, Yu Li, et al. 2025. CTR-Driven Advertising Image Generation with Multimodal Large Language Models. In *Proceedings of the ACM on Web Conference 2025*. 2262–2275.
- [6] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishii Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ipsir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [7] Ádám Tibor Czapp, Máttyás Jani, Bálint Domián, and Balázs Hidasi. 2024. Dynamic product image generation and recommendation at scale for personalized e-commerce. In *Proceedings of the 18th ACM Conference on Recommender Systems*. 768–770.
- [8] Yifan Gao, Zihang Lin, Chuanbin Liu, Min Zhou, Tiezheng Ge, Bo Zheng, and Hongtao Xie. 2025. Postermaker: Towards high-quality product poster generation with accurate text rendering. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 8083–8093.
- [9] Xiaotao Gu, Yuning Mao, Jiawei Han, Jialu Liu, You Wu, Cong Yu, Daniel Finnie, Hongkun Yu, Jiaqi Zhai, and Nicholas Zukoski. 2020. Generating representative headlines for news stories. In *Proceedings of The Web Conference 2020*. 1773–1784.
- [10] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shiroong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
- [11] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).
- [12] Yashal Shakti Kanungo, Gyanendra Das, Pooja A, and Sumit Negi. 2022. Cobart: controlled, optimized, bidirectional and auto-regressive transformer for ad headline generation. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*. 3127–3136.
- [13] Yashal Shakti Kanungo, Sumit Negi, and Aruna Rajan. 2021. Ad headline generation using self-critical masked language model. In *Proceedings of the 2021 conference of the north American chapter of the association for computational linguistics: human language technologies: industry papers*. 263–271.
- [14] Zeyang Lei, Chao Zhang, Xinchao Xu, Wenquan Wu, Zheng-Yu Niu, Hua Wu, Haifeng Wang, Yi Yang, and Shuanglong Li. 2022. Plato-ad: a unified advertisement text generation framework with multi-task prompt learning. In *Proceedings of the 2022 conference on empirical methods in natural language processing: industry track*. 512–520.
- [15] Zhaochen Li, Fengheng Li, Wei Feng, Honghe Zhu, Yaoyu Li, Zheng Zhang, Jingjing Lv, Junjie Shen, Zhangang Lin, Jingping Shao, et al. 2023. Planning and rendering: Towards product poster generation with diffusion models. *arXiv preprint arXiv:2312.08822* (2023).
- [16] Jianghao Lin, Peng Du, Jiaqi Liu, Weite Li, Yong Yu, Weinan Zhang, and Yang Cao. 2025. Sell It Before You Make It: Revolutionizing E-Commerce with Personalized AI-Generated Items. *arXiv preprint arXiv:2503.22182* (2025).
- [17] Haoran Liu, Amir Tahmasbi, Ehtesham Sam Haque, and Purak Jain. 2025. LLMs for Customized Marketing Content Generation and Evaluation at Scale. *arXiv preprint arXiv:2506.17863* (2025).
- [18] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. 2024. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177* (2024).
- [19] 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.
- [20] Elyas Meguellati, Lei Han, Abraham Bernstein, Shazia Sadiq, and Gianluca Demartini. 2024. How good are llms in generating personalized advertisements?. In *Companion Proceedings of the ACM Web Conference 2024*. 826–829.
- [21] Masato Mita, Soichiro Murakami, Akihiko Kato, and Peinan Zhang. 2023. Striking gold in advertising: Standardization and exploration of ad text generation. *arXiv preprint arXiv:2309.12030* (2023).
- [22] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [23] OpenAI. 2022. Introducing ChatGPT. <https://openai.com/blog/chatgpt>.
- [24] OpenAI. 2025. Introducing GPT-4.1 in the API. <https://openai.com/index/gpt-4-1/>. Accessed: 2025-05-14.
- [25] Qi Pi, Guorui Zhou, Yujing Zhang, Zhe Wang, Lejian Ren, Ying Fan, Xiaoqiang Zhu, and Kun Gai. 2020. Search-based user interest modeling with lifelong sequential behavior data for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2685–2692.
- [26] Yilun Qiu, Xiaoyan Zhao, Yang Zhang, Yimeng Bai, Wenjie Wang, Hong Cheng, Fuli Feng, and Tat-Seng Chua. 2025. Measuring what makes you unique: Difference-aware user modeling for enhancing llm personalization. *arXiv preprint arXiv:2503.02450* (2025).
- [27] Steffen Rendle. 2010. Factorization machines. In *2010 IEEE International conference on data mining*. IEEE, 995–1000.
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [29] Xiaoteng Shen, Rui Zhang, Xiaoyan Zhao, Jieming Zhu, and Xi Xiao. 2024. Pmg: Personalized multimodal generation with large language models. In *Proceedings of the ACM Web Conference 2024*. 3833–3843.
- [30] Xiangru Tang, Xingyao Zhang, Yanjun Shao, Jie Wu, Yilun Zhao, Arman Cohen, Ming Gong, Dongmei Zhang, and Mark Gerstein. 2024. Step-back profiling: Distilling user history for personalized scientific writing. *arXiv preprint arXiv:2406.14275* (2024).
- [31] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. 2021. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems* 34 (2021), 24261–24272.
- [32] Varun Vasudevan, Faezeh Akhavanizadeh, Abhinav Prakash, Yokila Arora, Jason Cho, Tanya Mendiratta, Sushant Kumar, and Kannan Achan. 2025. LLM-driven Constrained Copy Generation through Iterative Refinement. *arXiv preprint arXiv:2504.10391* (2025).
- [33] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [34] Penghui Wei, Xuanhua Yang, Shaoqun Liu, Liang Wang, and Bo Zheng. 2022. CREATER: CTR-driven advertising text generation with controlled pre-training and contrastive fine-tuning. *arXiv preprint arXiv:2205.08943* (2022).
- [35] Yiyan Xu, Wenjie Wang, Fuli Feng, Yunshan Ma, Jizhi Zhang, and Xiangnan He. 2024. Diffusion models for generative outfit recommendation. In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*. 1350–1359.
- [36] Yiyan Xu, Wenjie Wang, Yang Zhang, Biao Tang, Peng Yan, Fuli Feng, and Xiangnan He. 2025. Personalized image generation with large multimodal models. In *Proceedings of the ACM on Web Conference 2025*. 264–274.
- [37] Yiyan Xu, Wuyang Zheng, Wenjie Wang, Fengbin Zhu, Xinting Hu, Yang Zhang, Fuli Feng, and Tat-Seng Chua. 2025. DRC: Enhancing Personalized Image Generation via Disentangled Representation Composition. *arXiv preprint arXiv:2504.17349* (2025).
- [38] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).
- [39] Hao Yang, Jianxin Yuan, Shuai Yang, Linhe Xu, Shuo Yuan, and Yifan Zeng. 2024. A new creative generation pipeline for click-through rate with stable diffusion model. In *Companion Proceedings of the ACM Web Conference 2024*. 180–189.
- [40] Jingying Zeng, Jaewon Yang, Waleed Malik, Xiao Yan, Richard Huang, and Qi He. 2023. Let AI Entertain You: Increasing User Engagement with Generative AI and Rejection Sampling. *arXiv preprint arXiv:2312.12457* (2023).
- [41] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385* (2024).
- [42] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1059–1068.
- [43] Jianghui Zhou, Ya Gao, Jie Liu, Xuemin Zhao, Zhaohua Yang, Yue Wu, and Lirong Shi. 2024. GCOF: Self-iterative Text Generation for Copywriting Using Large Language Model. *arXiv preprint arXiv:2402.13667* (2024).

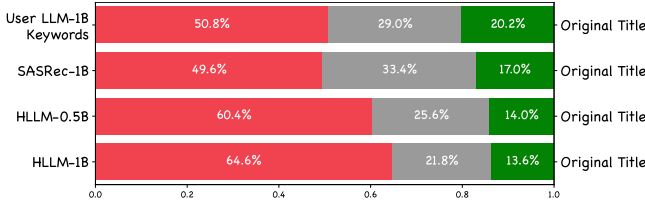


Figure 9: Comparison of HLLM-Creator with other user modeling methods on academic dataset.

Table 5: Ablation study of auxiliary losses on the academic dataset.

align	cls	recon	Good	Same	Bad	Advantage
			64.6%	21.8%	13.6%	51.0%
✓			67.4%	20.6%	12.0%	55.4%
	✓		67.6%	21.6%	10.8%	56.8%
		✓	65.2%	22.0%	12.8%	52.4%
✓	✓		67.6%	22.8%	9.6%	58.0%
✓	✓	✓	70.6%	18.2%	11.2%	59.4%

A Ablation Study on Academic Dataset

We constructed personalized title data based on the academic dataset Amazon Book Reviews [19] and conducted a series of experiments to demonstrate the effectiveness of HLLM-Creator. We use original title and description as item-side constraints. It should be noted that the constructed data was not subjected to hallucination filtering and was used solely for modeling validation of personalized generation, without practical significance.

A.1 User Modeling Method

In the main text, we have verified the necessity of personalized user modeling. Here, we further validate this conclusion on the academic dataset, showing that introducing personalized modeling leads to significant improvements compared to using original titles. Additionally, we demonstrate that the quality of user modeling also has an important impact on personalized generation. As shown in Figure 9, we first replace the Item LLM with keywords extracted from titles by LLM or with item IDs, labeled as User LLM-1B Keywords and SASRec-1B, respectively. Compared to the original titles, these approaches yield improvements of 30.6% and 32.6%, respectively, but are less advantageous than modeling user features with HLLM (51%), denoted as HLLM-1B. Furthermore, the size of the HLLM model also affects the results: when the model size is reduced from 1B to 0.5B, the advantage decreases to 46.4%.

A.2 Auxiliary Loss

We also validate the effectiveness of the auxiliary loss on the academic dataset, with the results presented in Table 5, where the baseline is the original title. Consistent with the findings on the industrial dataset, HLLM-Creator’s advantage increases from 51.0% (without any auxiliary loss) to 59.4%.

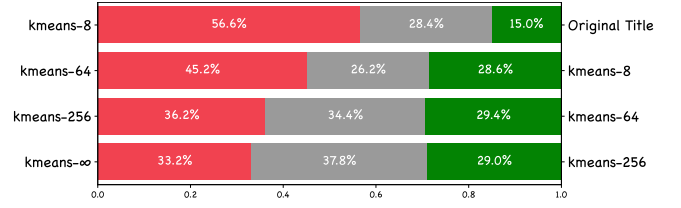


Figure 10: Ablation study of the impact of clustering on the academic dataset. ∞ indicates no clustering applied.

A.3 Clustering

Table 10 shows the impact of user clustering on the academic dataset. Similarly to the conclusions on the industrial dataset, user clustering is detrimental to the performance of generating personalized titles. As the number of clusters increases, the model performance gradually improves, with the best results achieved when no clustering is applied.

B Prompt Engineering

In Section 3.2 and Section 4.1.5 of the main text, we designed several prompts. Here, we provide the specific content of these prompts. For ease of reading, we have translated the Chinese into English.

Figure 11 shows the prompt of User Interest Profiling, which covers user interests across multiple dimensions.

Figure 12 displays the prompt of Interest-Driven Title Generation.

Figure 13 presents the prompt of Hallucination-Free Title Filtering, which is also used in hallucination evaluation. By filtering out hallucination-containing data, we ensure the reliability of the personalized title generation model in industrial scenarios.

Figure 14 shows prompts for offline model performance evaluation. By evaluating (model A, model B) and (model B, model A) pairs twice, we ensure the reliability of the evaluation results.

C Case Study

Figure 15 and Figure 16 present the intermediate results of the CoT process (including user interests, user interests, and selling points suitable for inclusion in the title), ad information (original title and ad selling points), and the final generated personalized titles. For privacy protection, the cases shown here are generated by large language models and are not real data; they are for demonstration purposes only. However, the conclusions are generally consistent with those observed in real-world scenarios.

You are a senior short video APP operation expert, and you need to deeply depict users' interests and explore their needs based on the ad titles that users have clicked on in the past. The following is the user's historical behavior (arranged in chronological order from oldest to most recent):

{user_historical_sequence}

Please analyze according to the following steps:

1. Carefully study the ad titles that the user has interacted with in history, comprehensively extract key information such as themes, products, services involved, accurately summarize the user's interests and hobbies, and avoid being too specific or fragmented in expression.
2. Based on the key information, conduct in-depth analysis and explore the user's long-term and short-term interests and hobbies.
3. Consider the user's purchasing preferences, such as whether they tend to buy specific types of products or services.
4. According to the above analysis, put forward suggestions for ad delivery that are in line with the user's interests and needs.
5. Finally, based on all your previous analysis content, depict the target user's interests as comprehensively as possible, and highlight the user characteristics that can be used for advertising and marketing.

It is prohibited to add any content that cannot be determined from user behavior. The answer shall be output in the following json format:

{"long-term interests": xx,"short-term interests": xx,"purchasing preferences": xx,"user needs": xx,"delivery suggestions": xx,"description of user interests": xx}

The following is your output:

Figure 11: The prompt used for User Interest Profiling.

You are a senior short video APP operation expert, proficient in consumer psychology, and skilled in generating ad titles that attracts users to click based on user profiles and ad information.

Your task is to rewrite an ad title that is more appealing to the user based on user profiles, historically clicked ad titles, user queries (which may not exist), original ad titles, and ad selling points.

By deeply understanding the content of user profiles, to flexibly integrate user information into the title to generate personalized content.

Requirements:

- Be completely based on the provided user information, ad titles, and ad selling points; do not introduce any other new information, and it is prohibited to fabricate content out of thin air or change the style of the original title.
- Select ad selling points that the user is more likely to be interested in and integrate them into the title.
- Fully consider user information, attract the user's interest as much as possible, and highlight the user in the first half of the title, but there is no need to forcefully integrate the user.
- Emphasize the content of the user query in the generated title.
- Add at most 20 additional characters to the original title; it is prohibited to include content related to the user's age or location, and the generated title should be vivid, smooth, and not stiff.

Please generate according to the following steps:

- Determine the information in the user profile that is suitable for integrating into the title.
- Extract selling points that the user may be more interested in (can be an empty string).
- Naturally integrate the above information into the title, output format: **{"User information suitable for integrating into the ad": "xx", "Selling points the user may be interested in": "xx", "Ad title combined with the user": "xx"}**

The following is the target user profile:

{user_profile}

The following are the ad titles that the user has historically clicked on:

{user_historical_sequence}

The following is the original ad title:

{original_title}

The following are the ad selling points:

{selling_points}

The following is the user query:

{user_query}

The following is your output:

Figure 12: The prompt used for Interest-Driven Title Generation.

This is a search advertising scenario. We will rewrite ad title using advertiser-provided ad information, selling points, user query, and user profile. However, the rewritten content may have hallucinations, i.e., new information that cannot be inferred from the original text. You do not need to identify some advertising marketing phrases as new information. Marketing content in the ad title tailored the user profile also does not need to be identified as new information, but it is necessary to exclude fabricated selling points introduced by the user profile.

Current user query term: **{user_query}**

Current user profile:
{user_profile}

The following is original ad information (title, and selling points):

Original Title:
{original_title}

Selling points:
{selling_points}

Please help me detect completely fabricated new information in the generated content.

The following is a generated ad title: **{generated_title}**

Does the generated ad title contain new information that cannot be inferred from the original text at all? Unless you are highly confident, uncertain cases should be regarded as rewriting hallucinations.

If yes, please provide it in json format: **{"new_information": [new_information1, new_information2, ...], "conflicting_information": [conflicting_information1, conflicting_information2, ...]}**. Otherwise, just output **{}**.

Do not include any other content. The following is your output:

Figure 13: The prompt used for Hallucination-Free Title Filtering and hallucination evaluation. Only the output "{}" is considered to pass the evaluation.

You are an expert in the field of video recommendation, highly skilled at capturing users' interest preferences.

Below, I will provide some video titles that the user has clicked on in the past (sorted in chronological order from oldest to most recent) and the user profile. Please deeply understand the user's interests based on these information.

I will also provide the selling points and title of the target advertising, as well as the user's current search query.

In addition, I will provide two titles, a and b. Please comprehensively consider all the information to judge which title the user is more likely to be interested in, and respond with (a, b, or same).

Input is as follows:

The user's query term is: **{user_query}**

The following are the advertising selling points:
{selling_points}

The following is the user's historical behavior:
{user_historical_sequence}

The following is the user profile:
{user_profile}

Title a: **{generated_title_a}**

Title b: **{generated_title_b}**

Please note that the output should only contain the answer (a, b, or same).

Output:

Figure 14: The prompt used for GSB performance evaluation.

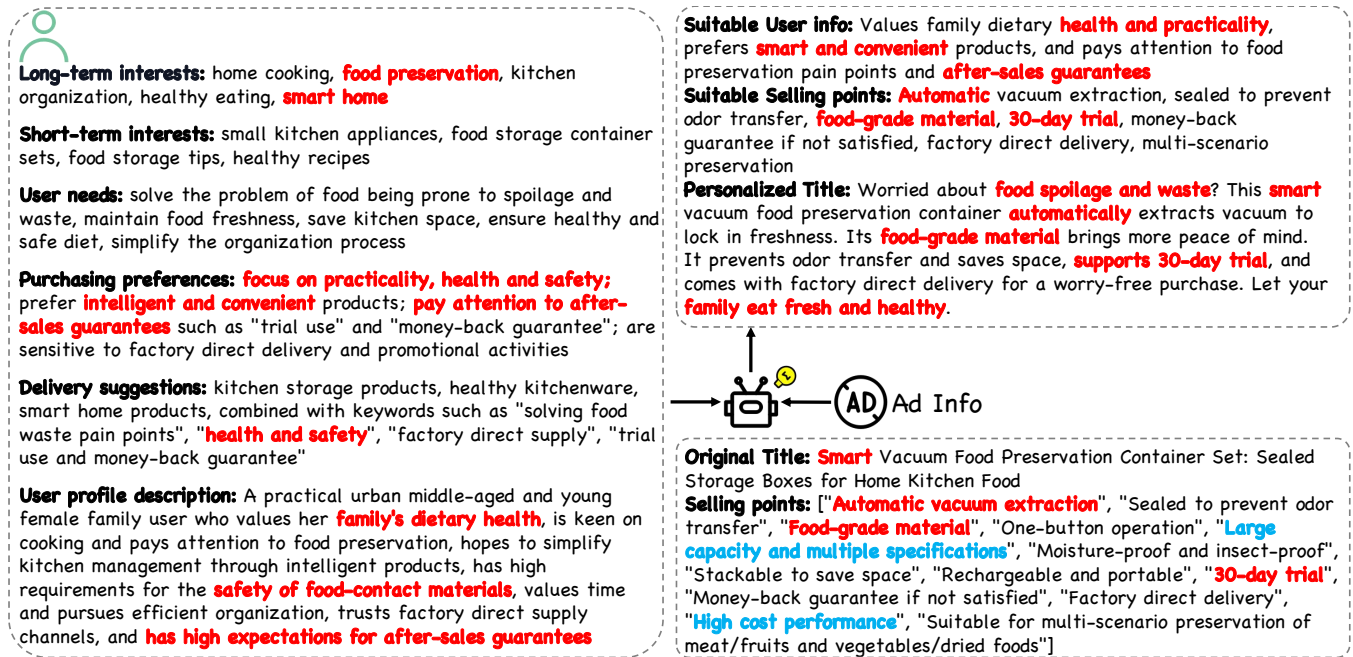


Figure 15: CoT-based personalized title, which extracts selling points more aligned with user interests such as "health", "smart", and "after-sales service". Red text indicates content that matches the user, while Blue text indicates content that does not match the user.

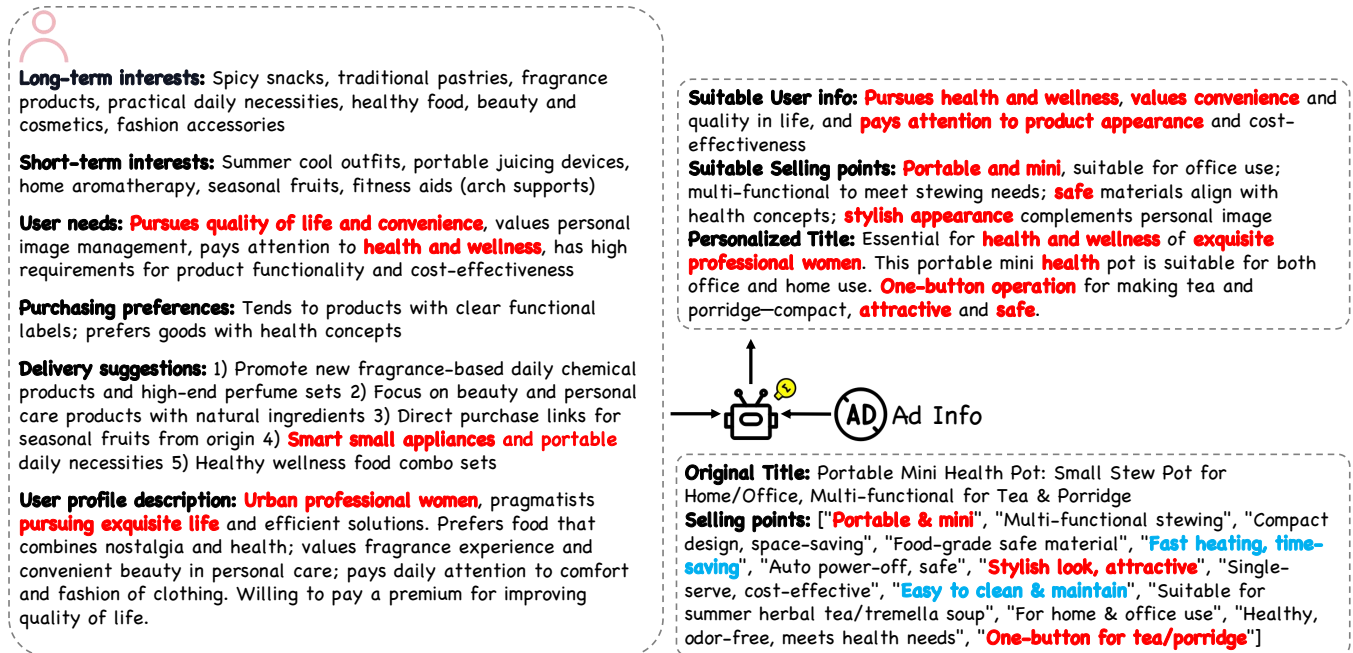


Figure 16: CoT-based personalized title, which explicitly emphasizes the user feature description "professional women". Red text indicates content that matches the user, while Blue text indicates content that does not match the user.