

# Echoes in Filter Bubble: Diagnosing and Curing Popularity Bias in Generative Recommenders

Jun Yin\*, Bangguo Zhu\*, Peng Huo, Ruochen Liu, Hao Chen, Senzhang Wang,  
Shirui Pan<sup>†</sup>, *Senior Member, IEEE* Chengqi Zhang<sup>†</sup>, *Fellow, IEEE*

**Abstract**—Recently, Generative Recommenders (GRs), characterized by a unified end-to-end framework, have exhibited astonishing potential in transforming the recommendation paradigm. Despite their effectiveness, we recognize that GRs are still susceptible to the long-standing issue of popularity bias that has pervaded the recommendation community. Although a few studies have attempted to extend traditional debiasing methods to GRs, their effectiveness is marginal, and the fundamental reason why GRs suffer from popularity bias remains under-explored. To bridge this gap, this study focuses on two core aspects in GRs: the optimization of generative framework and the item tokenization based on semantic index. Based on theoretical analyses, we identify that the severe popularity bias emerges from the confluence of a token-level optimization flaw and the undifferentiated property of item tokenization. Accordingly, this study develops a novel generative recommender system, called Ghost, by designing the asymmetric unlikelihood optimization and the skeleton-founded tokenization. Extensive empirical evaluations across three datasets, alongside multiple SOTA baselines, reveal that Ghost substantially alleviates popularity bias and promotes fairer recommendations, while incurring slight degradation to the overall recommendation utility.

**Index Terms**—Generative Recommender Systems, Popularity Bias, Large Language Models.

## I. INTRODUCTION

RECOMMENDER systems, due to the ability to capture user preferences by analyzing historical interactions, are essential for enhancing user experiences across platforms such as e-commerce [1], video streaming [2], and social networks [3]. Recently, generative recommenders (GRs) have emerged as a transformative paradigm [4]–[8] by replacing traditional item IDs [9], [10] with semantic indices (SIDs) [11], [12]. Driven by the adoption of large language models (LLMs) [13]–[15] as the backbone, the architecture of state-of-the-art (SOTA) GRs appears to become increasingly well-established, ushering in an era of rapid advancement.

However, despite their significant improvements in recommendation performance, we identify that current leading

GRs remain plagued by the persistent problem of popularity bias [16]–[18], which has long affected the recommender systems field [19]. As illustrated in Figure 1, three state-of-the-art (SOTA) GRs, LETTER [6], LC-Rec [5], and ED<sup>2</sup> [7], all exhibit an extreme over-recommendation phenomenon for popular head items while struggling to accurately predict the niche tail items. In particular, Figure 1a) showcases the significant gap between the performance of head items and tail items, which is up to 45 times greater for LC-Rec. Figure 1b) presents the number of head and tail items in recommendation lists and reveals an obvious preference for popular items, which account for more than 97 percent of the recommendation lists. Furthermore, according to Figure 1c), as the backbone scale of LC-Rec increases from 0.6B to 8B, the popularity bias is not mitigated spontaneously. Consequently, the popularity bias in GRs precipitates a filter bubble, wherein trending content monopolizes visibility while high-quality, long-tail items are severely marginalized.

Although a few preliminary studies [19] have attempted to extend conventional methods of popularity debiasing, such as item re-weighting [20], [21] and result re-ranking [22], [23], to GRs, their effectiveness is relatively marginal. As illustrated by Figure 2, existing methods (i.e., IFair-RW and IFair-RR) [19] still struggle to achieve appropriate *Pareto Optimality* [24], [25] while collectively considering the overall recommendation performance, tail item recommendation performance, and recommendation fairness. A non-negligible distance persists from the ideal region. On the other hand, the fundamental reason why GRs suffer from severe popularity bias is still under-explored. Compared with the discriminative paradigm [9], [10], [26], two distinctive characteristics of GRs are (i) the end-to-end generative framework and (ii) the item tokenization based on SIDs. However, these two distinct aspects are completely overlooked when designing popularity debiasing methods for GRs, rendering them devoid of strategic principles.

To bridge the current gap, this study is anchored in the optimization dynamics of the generative framework and the SID structures under current tokenization. *Firstly*, regarding the optimization based on maximum likelihood estimation (MLE), which is widely adopted by current GRs [4], [5], [7], we start with a gradient analysis and identify that the SID tokens of tail items mostly suffer from a gradient starvation issue. Due to the heavily long-tailed training distribution [16], [27], the tail item tokens are pathologically pushed away from user preference. Hence, during the recommendation process, when head and tail items compete, the optimization flaw leads

Jun Yin and Chengqi Zhang are with the Department of Data Science and Artificial Intelligence, Hong Kong Polytechnic University, Hong Kong SAR, China. Email: Junmay.yin@connect.polyu.hk, Chengqi.zhang@polyu.edu.hk

Bangguo Zhu, Ruochen Liu, and Senzhang Wang are with the School of Computer Science and Engineering, Central South University, Changsha, China. Email: {8210231132, ruochen, szwang}@csu.edu.cn

Peng Huo is with the National Super Computing Center, Tianjin, China. Email: huopeng@nssc-tj.cn

Hao Chen is with the Faculty of Data Science, City University of Macau, Macau, China. Email: haochen@cityu.edu.mo

Shirui Pan is with the School of Information and Communication Technology, Griffith University, Brisbane, Australia. Email: s.pan@griffith.edu.au

\* Equal Contribution. † Corresponding Authors.

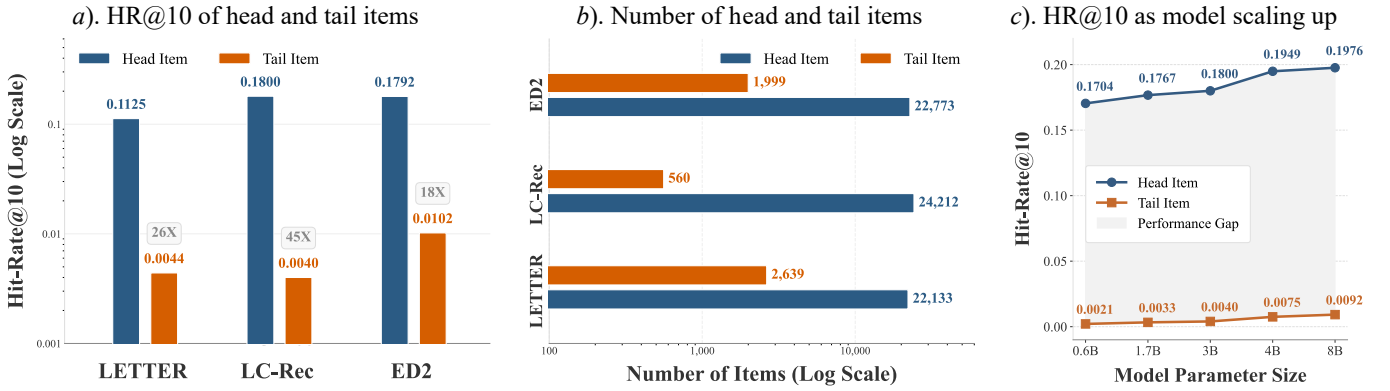


Fig. 1. *a*). Comparison of Hit-Rate@10 (i.e., HR@10) between head and tail items. *b*). Comparison between the number of head and tail items in the recommendation list provided by three GRs. *c*). Tendency of HR@10 as the backbone parameters of LC-Rec scaling up.

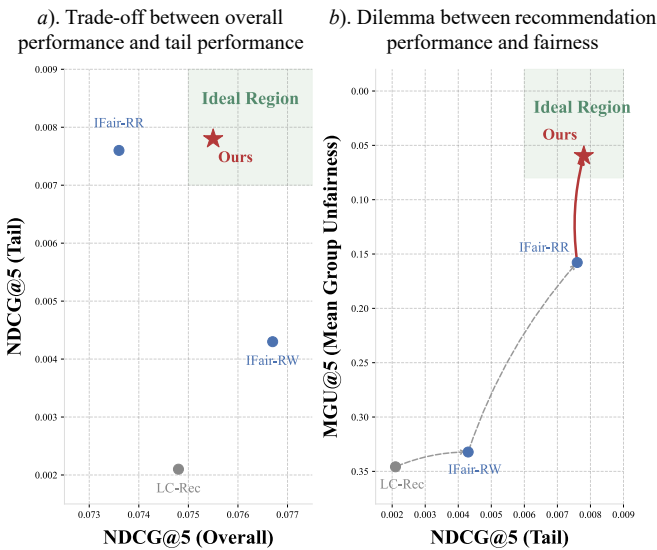


Fig. 2. Limitations of current popularity debiasing methods on GRs. *a*). Trade-off between overall recommendation performance and that of tail items. *b*). Dilemma between performance and fairness of recommendation results.

to the dominance of head item tokens. *Secondly*, this study sheds light on the fact that the current item tokenization is undifferentiated for head and tail items, without accounting for disparities in item popularity. It induces unpredictable item competition and continuously amplifies the token-level bias between head and tail items. Eventually, the undifferentiated tokenization renders the probability of tail items severely hijacked by their popular counterparts, and thus the whole GR model struggles to cast off the popularity bias.

Based on the crucial insights above, we develop a novel generative recommender named Ghost<sup>1</sup>, which is equipped with the asymmetric unlikely optimization (AUO) and the skeleton-founded tokenization (SKT). In detail, to mitigate the gradient starvation issue, AUO introduces asymmetric token-level unlikely between the tail items and their head counterparts. By constructing a reasonable undesired collection, AUO succeeds in rescuing the ineffective gradients of tail item

<sup>1</sup>Ghost denotes a Generative recommender with asymmetric unlikely optimization and skeleton-founded tokenization.

tokens. To inhibit the bias amplification effect caused by undifferentiated tokenization, SKT first designates the head item SIDs as the skeleton of the whole SID system. Subsequently, tail item SIDs are cultivated along the skeleton structure, with a dedicated focus on capturing the distinctiveness between tail and head items. The main contributions of this study are summarized as follows.

- This study diagnoses the troublesome issue of popularity bias in generative recommenders and identifies two fundamental factors, (i) the gradient starvation problem of tail item tokens and (ii) the bias amplification effect of undifferentiated item tokenization.
- This study develops Ghost, a novel GR that exhibits resilience to popularity bias, by designing the asymmetric unlikely optimization (AUO) and the skeleton-founded item tokenization (SKT). Specifically, AUO integrates effective negative feedback to calibrate the optimization of tail item tokens. Based on the skeleton structure, SKT ingeniously reduces disordered head-tail competition and characterizes the distinctiveness of tail items.
- Extensive empirical evaluations across three datasets, alongside multiple SOTA baselines, reveal that Ghost substantially alleviates popularity bias and promotes fairer recommendations, ultimately achieving the desired Pareto optimality of generative recommenders.

## II. PRELIMINARY

**Problem Formulation.** This study focuses on the *sequential recommendation* task [5], [7], [10], which aims to predict the next most suitable item based on the user historical behavior. Considering a system of  $K$  items  $\{v_k | k = 1, 2, \dots, K\}$  and  $J$  users  $\{u_j | j = 1, 2, \dots, J\}$ , the behavior of  $u_j$  can be represented by an item sequence  $h_{u_j} = [v_{k_1}, v_{k_2}, \dots, v_{k_l}]$ , where  $l$  is the sequence length.

**Generative Recommender Systems.** Existing efforts towards developing GRs [4]–[7], [28] primarily concentrate on the paradigm based on SIDs. Generally, SID-based GRs consist of two main components, i.e., item tokenization and recommendation-oriented finetuning [5], [7]. For item tokenization, SID-based GRs usually introduce vector quantization techniques, such as VQ-VAE [11], RQ-VAE [12], and

RQ-KMeans [29], to convert the continuous embeddings of items into discrete indices. Taking RQ-VAE as an example, and assuming that  $X_v$  is the item embedding and the collection  $\{\mu_n^{(i)}\}_{n=1}^N$  denotes the  $i$ -th codebook within the RQ-VAE, the SID generation process can be represented as,

$$\begin{aligned} c_v^{(i)} &= \arg \min_{n \in \{1, 2, \dots, N\}} \|r_v^{(i)} - \mu_n^{(i)}\|_2^2, \\ r_v^{(i+1)} &= r_v^{(i)} - \mu_{c_v^{(i)}}^{(i)}, \quad \text{for } i = 1, 2, \dots, L. \end{aligned} \quad (1)$$

where  $r_v^{(1)} = X_v$  and  $L$  denotes the SID length. Afterwards, item  $v$  can be indexed as an SID  $\Omega_v = (c_v^{(1)}, c_v^{(2)}, \dots, c_v^{(L)})$ . For recommendation-oriented finetuning, the sequential recommendation task can be reformulated as a language generation task based on SIDs. In particular, for historical behavior  $h_u$  and target item  $v$  with SID  $\Omega_v$ , the optimization objective of SID-based GRs follows Maximum Likelihood Estimation (MLE) and usually adopts the negative log-likelihood (NLL) loss, as follows,

$$\mathcal{L}_{\text{NLL}} = - \sum_i \log \mathcal{P}_\theta(c_v^{(i)} | h_u, c_v^{<i}), \quad (2)$$

where  $\theta$  denotes the GR parameters and  $c_v^{<i}$  denotes the subsequence of SID  $\Omega_v$  before the  $i$ -th position. *This study focuses on SID-based GRs, as they are currently the most influential paradigm with the greatest performance potential.*

**Popularity Bias.** In the domain of recommender systems, the interaction frequency of items often follows a long-tail distribution [16]–[18], [30], especially a power-law distribution. Popularity bias refers to the algorithmic preference to disproportionately favor frequently interacted short-head items at the expense of highly relevant but sparse long-tail items [16]. *This study follows a common practice [27] in popularity debiasing research, which groups the items into the head set (i.e., the top 20% most popular items) and the tail set (i.e., the remaining 80% of items) based on item popularity.*

**Pareto Optimality.** In multi-objective optimization, a model is simultaneously evaluated across  $M$  distinct, often conflicting metrics [24]. Assuming that all objectives  $\{f_m\}_{m=1}^M$  are to be minimized, a solution  $s$  is said to strictly precede another solution  $s'$ , denoted as  $s \prec s'$ , if and only if  $\forall m_1, f_{m_1}(s) \leq f_{m_1}(s')$  and  $\exists m_2, f_{m_2}(s) < f_{m_2}(s')$ . A solution  $s^*$  is defined as Pareto optimal if there exists no other feasible solution  $s$  such that  $s \prec s^*$ . Consequently, Pareto optimality represents the optimal trade-off among the multiple evaluation metrics, indicating that no single objective can be further improved without fundamentally degrading at least one other objective [24], [25].

### III. DIAGNOSE POPULARITY BIAS IN GRs

To investigate the popularity bias in current GRs, this study starts with an analysis of (i) the gradient of MLE optimization and (ii) the SID structure under undifferentiated tokenization.

#### A. Gradient Analysis of MLE Optimization

Under the MLE objective  $\mathcal{L}_{\text{NLL}}$  defined in Eq.(2), the optimization conditioned on user historical behavior  $h_u$  is governed by the Softmax derivative. Let  $\mathcal{D}$  denote the training

distribution of user-item interaction pairs  $(h_u, v)$ , where the history  $h_u$  is encoded into a representation  $X_{h_u}$  and the target item  $v$  is tokenized as a SID  $\Omega_v = (c_v^{(1)}, c_v^{(2)}, \dots, c_v^{(L)})$ .

**LEMMA 1 (Gradient Starvation in MLE).** *For an arbitrary SID token  $c$ , whose embedding is  $e_c$ , the expected MLE gradient update  $\mathbb{E}_{\mathcal{D}}[\Delta e_c]$  over the training distribution  $\mathcal{D}$  conforms to*

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[\Delta e_c] \propto \mathbb{E}_{\mathcal{D}} \left[ \sum_i \left( \mathbb{I}\{c = c_v^{(i)}\} \cdot (1 - \mathcal{P}_\theta(c | h_u, c_v^{<i})) \right. \right. \\ \left. \left. - \mathbb{I}\{c \neq c_v^{(i)}\} \cdot \mathcal{P}_\theta(c | h_u, c_v^{<i}) \right) \cdot X_{h_u} \right]. \end{aligned} \quad (3)$$

*However, for a tail token  $c_{\text{tail}}$  that predominantly composes tail items, the expected gradient update is heavily skewed in the negative direction, formulated as,*

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[\langle \Delta e_{c_{\text{tail}}}, X_{h_u} \rangle] \\ \approx -\mathbb{E}_{\mathcal{D}} \left[ \sum_i \mathcal{P}_\theta(c_{\text{tail}} | h_u, c_v^{<i}) \cdot \|X_{h_u}\|_2^2 \right] \leq 0. \end{aligned} \quad (4)$$

In a long-tailed distribution  $\mathcal{D}$ , tokens composing head items frequently act as the targets, receiving massive positive gradient updates that align the embedding  $e_c$  with user preference  $X_{h_u}$ . Conversely, tokens that are specific to tail items mainly act as trivial negative samples in the denominator of the Softmax operation. This leads to the *Gradient Starvation* issue [31] for tail item tokens, in which the token embeddings are consistently pushed away from the user preference space.

#### B. SID Branching Points under Undifferentiated Tokenization

Subsequently, we further investigate how this token-level optimization flaw propagates into item-level bias during the recommendation process. Let  $c_{\text{head}}^{(i)}$  and  $c_{\text{tail}}^{(i)}$  be candidate tokens competing at the  $i$ -th item generation step, and step  $i$  is denoted as a branching point. Due to the asymptotically repulsive gradient updates stemming from gradient starvation, the generation probability of the tail token  $c_{\text{tail}}^{(i)}$  is systematically penalized, resulting in a bias amplification factor.

**COROLLARY 1 (Head Token Dominance at Branching Point).** *At step  $i$ , where head and tail paths compete, the token generation probability ratio diverges from the true data distribution  $\mathcal{P}_d$  by a local amplification factor  $\gamma_i > 1$ :*

$$\frac{\mathcal{P}_\theta(c_{\text{head}}^{(i)} | h_u, c^{<i})}{\mathcal{P}_\theta(c_{\text{tail}}^{(i)} | h_u, c^{<i})} = \gamma_i \cdot \frac{\mathcal{P}_d(c_{\text{head}}^{(i)} | h_u, c^{<i})}{\mathcal{P}_d(c_{\text{tail}}^{(i)} | h_u, c^{<i})}. \quad (5)$$

Eq.(5) implies that the generation process becomes pathologically overconfident in head tokens, regardless of the current context (i.e.,  $h_u$  and  $c^{<i}$ ). Furthermore, current GRs mostly tokenize items into SIDs in an undifferentiated style, which treats items identically without accounting for their inherent popularity disparities. The undifferentiated nature induces that there is no predictable structural branching point between head and tail items. Within the undifferentiated tokenization, a tail item  $v_{\text{tail}}$  shares prefixes of varying lengths with numerous popular items and is bound to encounter a sequence of branching points with head token dominance.

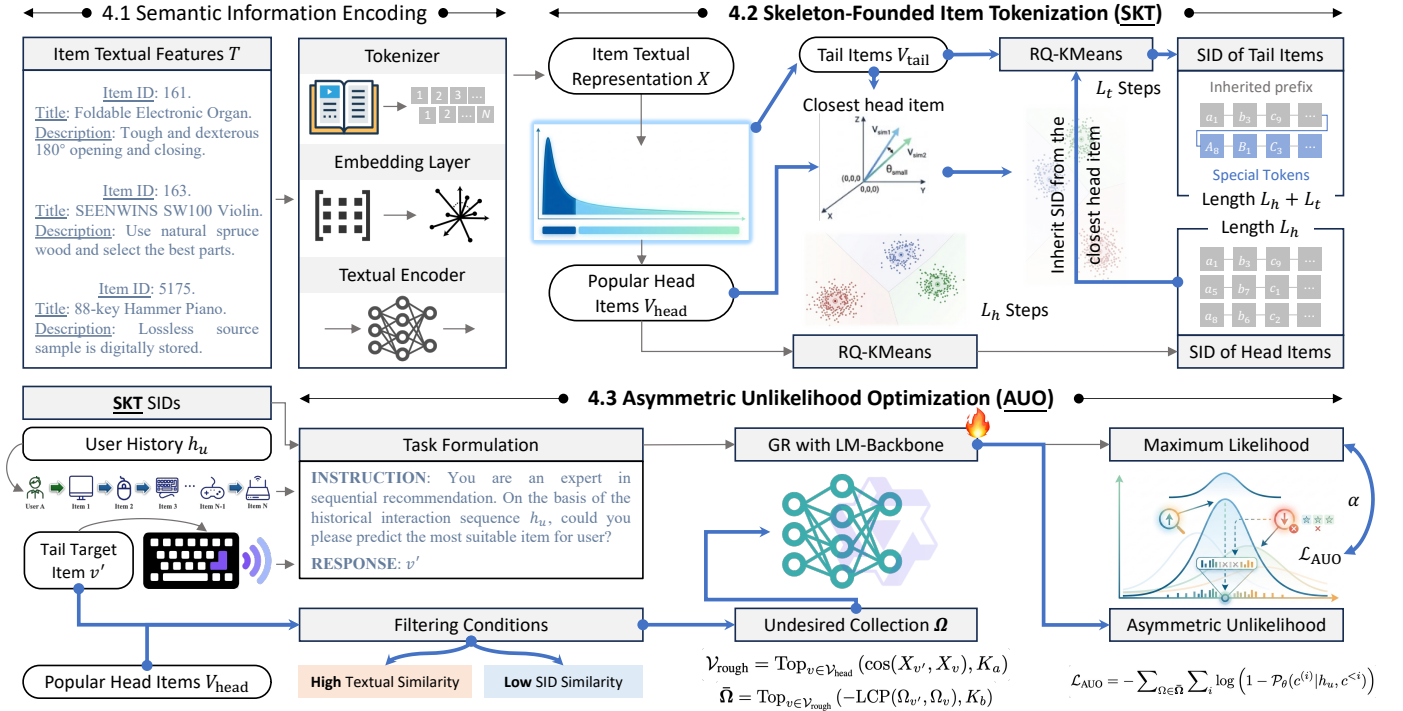


Fig. 3. Overview of the Ghost model. First, textual representations are encoded based on item features. After categorizing items into head and tail sets, SKT assigns SIDs via RQ-KMeans, allowing tail items to inherit prefixes from their closest head items. At last, AUO optimizes the GR model by penalizing an undesired collection customized for each tail item, alongside standard MLE training.

**LEMMA 2 (Bias Amplification via Undifferentiated Tokenization).** *Let  $\mathcal{Z}$  (where  $|\mathcal{Z}| = z \leq L$ ) be the set of steps along the generation process of  $v_{\text{tail}}$ , where it must compete against head candidate tokens. The probability of successfully navigating these intersections without being hijacked by head items is geometrically suppressed during recommendation,*

$$\begin{aligned} \mathcal{P}_{\theta}(v_{\text{tail}} | h_u) &= \prod_j \mathcal{P}_{\theta}(c_{\text{tail}}^{(j)} | h_u, c^{<j}) \\ &\leq (\gamma_{\min})^{-z} \cdot \prod_j \mathcal{P}_d(c_{\text{tail}}^{(j)} | h_u, c^{<j}). \end{aligned} \quad (6)$$

where  $\gamma_{\min} = \min_{j \in \mathcal{Z}}(\gamma_j) > 1$ . Therefore, the probability ratio of generating a competing head item over the tail item  $v_{\text{tail}}$  cascades geometrically by a factor of at least  $\mathcal{O}(\gamma_{\min}^z)$ .

**Discussion.** According to the analyses above, the severe popularity bias emerges from the confluence of (i) gradient starvation in MLE optimization and (ii) bias amplification of undifferentiated tokenization. Detailed derivations are presented in Appendix F. Intuitively, under MLE optimization, tail item tokens fail to receive effective gradient signals, trapping them in gradient starvation. Then, the gradient starvation of tail item tokens results in the head token dominance at branching point. Finally, during the recommendation process, a sequence of head-dominated branching points are brought by the undifferentiated item tokenization, leading to the severe popularity bias.

## IV. METHODOLOGY

With a focused effort to the two diagnosed fundamental factors, the Ghost model is developed by designing the

skeleton-founded item tokenization (SKT) and the asymmetric unlikelihood optimization (AUO). Procedurally, the overview of Ghost is illustrated in Figure 3. In particular, SKT subtly reduces unpredictable branching points by formulating a mechanism for SID inheritance from head to tail items. Hence, the bias amplification effect of undifferentiated tokenization is suppressed. Afterwards, based on the undesired collection of head item tokens, AUO is able to reallocate the supervision signals and thus rescue the gradient starvation issue of tail item tokens.

### A. Skeleton-Founded Item Tokenization

Current GRs [4], [5], [7] mostly adopt standard vector quantization techniques, such as RQ-VAE [12] and RQ-KMeans [29], to allocate item SIDs. These standard approaches are agnostic to item popularity disparities, where head and tail items are indiscriminately processed. Therefore, the generated SIDs contain unstructured branching points where head and tail item tokens compete. To overcome this limitation, we propose the skeleton-founded item tokenization (SKT). In detail, SKT circumvents the chaotic competition between head and tail item tokens by specifying the position of the branching point. This approach makes it possible to capture the occurring patterns of tail items without having them drowned out by their popular counterparts.

To be more specific, SKT assigns the SIDs for head and tail items asynchronously. Firstly, given the head item set  $V_{\text{head}}$  and the corresponding representations, SKT adopts the RQ-KMeans algorithm to generate the head item SIDs. For head

item  $v$  with representation  $X_v$ , the SID generation follows,

$$\begin{aligned} c_v^{(i)} &= \arg \min_{n \in \{1, 2, \dots, N\}} \|r_v^{(i)} - \mu_n^{(i)}\|_2^2, \\ r_v^{(i+1)} &= r_v^{(i)} - \mu_{c_v^{(i)}}^{(i)}, \quad \text{for } i = 1, 2, \dots, L^h. \end{aligned} \quad (7)$$

$L^h$  is the SID length for head items. Initialized with  $r_v^{(1)} = X_v$ , the head item  $v$  is assigned SID  $\Omega_v = (c_v^{(1)}, c_v^{(2)}, \dots, c_v^{(L^h)})$ . Afterwards, for the tail item  $v'$  with representation  $X_{v'}$ , SKT begins by retrieving the head item with the highest semantic similarity, denoted as  $v^*$ . Subsequently, the tail item  $v'$  inherits the first  $L^h$  tokens from  $\Omega_{v^*}$  and then additionally generates  $L^t$  SID tokens. Similar to Eq.(7), the additional SID generation for tail item  $v'$  is represented below,

$$\begin{aligned} c_{v'}^{(j)} &= \arg \min_{n \in \{1, 2, \dots, N\}} \|r_{v'}^{(j)} - \mu_n^{(j)}\|_2^2, \\ r_{v'}^{(j+1)} &= r_{v'}^{(j)} - \mu_{c_{v'}^{(j)}}^{(j)}, \quad \text{for } j = 1, 2, \dots, L^t. \end{aligned} \quad (8)$$

Differing from the SID generation of head items, Eq.(8) is initialized with

$$r_{v'}^{(1)} = X_{v'} - \sum_i \mu_{c_{v^*}^{(i)}}^{(i)},$$

and thus the tail item  $v'$  is indexed with SID  $\Omega_{v'} = (c_{v^*}^{(1)}, \dots, c_{v^*}^{(L^h)}, c_{v'}^{(1)}, \dots, c_{v'}^{(L^t)})$ .

Conceptually, the SID collection of head items  $\{\Omega_v | v \in \mathcal{V}_{\text{head}}\}$  serves as the skeleton of the whole SID system and essentially tessellates the SID space into several balanced partitions. On one hand, by inheriting the tail item SID from that of the closest head item, SKT maintains the correlation between items. Therefore, the basic principle of SID-based GRs, that similar items are prone to share similar SIDs, still holds. On the other hand, SKT unifies the branching point between head and tail item at step  $(L^h + 1)$ , subtly reducing the bias amplification effect of undifferentiated item tokenization. Moreover, regarding the expressiveness of our SIDs, SKT captures the distinctiveness of the tail item  $v'$  compared to its closest head item  $v^*$  through  $L^t$  additional SID tokens, which elegantly leverages the strong modeling capability on head items towards tail item improvement.

### B. Asymmetric Unlikelihood Optimization

Based on the novel SIDs generated by SKT, Ghost further designs the asymmetric unlikelihood optimization (AUO) to rectify the gradient starvation issue of tail item tokens. As analyzed in Section III-A, the optimization of current GRs relies on MLE, where the tail items seldom serve as right answers. As a result, the optimized GR models prefer repeating head items, since they are never explicitly penalized for over-estimating statistically popular but contextually inappropriate head item tokens.

Inspired by the philosophy of unlikelihood training [32], [33], AUO modulates the supervision signals by actively imposing penalties if the GR model generates notoriously undesirable tokens. Opposite to the standard NLL loss  $\mathcal{L}_{\text{NLL}}$ , the AUO loss  $\mathcal{L}_{\text{AUO}}$  can be concisely defined as follows,

$$\mathcal{L}_{\text{AUO}} = - \sum_{\Omega \in \bar{\Omega}} \sum_i \log \left( 1 - \mathcal{P}_\theta(c^{(i)} | h_u, c^{<i}) \right), \quad (9)$$

where  $\bar{\Omega}$  denotes the SID collection of undesired items. Notably, the AUO loss  $\mathcal{L}_{\text{AUO}}$  is customized for tail items rather than the whole item set. For tail item  $v'$  with representation  $X_{v'}$  and SID  $\Omega_{v'}$ , Ghost first retrieves the top similar items from the popular head collection  $\mathcal{V}_{\text{head}}$ , which functions as a rough selection of undesired items. Then, Ghost further refines the rough candidates  $\mathcal{V}_{\text{rough}}$  according to their SID distance from  $\Omega_{v'}$ , corresponding to a finer filtering of notoriously undesired items. Formally, the construction of  $\bar{\Omega}$  can be represented as follows,

$$\begin{aligned} \mathcal{V}_{\text{rough}} &= \text{Top}_{v \in \mathcal{V}_{\text{head}}} (\cos(X_{v'}, X_v), K_a), \\ \bar{\Omega} &= \text{Top}_{v \in \mathcal{V}_{\text{rough}}} (-\text{LCP}(\Omega_{v'}, \Omega_v), K_b), \end{aligned} \quad (10)$$

where the  $\text{Top}(f(\cdot, \cdot), K)$  operator returns the  $K$  objects with the highest  $f$  function value, the  $\text{LCP}(\cdot, \cdot)$  operator returns the length of the *Longest Common Prefix* between the two input sequences, and  $K_a, K_b$  are two hyper-parameters controlling the candidate scale of the undesired item selection.

For tail item  $v'$ , the undesired collection  $\bar{\Omega}_{v'}$  includes several head items whose representations are similar to  $X_{v'}$ , while their SIDs are divergent from  $\Omega_{v'}$ . Hence, AUO is able to inhibit the GR model from being hijacked by the popular head items that are similar to the ground-truth tail item in terms of semantic representation. In summary, the specialized asymmetric unlikelihood optimization is combined with MLE optimization, and the overall objective function of Ghost is defined as,

$$\mathcal{L}_{\text{All}} = \mathcal{L}_{\text{NLL}} + \alpha \cdot \mathcal{L}_{\text{AUO}}, \quad (11)$$

where  $\alpha$  is the weighted parameter. Implementation details of Ghost are introduced in Appendix C.

### C. Theoretical Analysis

Within SKT, since a tail item  $v'$  inherits the prefix skeleton of its closest head item  $v^*$ , the branching point of item recommendation is uniformly deferred to step  $(L^h + 1)$ . By establishing a singular, predictable branching point, SKT exponentially restricts the local amplification factor in LEMMA 2.

**LEMMA 3 (Mitigation of Bias Amplification).** *The generative probability of the tail item  $v'$  is insulated from the multi-step geometric suppression, and its deviation from the true data distribution  $\mathcal{P}_d$  is governed exclusively by a localized factor,*

$$\begin{aligned} \mathcal{P}_\theta(v' | h_u) &= \prod_{j=1}^{L^h + L^t} \mathcal{P}_\theta(c_{v'}^{(j)} | h_u, c^{<j}) \\ &\approx (\gamma_{\text{EOS}})^{-1} \cdot \prod_{j=1}^{L^h + L^t} \mathcal{P}_d(c_{v'}^{(j)} | h_u, c^{<j}). \end{aligned} \quad (12)$$

where  $\gamma_{\text{EOS}} \geq 1$  represents the head-dominance factor confined to the  $(L^h + 1)$ -th generative recommendation step, in which the tokens of tail item  $v'$  compete exclusively against the EOS token.

Contrasting Eq.(12) with Eq.(6), we can notice that by transforming a geometric suppression  $\mathcal{O}(\gamma_{\min}^z)$  into a single-step discrepancy  $\mathcal{O}(\gamma_{EOS})$ , SKT effectively halts the amplification of popularity bias during the generative recommendation process. Regarding the overall objective defined in Eq.(11), a comprehensive gradient analysis is conducted to uncover its rationale. While MLE subjects tail tokens to gradient starvation through the suppressive  $-\mathcal{P}_\theta(c_{\text{tail}}^{(i)}|h_u, c^{<i}) \cdot X_{h_u}$  term, the Ghost model provides a powerful *Rescue Force* on the basis of AUO.

**LEMMA 4 (Gradient Rescue based on AUO).** *For false positive head tokens  $c_{\text{head}}^- \in \bar{\Omega}$ , the gradient correctly cancels the denominator via the softmax derivative for its own direct penalty. Crucially, it also systematically absorbs the cross-penalization dynamics originating from other tokens within the undesired set. Formally, the gradient is,*

$$\begin{aligned} \Delta e_{c_{\text{head}}^-} &\propto - \underbrace{(1 + \alpha) \mathcal{P}_\theta(c_{\text{head}}^-) \cdot X_{h_u}}_{\text{Targeted Repulsion}} \\ &+ \alpha \underbrace{\sum_{c_j \in \bar{\Omega} \setminus \{c_{\text{head}}^-\}} \frac{\mathcal{P}_\theta(c_j) \mathcal{P}_\theta(c_{\text{head}}^-)}{1 - \mathcal{P}_\theta(c_j)} \cdot X_{h_u}}_{\text{Cross-Penalization Offset}}. \end{aligned} \quad (13)$$

For a token  $c_{\text{tail}}$  of the target tail item, the expected update gradient is reformulated as,

$$\begin{aligned} \Delta e_{c_{\text{tail}}} &\propto - \underbrace{\mathcal{P}_\theta(c_{\text{tail}}) \cdot X_{h_u}}_{\text{MLE Push}} \\ &+ \alpha \underbrace{\sum_{c_j \in \bar{\Omega}} \frac{\mathcal{P}_\theta(c_j) \mathcal{P}_\theta(c_{\text{tail}})}{1 - \mathcal{P}_\theta(c_j)} \cdot X_{h_u}}_{\text{AUO Rescue}}. \end{aligned} \quad (14)$$

By actively penalizing the undesired head tokens in  $\bar{\Omega}$ , the Jacobian of the AUO loss systematically redistributes positive probability mass to the remaining tokens. This structural rescue ensures that tail tokens receive active, rational parameter updates, preventing them from being indefinitely pushed away from the user intent space as mere trivial negative samples. Detailed derivations of LEMMA 4 and LEMMA 5 are presented in Appendix F as well.

## V. EXPERIMENT

• **Dataset.** This study evaluates Ghost and baselines on three public benchmarks extracted from the Amazon platform, i.e., "*Musical Instruments*", "*Arts, Crafts and Sewing*", and "*Video Games*".

• **Baseline.** The evaluation introduces three SOTA GRs (LETTER [6], LC-Rec [5], and ED<sup>2</sup> [7]) as the standard baselines. IFairLRS [19], which pioneers the investigation of popularity bias in GRs, serves as a competitive baseline. In detail, *RW* and *RR* refer to the *re-weighting only* and *re-weight & re-ranking* strategies. Moreover, to provide a comprehensive comparison, this study replaces each popular head item with the most similar tail item according to a pre-defined probability. The modified interaction sequence is either appended to the original training set or substituted for the original sequence, corresponding to the baselines denoted as

*Augmentation* and *Substitution* (The details of baselines are presented in Appendix C-B).

• **Metric.** Following well-established benchmarks [5], [34], this study adopts Hit-Rate (HR) and normalized discounted cumulative gain (NDCG) to evaluate the recommendation performance. Besides, this study adopts mean group unfairness (MGU) and average recommendation popularity (ARP) to quantify the fairness and the average popularity of the recommendation results [27]. To reflect the proximity to Pareto optimality, this study defines the comprehensively normalized score (CNS) as follows, based on the Min-Max normalized values of HR, NDCG, MGU, and ARP metrics,

$$\begin{aligned} \text{CNS} = & \frac{(\overline{\text{HR}}_{\text{All}} + \overline{\text{HR}}_{\text{Tail}} + \overline{\text{NDCG}}_{\text{All}}}{+ \overline{\text{NDCG}}_{\text{Tail}} + \overline{\text{MGU}} + \overline{\text{ARP}}) / 6. \end{aligned} \quad (15)$$

The detailed introduction of the datasets and metrics is presented in Appendix C.

### A. Main Result

To demonstrate the effectiveness of Ghost in mitigating popularity bias, we compare its performance against SOTA baselines across three public datasets. The results are summarized in Table I.

In general, Ghost can effectively approach the Pareto frontier while jointly considering overall recommendation performance, tail recommendation performance, and recommendation fairness. In detail, we can draw the following three observations. **First**, Ghost is primarily characterized by its significant performance gains in tail-item recommendation. Across three datasets, Ghost outperforms standard GRs (i.e., LETTER, LC-Rec, and ED<sup>2</sup>) by increasing tail HR and NDCG by 63.91% and 70.66% on average. Compared to existing popularity debiasing methods, it delivers an average increase of 28.39% and 15.04% for tail HR and NDCG, with maximum improvements up to 57.24%. **Second**, Ghost effectively suppresses the over-recommendation of head items. It lowers the MGU by an average of 55.76% compared to standard GRs, and outperforms the most competitive popularity debiasing baseline, IFairLRS-RR, by further reducing the MGU by 16.68% on average (up to a maximum of 66.60%). **Third**, with respect to overall recommendation performance, Ghost exhibits an acceptable degradation compared to standard models. Relative to the strongest baseline, LC-Rec, the overall HR and NDCG of Ghost drop by only 7.46% and 6.34% on average. For context, the performance drops among all popularity debiasing methods in Table I range between 2.46% and 18.90% for HR, and between 2.24% and 21.28% for NDCG. Taken together, Ghost yields an average CNS gain of 16.81%, peaking at 22.06%, which indicates its capability to substantially approach the Pareto optimal state for the GRs popularity bias issue.

### B. Ablation Study

To investigate the contributions of the core components in Ghost, an ablation study is conducted by removing specific modules. Based on Table II, the following insights can be

TABLE I  
PERFORMANCE COMPARISON OF GHOST AND BASELINES ACROSS THREE DATASETS. THE BEST, THE RUNNER-UP, AND THE THIRD-BEST RESULTS ARE HIGHLIGHTED IN **BOLD**, UNDERLINED, AND COLORED FONTS.

Dataset	Metric Model	HR@5 $\uparrow$		HR@10 $\uparrow$		NDCG@5 $\uparrow$		NDCG@10 $\uparrow$		MGU@5 $\downarrow$	MGU@10 $\downarrow$	ARP@5 $\downarrow$	ARP@10 $\downarrow$	CNS@5 $\uparrow$	CNS@10 $\uparrow$
		All	Tail	All	Tail	All	Tail	All	Tail						
Ins	LETTER	0.0593	0.0025	0.0662	0.0044	0.0530	0.0016	0.0553	0.0023	0.1965	0.1442	<b>181.2628</b>	<b>151.7149</b>	0.2457	0.2714
	LC-Rec	0.0870	0.0027	0.1046	0.0040	0.0748	0.0021	0.0804	0.0025	0.3458	0.3380	365.4217	338.5304	0.3214	0.3134
	ED <sup>2</sup>	<b>0.0898</b>	0.0067	<u>0.1068</u>	0.0102	<u>0.0765</u>	0.0043	<u>0.0820</u>	0.0055	0.2551	0.2305	279.8572	237.9431	0.6067	0.6217
	Augmentation	0.0875	0.0071	0.1039	0.0105	0.0759	0.0051	0.0812	0.0061	0.2383	0.2049	271.1657	242.7466	0.6327	<u>0.6323</u>
	Substitution	0.0780	0.0058	0.0914	0.0084	0.0645	0.0039	0.0689	0.0048	<b>0.0310</b>	<b>0.0114</b>	253.3299	229.5070	0.5714	0.5601
	IFairLRS-RW	<u>0.0888</u>	0.0058	<b>0.1072</b>	<u>0.0081</u>	<b>0.0767</b>	0.0043	<b>0.0826</b>	0.0050	0.3322	0.3226	328.5940	299.5196	0.5007	0.4887
	IFairLRS-RR	0.0853	<u>0.0100</u>	0.1053	<u>0.0110</u>	0.0736	<b>0.0078</b>	0.0800	<u>0.0081</u>	0.1578	0.2330	301.2252	286.8023	<u>0.7462</u>	0.6286
	Ghost (Ours)	0.0864	<b>0.0117</b>	0.1017	<b>0.0173</b>	0.0755	<b>0.0078</b>	0.0805	<b>0.0097</b>	<u>0.0596</u>	<u>0.0958</u>	<u>248.8179</u>	<u>209.8877</u>	<b>0.8974</b>	<b>0.8694</b>
	LETTER	0.0448	0.0062	0.0521	0.0090	0.0378	0.0040	0.0401	0.0050	0.1768	0.1283	<u>85.6235</u>	<u>73.2413</u>	0.2108	0.2693
	LC-Rec	<b>0.0885</b>	0.0188	<b>0.1095</b>	0.0269	<b>0.0737</b>	0.0133	<b>0.0805</b>	0.0159	0.2992	0.2905	135.2274	119.2569	0.5518	0.5548
ED <sup>2</sup>	0.0796	0.0105	0.0993	0.0172	0.0656	0.0073	0.0719	0.0095	0.2960	0.2879	148.5051	129.2274	0.3331	0.3548	
Arts	Augmentation	0.0839	<u>0.0218</u>	0.1021	<u>0.0306</u>	0.0697	<u>0.0159</u>	0.0756	<u>0.0187</u>	<u>0.1572</u>	0.1299	112.3944	98.8012	<u>0.6946</u>	0.7258
	Substitution	<u>0.0732</u>	0.0215	0.0915	0.0305	0.0581	0.0155	0.0640	0.0184	<b>0.0490</b>	<b>0.0531</b>	96.9228	88.0055	0.6910	<u>0.7269</u>
	IFairLRS-RW	<u>0.0845</u>	0.0138	<u>0.1022</u>	0.0191	<u>0.0706</u>	0.0102	<u>0.0763</u>	0.0119	0.3076	0.2925	141.4042	122.9239	0.4359	0.4270
	IFairLRS-RR	0.0827	0.0179	0.1005	0.0248	0.0696	0.0130	0.0753	0.0152	0.2255	0.2103	134.5895	117.1711	0.5447	0.5485
	Ghost (Ours)	0.0831	<b>0.0296</b>	0.1000	<b>0.0393</b>	<u>0.0706</u>	<b>0.0213</b>	0.0760	<b>0.0245</b>	0.1763	0.1621	<b>70.5222</b>	<b>65.1832</b>	<b>0.8828</b>	<b>0.8780</b>
	LETTER	0.0264	0.0063	0.0375	0.0096	0.0187	0.0044	0.0222	0.0055	0.1727	0.1396	<u>119.9021</u>	<u>108.3559</u>	0.2458	0.2652
LC-Rec	<b>0.0636</b>	0.0148	<b>0.0938</b>	0.0244	<b>0.0438</b>	0.0091	<b>0.0536</b>	0.0122	0.2907	0.2740	164.9579	150.6237	0.5352	0.5499	
ED <sup>2</sup>	0.0572	0.0086	0.0854	0.0150	0.0393	0.0052	0.0484	0.0072	0.3110	0.2925	188.5855	170.8846	0.3052	0.3284	
Games	Augmentation	0.0559	0.0173	0.0822	0.0282	0.0384	0.0113	0.0469	0.0148	0.1526	0.1210	138.5779	125.0857	0.6701	0.6978
	Substitution	0.0448	0.0199	0.0683	0.0311	0.0309	0.0131	0.0384	0.0167	<u>0.0783</u>	<u>0.0824</u>	121.4722	111.5519	0.7074	0.7093
	IFairLRS-RW	<u>0.0605</u>	0.0134	<u>0.0907</u>	0.0228	<u>0.0416</u>	0.0087	0.0513	0.0117	0.2869	0.2710	166.8988	153.0392	0.4878	0.5103
	IFairLRS-RR	0.0581	<u>0.0229</u>	0.0877	<u>0.0315</u>	0.0402	<u>0.0153</u>	<u>0.0497</u>	<u>0.0181</u>	0.1051	0.1449	148.2428	141.5092	<u>0.8073</u>	<u>0.7229</u>
	Ghost (Ours)	0.0562	<b>0.0257</b>	0.0840	<b>0.0392</b>	0.0390	<b>0.0167</b>	0.0479	<b>0.0218</b>	<b>0.0743</b>	<b>0.0484</b>	<b>111.4375</b>	<b>106.0538</b>	<b>0.9349</b>	<b>0.9406</b>

TABLE II  
PERFORMANCE COMPARISON OF ABLATION STUDY ON *Ins* DATASET. *RQK-4/6* DENOTES A LC-REC MODEL ADOPTED 4/6-TOKENS SIDS PROVIDED BY RQ-KMEANS, RESPECTIVELY. *RQK-4-6* ADOPTS 4-TOKENS SIDS FOR HEAD ITEMS WHILE 6-TOKENS SIDS FOR TAIL ITEMS, WHILE WITHOUT THE INHERITING PROCESS IN SKT.

Model	HR@5 $\uparrow$		HR@10 $\uparrow$		NDCG@5 $\uparrow$		NDCG@10 $\uparrow$		MGU@5 $\downarrow$	MGU@10 $\downarrow$	ARP@5 $\downarrow$	ARP@10 $\downarrow$
	All	Tail	All	Tail	All	Tail	All	Tail				
Ghost	0.0864	<b>0.0117</b>	0.1017	<b>0.0173</b>	0.0755	<b>0.0078</b>	0.0805	<b>0.0097</b>	<u>0.0596</u>	<u>0.0958</u>	<b>248.8179</b>	<b>209.8877</b>
w/o AUO	0.0849	<u>0.0111</u>	0.1003	<u>0.0161</u>	0.0733	<u>0.0075</u>	0.0782	<u>0.0091</u>	<b>0.0035</b>	<b>0.0112</b>	<u>270.4236</u>	<u>235.9759</u>
w/o SKT	0.0920	0.0060	0.1102	0.0095	<b>0.0798</b>	0.0047	<u>0.0856</u>	0.0058	0.3211	0.3056	387.0114	357.7677
RQK-4	<u>0.0929</u>	0.0059	<u>0.1116</u>	0.0095	0.0796	0.0042	<u>0.0856</u>	0.0054	0.3037	0.2885	304.3177	267.9522
RQK-6	<b>0.0932</b>	0.0052	<b>0.1133</b>	0.0088	<u>0.0797</u>	0.0038	<b>0.0862</b>	0.0049	0.3073	0.2902	303.1651	269.7153
RQK-4-6	0.0883	0.0070	0.1080	0.0094	0.0762	0.0049	0.0826	0.0057	0.2268	0.2468	320.2727	288.6501

drawn. *First*, bias amplification stemming from undifferentiated tokenization constitutes the critical cause of the GRs susceptibility to popularity bias. A performance comparison between Ghost and the *w/o SKT* variant reveals that the inability to inhibit the unpredictable competition between head and tail tokens leads to a substantial decline in tail recommendation performance. The results from *RQK-4*, *RQK-6*, and *RQK-4-6* serve as additional evidence supporting this finding. *Second*, the supervisory signals corrected by AUO must be absorbed by the corresponding learnable parameters. By comparing the performance of Ghost, *w/o AUO*, and *w/o SKT*, one can notice that the model ability to mitigate popularity bias is limited when relying solely on AUO. However, once SKT introduces additional tokens for tail items, the corrective effect provided by AUO drives the model closer to Pareto optimality.

### C. SID Length Analysis

As shown by Figure 4, we vary the skeleton length  $L^h \in \{3, 4, 5\}$  and the additional tail-specific length  $L^t \in \{1, 2, 3\}$ . The following insights can be drawn. A larger  $L^t$  effectively mitigates popularity bias, as evidenced by the improvement

in tail HR@10 and the reduction in ARP@10. However, this exploration comes at the expense of overall performance and fairness. This reveals a clear trade-off in identifier length allocation. Assigning longer sequences to tail items enhances their representational capacity and retrieval probability, but it simultaneously introduces noise that dilutes the learning of head items, thereby compromising overall performance. Furthermore, the skeleton length  $L^h$  dictates the extent to which tail items inherit semantic prefixes from head items. Forcing tail items to strictly mirror head items over too many generative steps suppresses their unique semantic identities.

### D. Scaling Pattern Analysis

To further investigate the scalability of our proposed framework, we compare the performance of Ghost against various baselines across different LLM backbone scales (ranging from 0.6B to 8B). The results are summarized in Table III.

In general, Ghost consistently demonstrates superior popularity bias mitigation and long-tail item excavation capabilities regardless of the underlying backbone size. In detail, we can draw the following three observations. *First*, Ghost is

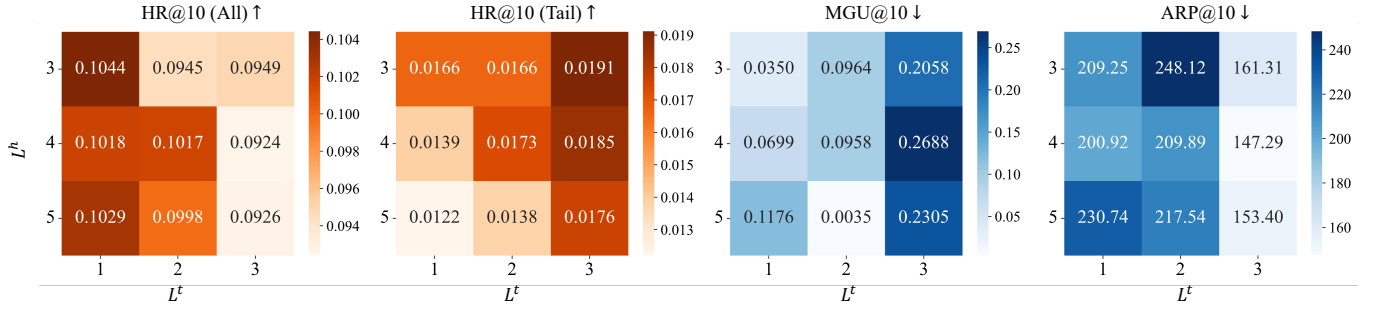
Fig. 4. Analysis of SID lengths, including head length  $L^h$  and additional length  $L^t$  for tail items.

TABLE III  
PERFORMANCE COMPARISON ACROSS DIFFERENT BACKBONE SCALES. THE BEST AND THE RUNNER-UP ARE HIGHLIGHTED IN **BOLD**, UNDERLINED FONTS.

Scale	Model	HR@5 ↑		HR@10 ↑		NDCG@5 ↑		NDCG@10 ↑		MGU@5 ↓	MGU@10 ↓	ARP@5 ↓	ARP@10 ↓
		All	Tail	All	Tail	All	Tail	All	Tail				
0.6B	LC-Rec	0.0794	0.0015	0.0983	0.0021	0.0671	0.0010	0.0731	0.0011	0.3568	0.3550	433.3062	404.2913
	IFairLRS-RW	<u>0.0847</u>	0.0025	<b>0.1030</b>	0.0036	<u>0.0725</u>	0.0018	<u>0.0784</u>	0.0021	0.3483	0.3431	393.1610	360.2472
	IFairLRS-RR	0.0819	<u>0.0053</u>	0.1010	<u>0.0060</u>	0.0705	<u>0.0045</u>	0.0766	<u>0.0047</u>	<u>0.2414</u>	<u>0.2918</u>	<u>378.3699</u>	<u>353.8215</u>
	Ghost	<b>0.0848</b>	<b>0.0101</b>	<u>0.1025</u>	<b>0.0148</b>	<b>0.0748</b>	<b>0.0069</b>	<b>0.0805</b>	<b>0.0084</b>	<b>0.0136</b>	<b>0.0208</b>	<b>239.9323</b>	<b>209.8331</b>
1.7B	LC-Rec	0.0855	0.0027	0.1025	0.0033	0.0731	0.0022	0.0786	0.0024	0.3531	0.3512	418.5129	386.3999
	IFairLRS-RW	<b>0.0879</b>	0.0038	<b>0.1070</b>	0.0057	<b>0.0764</b>	0.0026	<b>0.0825</b>	0.0032	0.3459	0.3422	373.8739	341.3491
	IFairLRS-RR	0.0847	<u>0.0084</u>	0.1047	<u>0.0090</u>	0.0734	<u>0.0070</u>	0.0798	<u>0.0071</u>	<u>0.2052</u>	<u>0.2676</u>	<u>346.5825</u>	<u>326.8923</u>
	Ghost	<u>0.0870</u>	<b>0.0110</b>	<u>0.1059</u>	<b>0.0174</b>	<u>0.0761</u>	<b>0.0081</b>	<u>0.0821</u>	<b>0.0102</b>	<b>0.0262</b>	<b>0.0508</b>	<b>248.1768</b>	<b>224.2784</b>
4B	LC-Rec	<b>0.0937</b>	0.0057	<u>0.1112</u>	0.0074	<b>0.0811</b>	0.0047	<b>0.0867</b>	0.0053	0.3363	0.3301	330.5470	301.0089
	IFairLRS-RW	<b>0.0937</b>	0.0067	<b>0.1125</b>	0.0097	<u>0.0800</u>	0.0051	<u>0.0860</u>	0.0061	0.3160	0.3033	293.3685	264.1469
	IFairLRS-RR	0.0884	<u>0.0124</u>	0.1097	<u>0.0148</u>	0.0764	<u>0.0097</u>	0.0832	<u>0.0105</u>	<u>0.0871</u>	<u>0.1836</u>	<u>256.1628</u>	<u>247.6822</u>
	Ghost	<u>0.0898</u>	<b>0.0136</b>	0.1074	<b>0.0190</b>	0.0782	<b>0.0101</b>	0.0838	<b>0.0118</b>	<b>0.0356</b>	<b>0.0211</b>	<b>239.5142</b>	<b>215.0566</b>
8B	LC-Rec	<b>0.0954</b>	0.0062	<b>0.1170</b>	0.0092	<b>0.0817</b>	0.0047	<b>0.0886</b>	0.0056	0.3233	0.3142	300.1242	<u>268.9937</u>
	IFairLRS-RW	0.0900	0.0066	0.1102	0.0088	0.0785	0.0048	0.0850	0.0055	0.3251	0.3177	317.3005	285.7003
	IFairLRS-RR	0.0872	<u>0.0120</u>	0.1076	<u>0.0140</u>	0.0764	<u>0.0095</u>	0.0830	<u>0.0101</u>	<u>0.1445</u>	<u>0.2212</u>	<u>289.5438</u>	273.9240
	Ghost	<u>0.0907</u>	<b>0.0150</b>	<u>0.1109</u>	<b>0.0228</b>	<u>0.0787</u>	<b>0.0106</b>	<u>0.0851</u>	<b>0.0131</b>	<b>0.0479</b>	<b>0.0410</b>	<b>207.7282</b>	<b>181.4011</b>

primarily characterized by its robust debiasing and tail-item recommendation performance across all scales. Compared to both the standard baseline LC-Rec and the fairness-aware baselines IFairLRS, Ghost consistently achieves the highest Tail HR and Tail NDCG, alongside the lowest ARP and MGU scores. For instance, at the 8B scale, Ghost substantially increases Tail HR@10 to 0.0228 compared to 0.0092 of LC-Rec, while drastically reducing ARP@10 to 181.4011 compared to 268.9937 of LC-Rec. *Second*, with respect to overall recommendation performance, Ghost maintains highly competitive utility. Despite aggressively promoting tail items, the overall HR and NDCG metrics of Ghost remain closely comparable to those of the standard generative recommendation models and generally surpass the IFairLRS baselines. This indicates that Ghost achieves an excellent trade-off, mitigating popularity bias without inducing unacceptable degradation in general recommendation accuracy. *Third*, Ghost exhibits strong scalability with respect to LLM backbone capacity. As the parameter size increases from 0.6B to 8B, Ghost effectively leverages the enhanced representation and reasoning capabilities of larger models to further boost its performance. Specifically, Ghost’s Tail HR@10 steadily increases from 0.0148 at 0.6B to 0.0228 at 8B, while its ARP@10 progressively drops from 209.8331 to 181.4011. This demonstrates that larger

backbones consistently empower the framework to deliver increasingly diverse, balanced, and fair recommendations.

## VI. HEAD-TO-TAIL RATIO ANALYSIS

In this section, we present the ratio of head items to tail items in the recommendation results provided by each GR model. Figure 5 illustrates the exposure distribution of head and tail items retrieved by Ghost and various baselines. While existing generative baselines like LETTER, LC-Rec, and ED2 exhibit severe popularity bias by disproportionately favoring head items. For example, LC-Rec retrieves 24,212 head items versus only 560 tail items. Ghost significantly alleviates this discrepancy, achieving a balanced distribution of 15,372 head and 9,400 tail items. Furthermore, evaluating Ghost across different Qwen3 backbone scales reveals a positive correlation between LLM capacity and long-tail item excavation. Upgrading the backbone from 0.6B to 4B parameters steadily increases the retrieval of tail items from 7,344 to 9,688, demonstrating that the enhanced representation capabilities of larger models effectively empower the framework to deliver diverse, balanced recommendations without overfitting to mainstream trends.

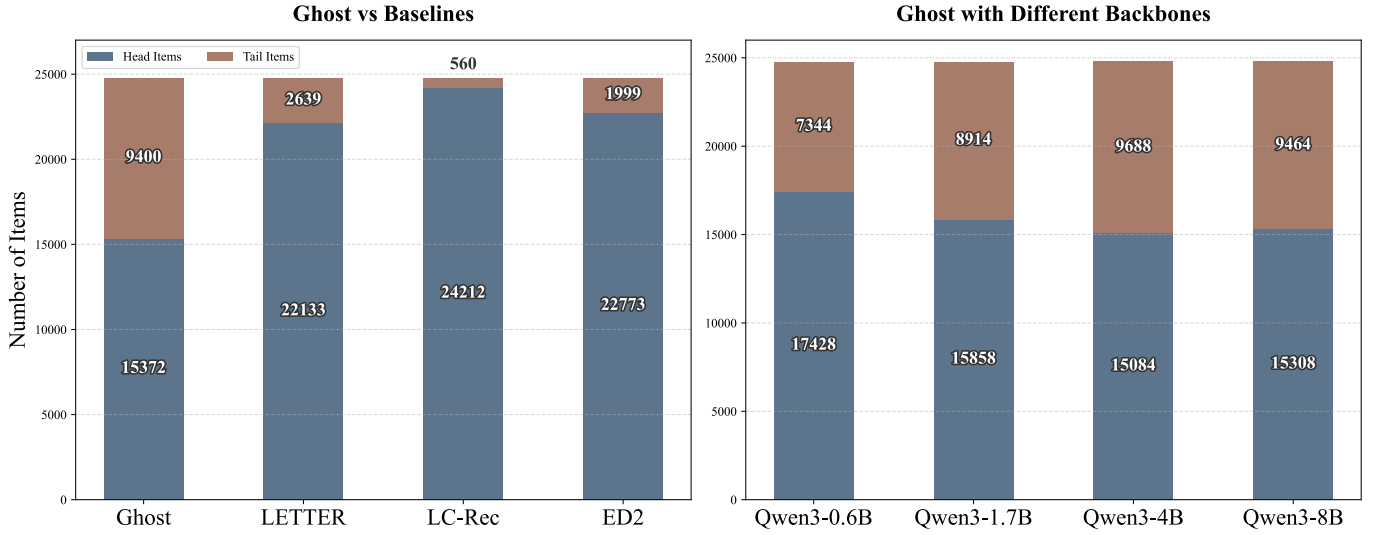


Fig. 5. Numbers of head and tail items in the recommendation results provided by (left) Ghost and baseline models, and (right) Ghost with different backbones.

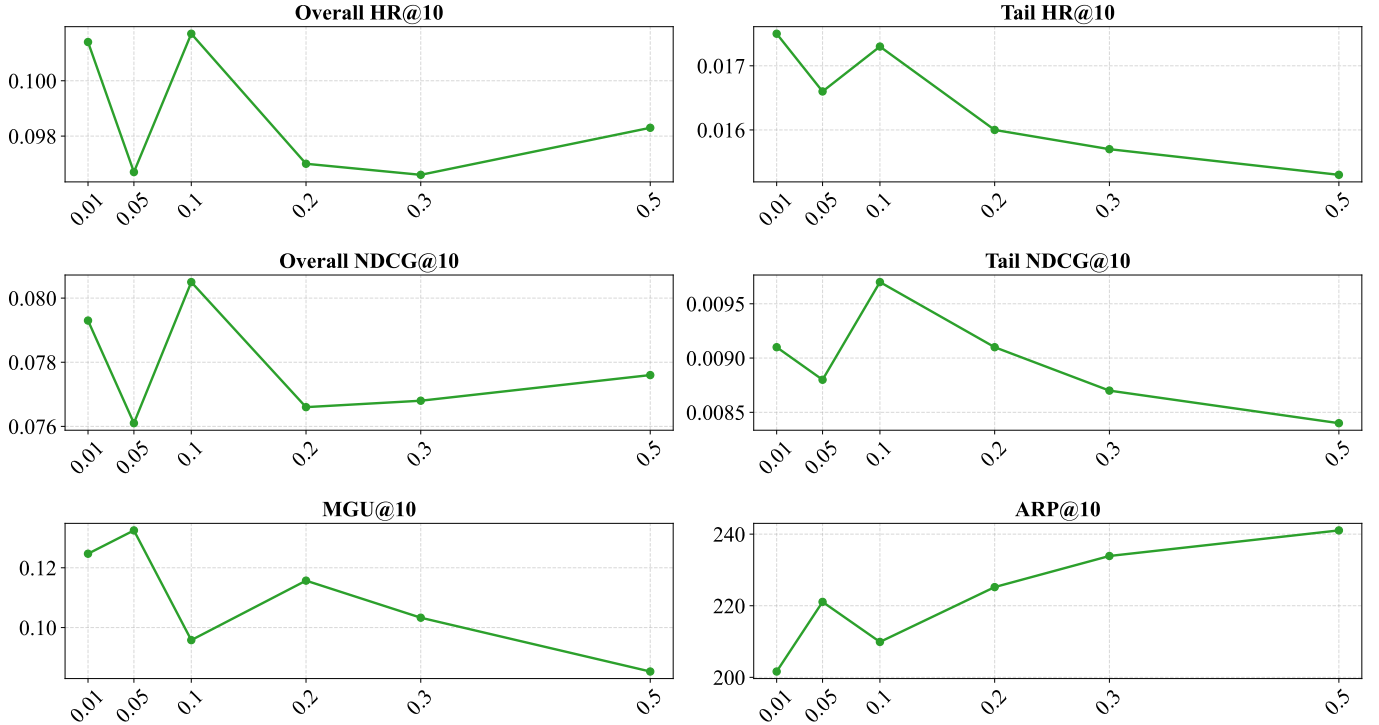


Fig. 6. Tendency of Ghost performance on Ins dataset, under different AUO weights  $\alpha$ . The  $x$ -axis denotes the values of  $\alpha$ , and the  $y$ -axis is the metric values.

## VII. HYPER-PARAMETER ANALYSIS

### A. AUO Weight $\alpha$

Here, we investigate the impact of the weight parameter  $\alpha$  of AUO. Figure 6 investigates the sensitivity of the Ghost model to the AUO weight  $\alpha$  on the Ins dataset. The results indicate that  $\alpha$  plays a critical role in balancing recommendation accuracy and popularity bias mitigation. Specifically,  $\alpha = 0.1$  emerges as the optimal setting, achieving the highest Overall NDCG@10 and Tail NDCG@10 while maintaining highly competitive Overall and Tail HR@10. Crucially, at this optimal

point, the ARP@10 experiences a local drop, confirming the model’s capability to effectively surface less popular items without sacrificing accuracy. Conversely, increasing the weight beyond 0.1 (e.g., 0.2 to 0.5) consistently degrades both overall and tail-specific accuracy metrics, significantly diminishes MGU@10, and steadily increases ARP@10. This clear trend demonstrates that an excessively large AUO weight forces the model to over-prioritize mainstream items, thereby exacerbating popularity bias and compromising both long-tail excavation and overall recommendation quality.

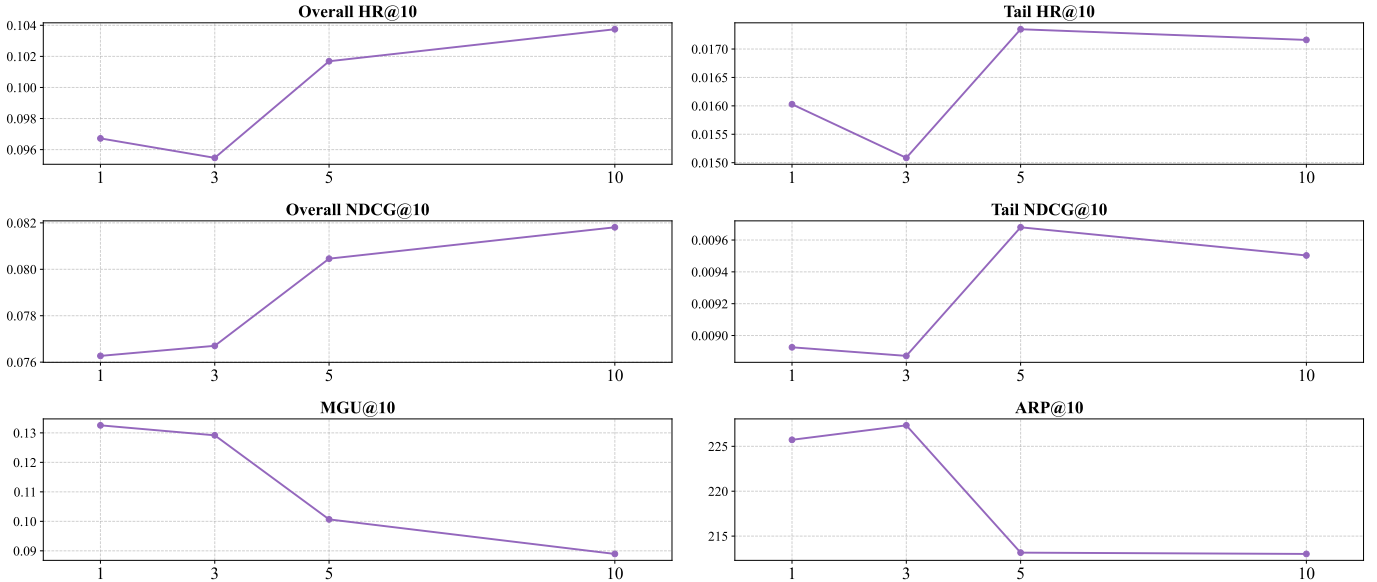


Fig. 7. Tendency of Ghost performance on Ins dataset, under different undesired collection sizes. The  $x$ -axis denotes the size of undesired collection, and the  $y$ -axis is the metric values.

### B. Size of Undesired Collection $|\bar{\Omega}|$

Here, we investigate the impact of the undesired collection size. Figure 7 investigates the impact of the undesired collection size  $|\bar{\Omega}|$  on the performance of the Ghost model using the Ins dataset. The empirical results demonstrate that the magnitude of  $|\bar{\Omega}|$  critically influences the trade-off between recommendation accuracy and popularity bias mitigation. Specifically, setting the collection size to 5 emerges as the optimal configuration; at this point, the model achieves peak effectiveness in excavating long-tail items, evidenced by the highest Tail HR@10 and Tail NDCG@10, alongside a sharp and favorable decline in ARP@10. Concurrently, overall accuracy metrics (Overall HR@10 and NDCG@10) experience substantial gains compared to smaller sizes. However, expanding the undesired collection further to 10 yields diminishing returns: while overall accuracy marginally increases, tail-specific metrics slightly degrade, and user coverage (MGU@10) continues a steady, negative decline. This indicates that a moderate undesired collection size effectively provides sufficient contrast against over-recommended popular items to surface relevant tail items, whereas an excessively large collection may introduce noise that compromises tail item retrieval and overall recommendation diversity across different users.

## VIII. CONCLUSION

In this paper, we theoretically demonstrate that standard MLE training and undifferentiated item tokenization of current GRs inherently cause a token-level optimization flaw and multi-step bias amplification. Accordingly, we propose Ghost equipped with asymmetric unlikelihood optimization (AUO) and skeleton-founded tokenization (SKT). AUO provides explicit negative supervision to rescue tail tokens from gradient starvation, while SKT establishes unified branching points to

halt item-level bias amplification. Extensive empirical evaluations confirm that Ghost effectively breaks the filter bubble and substantially promotes fairer long-tail recommendations with slight losses to overall recommendation utility, approaching Pareto optimality in GRs popularity debiasing.

### IMPACT STATEMENT

This work aims to enhance the fairness and diversity of Generative Recommender Systems. By mitigating the fundamental causes of popularity bias, our proposed model effectively breaks the pervasive filter bubble. Ethically, this approach addresses the algorithmic "Matthew Effect," which typically over-represents trending items while severely marginalizing niche content. By promoting fairer long-tail recommendations without significantly compromising overall utility, this research fosters the responsible deployment of LLM-based recommendation technologies.

## REFERENCES

- [1] G. Zhou, C. Song, X. Zhu, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, "Deep interest network for click-through rate prediction," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018.
- [2] P. Covington, J. Adams, and E. Sargin, "Deep neural networks for youtube recommendations," in *Proceedings of the ACM Conference on Recommender Systems*, 2016.
- [3] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang, and J. Tang, "Deepinf: Social influence prediction with deep learning," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018.
- [4] S. Rajput, N. Mehta, A. Singh, R. H. Keshavan, T. Vu, L. Heldt, L. Hong, Y. Tay, V. Q. Tran, J. Samost, M. Kula, E. H. Chi, and M. Sathiamoorthy, "Recommender systems with generative retrieval," in *Proceedings of the International Conference on Neural Information Processing Systems*, 2023.
- [5] B. Zheng, Y. Hou, H. Lu, Y. Chen, W. X. Zhao, and M. Chen, "Adapting large language models by integrating collaborative semantics for recommendation," in *Proceedings of the IEEE International Conference on Data Engineering*, 2024.
- [6] W. Wang, H. Bao, X. Lin, J. Zhang, Y. Li, F. Feng, S.-K. Ng, and T.-S. Chua, "Learnable item tokenization for generative recommendation," in *Proceedings of the ACM International Conference on Information and Knowledge Management*, 2024.
- [7] J. Yin, Z. Zeng, M. Li, H. Yan, C. Li, W. Han, J. Zhang, R. Liu, H. Sun, W. Deng, F. Sun, Q. Zhang, S. Pan, and S. Wang, "Unleash llms potential for sequential recommendation by coordinating dual dynamic index mechanism," in *Proceedings of the ACM on Web Conference*, 2025.
- [8] J. Zhai, Z.-F. Mai, C.-D. Wang, F. Yang, X. Zheng, H. Li, and Y. Tian, "Multimodal quantitative language for generative recommendation," in *Proceedings of the International Conference on Learning Representations*, 2025.
- [9] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, "Lightgcn: Simplifying and powering graph convolution network for recommendation," in *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020.
- [10] W.-C. Kang and J. McAuley, "Self-attentive sequential recommendation," in *Proceedings of the IEEE International Conference on Data Mining*, 2018.
- [11] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," in *Proceedings of the International Conference on Neural Information Processing Systems*, 2017.
- [12] D. Lee, C. Kim, S. Kim, M. Cho, and W.-S. Han, "Autoregressive image generation using residual quantization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [13] D. Guo, D. Yang, H. Zhang, J. Song, P. Wang *et al.*, "Deepseek-r1 incentivizes reasoning in llms through reinforcement learning," *Nature*, vol. 645, pp. 633–638, 2025.
- [14] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, and F. Azhar, "Llama: Open and efficient foundation language models," 2023.
- [15] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.
- [16] S. Lin, C. Gao, J. Chen, S. Zhou, B. Hu, Y. Feng, C. Chen, and C. Wang, "How do recommendation models amplify popularity bias? an analysis from the spectral perspective," in *Proceedings of the ACM International Conference on Web Search and Data Mining*, 2025.
- [17] T. Wei, F. Feng, J. Chen, Z. Wu, J. Yi, and X. He, "Model-agnostic counterfactual reasoning for eliminating popularity bias in recommender system," in *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021.
- [18] Z. Zhu, Y. He, X. Zhao, Y. Zhang, J. Wang, and J. Caverlee, "Popularity-opportunity bias in collaborative filtering," in *Proceedings of the ACM International Conference on Web Search and Data Mining*, 2021.
- [19] M. Jiang, K. Bao, J. Zhang, W. Wang, Z. Yang, F. Feng, and X. He, "Item-side fairness of large language model-based recommendation system," in *Proceedings of the ACM Web Conference*, 2024.
- [20] L. Wang, C. Ma, X. Wu, Z. Qiu, Y. Zheng, and X. Chen, "Causally debiased time-aware recommendation," in *Proceedings of the ACM Web Conference*, 2024.
- [21] F. Zhang and Q. Shen, "A model-agnostic popularity debias training framework for click-through rate prediction in recommender system," in *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2023.
- [22] X. Li, Y. Chen, B. Pettit, and M. D. Rijke, "Personalised reranking of paper recommendations using paper content and user behavior," *ACM Transactions on Information Systems*, vol. 37, no. 3, pp. 1–23, 2019.
- [23] D. Carraro and D. Bridge, "Enhancing recommendation diversity by re-ranking with large language models," *ACM Transactions on Recommender Systems*, vol. 4, no. 2, pp. 1–40, 2025.
- [24] K. Miettinen, *Nonlinear multiobjective optimization*. Springer Science & Business Media, 1999, vol. 12.
- [25] D. Mahapatra and V. Rajan, "Multi-task learning with user preferences: Gradient descent with controlled ascent in pareto optimization," in *Proceedings of the International Conference on Machine Learning*, 2020.
- [26] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua, "Neural collaborative filtering," in *Proceedings of the International Conference on World Wide Web*, 2017.
- [27] J. Li, H. Gu, S. Wang, Q. Zhang, S. Yu, C. Wang, X. Xu, and F. Chen, "Towards fair large language model-based recommender systems without costly retraining," in *Proceedings of the ACM Web Conference*, 2026.
- [28] X. Lin, P. Liu, W. Wang, Y. Hu, C. Xu, F. Feng, Q. Wang, and T.-S. Chua, "Bringing reasoning to generative recommendation through the lens of cascaded ranking," in *Proceedings of the ACM Web Conference*, 2026.
- [29] X. Luo, J. Cao, T. Sun, J. Yu, R. Huang, W. Yuan, H. Lin, Y. Zheng, S. Wang, Q. Hu, C. Qiu, J. Zhang, X. Zhang, Z. Yan, J. Zhang, S. Zhang, M. Wen, Z. Liu, and G. Zhou, "Qarm: Quantitative alignment multi-modal recommendation at kuaishou," in *Proceedings of the ACM International Conference on Information and Knowledge Management*, 2025.
- [30] W. Ren, L. Wang, K. Liu, R. Guo, L. E. Peng, and Y. Fu, "Mitigating popularity bias in recommendation with unbalanced interactions: A gradient perspective," in *Proceedings of the IEEE International Conference on Data Mining*, 2022.
- [31] M. Pezeshki, O. Kaba, Y. Bengio, A. C. Courville, D. Precup, and G. Lajoie, "Gradient starvation: A learning proclivity in neural networks," in *Proceedings of the International Conference on Neural Information Processing Systems*, 2021.
- [32] S. Welleck, I. Kulikov, S. Roller, E. Dinan, K. Cho, and J. Weston, "Neural text generation with unlikelihood training," in *Proceedings of the International Conference on Learning Representations*, 2020.
- [33] E. Lagutin, D. Gavrilo, and P. Kalaidin, "Implicit unlikelihood training: Improving neural text generation with reinforcement learning," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021.
- [34] C. M. Ju, L. Collins, L. Neves, B. Kumar, L. Y. Wang, T. Zhao, and N. Shah, "Generative recommendation with semantic ids: A practitioner's handbook," in *Proceedings of the ACM International Conference on Information and Knowledge Management*, 2025.
- [35] J. Chen, H. Dong, X. Wang, F. Feng, M. Wang, and X. He, "Bias and debias in recommender system: A survey and future directions," *ACM Transactions on Information Systems*, vol. 41, no. 3, pp. 1–39, 2023.
- [36] T. Schnabel, A. Swaminathan, A. Singh, N. Chandak, and T. Joachims, "Recommendations as treatments: Debiasing learning and evaluation," in *Proceedings of the International Conference on Machine Learning*, 2016.
- [37] C. Gao, Y. Zheng, W. Wang, F. Feng, X. He, and Y. Li, "Causal inference in recommender systems: A survey and future directions," *ACM Transactions on Information Systems*, vol. 42, no. 4, pp. 1–32, 2024.
- [38] Y. Zhang, F. Feng, X. He, T. Wei, C. Song, G. Ling, and Y. Zhang, "Causal intervention for leveraging popularity bias in recommendation," in *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021.
- [39] Y. Linde, A. Buzo, and R. Gray, "An algorithm for vector quantizer design," *IEEE Transactions on communications*, vol. 28, no. 1, pp. 84–95, 1980.
- [40] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 873–12 883.
- [41] R. He and J. McAuley, "Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering," in *Proceedings of the International Conference on World Wide Web*, 2016.
- [42] M. Cuturi, "Sinkhorn distances: lightspeed computation of optimal transport," in *Proceedings of the International Conference on Neural Information Processing Systems*, vol. 2, 2013, pp. 2292–2300.

- [43] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang *et al.*, “Qwen2 technical report,” *arXiv preprint arXiv:2407.10671*, 2024.
- [44] Q. Team, “Qwen2.5: A party of foundation models,” September 2024. [Online]. Available: <https://qwenlm.github.io/blog/qwen2.5/>
- [45] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv *et al.*, “Qwen3 technical report,” *arXiv preprint arXiv:2505.09388*, 2025.
- [46] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proceedings of the International Conference on Learning Representations*, 2015.
- [47] G. Zhou, H. Bao, J. Huang, J. Deng, J. Zhang, J. She, K. Cai, L. Ren, L. Ren, Q. Luo *et al.*, “Openorec technical report,” *arXiv preprint arXiv:2512.24762*, 2025.
- [48] C. Gao, R. Chen, S. Yuan, K. Huang, Y. Yu, and X. He, “Sprec: Self-play to debias llm-based recommendation,” in *Proceedings of the ACM on Web Conference*, 2025.

APPENDIX A  
NOTATION

The notations and corresponding descriptions are summarized in Table IV.

TABLE IV  
SUMMARY OF MATHEMATICAL AND MODEL NOTATIONS

Symbol	Description	Symbol	Description	Symbol	Description
$K, J$	Number of items and users in the system, respectively.	$L$	Length of SIDs.	$\mathcal{V}_{\text{head}}$	Set of popular head items.
$v_k, u_j$	The $k$ -th item and $j$ -th user.	$\theta$	Parameters of the Generative Recommender (GR) model.	$L^h$	SID length assigned for head items (the skeleton length).
$h_{u_j}$ or $h_u$	Historical behavior / item sequence of user $u_j$ .	$\mathcal{P}_\theta$	Token generation probability parameterized by $\theta$ .	$v'$	A target tail item.
$l$	Length of the item sequence.	$\mathcal{P}_d$	True data distribution.	$v^*$	The closest head item to tail item $v'$ based on highest semantic similarity.
$T_k$	Textual features (e.g., title, description) attached to item $v_k$ .	$\mathcal{L}_{\text{NLL}}$	Negative log-likelihood loss used for Maximum Likelihood Estimation (MLE).	$L^t$	Additional SID tokens length specifically for tail items.
$f_e$	Pre-trained textual encoder.	$\mathcal{D}$	Training distribution of user-item interaction pairs.	$\gamma_{EOS}$	Head-dominance factor at the $(L^h + 1)$ -th generative step against the EOS token.
$X_v, X_k$	Semantic textual representation of item $v$ or $k$ .	$X_{h_u}$	Encoded representation of the user's historical behavior $h_u$ .	$\bar{\Omega}$	SID collection of roughly selected undesired head items.
$\mu_n^{(i)}$	The $n$ -th embedding in the $i$ -th codebook within the RQ-VAE.	$e_c$	Embedding of token $c$ .	$\mathcal{V}_{\text{rough}}$	Rough candidate set of popular head items semantically similar to the tail item.
$r_v^{(i)}$	Residual embedding at step $i$ during tokenization.	$c_{\text{head}}^{(i)}, c_{\text{tail}}^{(i)}$	Candidate head and tail tokens competing at the $i$ -th step.	$K_a, K_b$	Hyper-parameters controlling the candidate scale for undesired item selection.
$c_v^{(i)}$	The $i$ -th SID token for item $v$ .	$\gamma_i$	Amplification factor at step $i$ .	$\mathcal{L}_{\text{AUO}}$	Asymmetric unlikelihood optimization (AUO) loss.
$\Omega_v$	Complete SID of item $v$ , indexed as $(c_v^{(1)}, c_v^{(2)}, \dots, c_v^{(L)})$ .	$\mathcal{Z}$	Set of steps during tail item generation where it competes against head tokens.	$\mathcal{L}_{\text{All}}$	Overall optimization objective function of the Ghost model.
$c_v^{<i}$	Sub-sequence of SID $\Omega_v$ before the $i$ -th position.	$\gamma_{\text{min}}$	Minimum factor across the competing steps in $\mathcal{Z}$ .	$\alpha$	Weighted parameter controlling the AUO loss.

APPENDIX B  
RELATED WORK

*A. Recommender Systems: Discriminative and Generative Paradigms*

Traditional recommender systems typically formulate recommendation as a discriminative task, relying heavily on discrete item IDs to capture user preferences by analyzing historical interactions [1], [2], [9], [10], [26]. At the core of this pipeline, both users and items are projected into a shared low-dimensional latent space to learn dense embeddings, allowing the system to determine the final ranking by computing their pairwise similarities or interaction scores. To efficiently handle massive item catalogs, these systems typically employ a multi-stage pipeline, primarily consisting of candidate generation and ranking. Specifically, the candidate generation stage rapidly retrieves a coarse-grained subset of relevant items from the entire corpus, which are subsequently evaluated and sorted by a more complex ranking model to produce the final recommendations.

Recently, fueled by the adoption of large language models (LLMs) [13]–[15] as the underlying backbone, Generative Recommenders (GRs) have emerged as a transformative paradigm [4], [34]. Instead of modeling interaction probabilities directly, GRs replace traditional item IDs with Semantic IDs (SIDs), reframing recommendation as a unified end-to-end sequential generation process [5]–[8]. To construct these SIDs, standard vector quantization techniques, such as VQ-VAE [11], RQ-VAE [12], and RQ-KMeans [29], are widely utilized to convert continuous item embeddings into discrete indices [4], [5]. However, these established tokenization strategies are fundamentally undifferentiated. They assign SIDs identically without accounting for inherent item popularity disparities, leading to unstructured and unpredictable branching points where head and tail item tokens compete.

*B. Popularity Bias and Fairness in Recommendation*

Traditional recommender systems frequently suffer from popularity bias, a phenomenon where a small fraction of highly interacted items disproportionately dominates algorithm exposure, leaving the long-tail of items largely ignored [16]. This algorithmic amplification, often referred to as the *rich-get-richer* effect (a.k.a. the Matthew Effect), not only degrades user experience by failing to capture niche or diverse interests but also introduces critical fairness concerns across the platform,

TABLE V  
STATISTICS OF THE EVALUATED DATASETS. *Avg.L* IS THE AVERAGE LENGTH OF THE USER INTERACTION SEQUENCES.

Dataset	#User	#Item	#Interaction	Sparsity	Avg.L
Instruments	24,772	9,922	206,153	99.92%	8.32
Games	50,546	16,859	452,989	99.95%	8.96
Arts	45,141	20,956	390,832	99.96%	8.66

particularly for minority content creators who receive inequitable visibility. To mitigate these disparities, early literature primarily focused on heuristic post-processing and re-ranking techniques that explicitly boosted the exposure of long-tail items, albeit often at the expense of overall accuracy [35]. More recently, the field has gravitated toward principled statistical and causal frameworks to achieve debiased learning. Schnabel et al. utilize inverse propensity weighting (IPW) [36] to directly correct data collection biases by re-weighting user-item interactions during the model training phase. Furthermore, recent advances leveraging causal inference [37] and adversarial learning have enabled systems to explicitly disentangle genuine user preferences from popularity-driven conformity, yielding robust representations that strive to optimize accuracy while maintaining multi-sided fairness for both users and providers [38].

### C. Vector Quantization

Vector Quantization (VQ) originated as a classical signal processing technique for data compression [39], but it has recently become a cornerstone of deep representation learning. The seminal VQ-VAE [11] successfully integrated discrete latent spaces into neural architectures, effectively resolving the posterior collapse issue common in continuous models. Building upon this, frameworks like VQGAN [40] utilized adversarial objectives to dramatically enhance reconstruction fidelity. Crucially, by mapping high-dimensional data into discrete codebook indices, VQ allows complex signals to be treated as sequential tokens for autoregressive transformers. Beyond visual and audio synthesis, this discrete tokenization paradigm has recently advanced generative recommender systems. As we mentioned above, by employing VQ techniques, such as residual quantization [12], [29], to discretize continuous item embeddings into semantic, categorical IDs, researchers have successfully reformulated recommendation as a sequence-to-sequence generation task [4], [5], [7]. This approach allows large language models to autoregressively predict next-item interactions using these discrete item tokens, seamlessly bridging traditional collaborative filtering with the powerful generative and reasoning capabilities of modern Transformer architectures.

## APPENDIX C EXPERIMENTAL DETAIL

### A. Dataset

The dataset statistics are presented in Table V. The three sequential recommendation datasets originate from the Amazon Product Review dataset [41], which contains user review data from May 1996 to October 2018. Particularly, three categories for the sequential recommendation task, including "*Musical Instruments*", "*Video Games*", and "*Arts, Crafts and Sewing*", are extracted and organized into individual datasets *Instruments*, *Games*, and *Arts*, respectively. Within the above datasets, each item is associated with a series of textual contents, including the item title, the detailed description, the item category, and so on. Similarly, the associated textual contents of the user entity include the user comment, the search query, and so on. Following standard procedure, inactive users/items with fewer than 5 interactions are filtered out, and the user interaction sequence is created in chronological order. In Figure 8, we present the distribution of item popularity for each of the three datasets, all of which are heavily long-tailed.

### B. Baseline

- **LETTER** [6], short for LEarnable Tokenizer for generaTivE Recommendation, is a learnable item tokenizer tailored for LLM-based generative recommendation. It addresses the limitations of existing ID, textual, and codebook-based identifiers by comprehensively integrating hierarchical semantics, collaborative signals, and code assignment diversity into item identifiers. Specifically, LETTER employs a RQ-VAE to encode item semantic information into hierarchical code sequences. To overcome the misalignment between semantics and collaborative signals, it introduces a contrastive alignment loss to align semantic quantized embeddings with collaborative filtering embeddings from well-trained models. Additionally, it applies a diversity loss based on constrained  $K$ -means clustering to regularize code embeddings, effectively mitigating code assignment and item generation biases. When instantiated on generative recommender models, LETTER further incorporates a ranking-guided generation loss to theoretically and empirically augment their top- $K$  ranking capabilities.

- **LC-Rec** [5] is a generative large language model based sequential recommendation approach designed to bridge the semantic gap between the language semantics of LLMs and the collaborative semantics of recommender systems. To represent

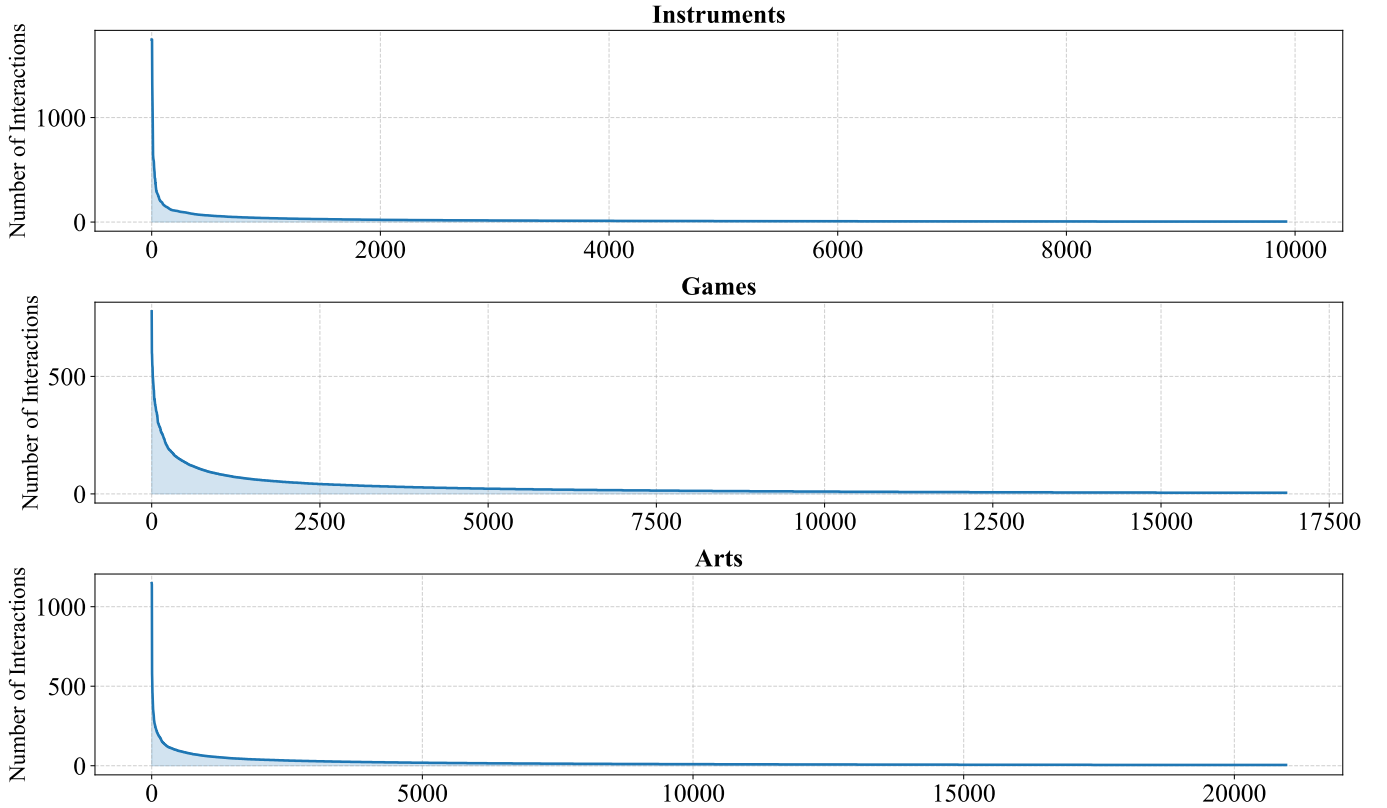


Fig. 8. Long tail distribution of the item popularity.

items effectively without vocabulary explosion, LC-Rec employs a tree-structured vector quantization method based on item text embeddings to construct discrete item indices. It further utilizes a uniform semantic mapping technique [42] during index allocation to eliminate potential index conflicts among items. Instead of relying on predefined candidate sets, LC-Rec can autoregressively generate target items from the entire item set. To achieve deep integration of language and collaborative semantics, the LLM is fine-tuned on a series of specialized semantic alignment tasks, including sequential item prediction, explicit index-language alignment, and implicit recommendation-oriented alignment.

- **ED<sup>2</sup>**, i.e., the End-to-End Dual Dynamic recommender [7], is an LLM-based sequential recommender system that introduces a dual dynamic index mechanism. It addresses the limitations of existing LLM-based models that typically separate index generation from the sequential recommendation process and neglect user-related information. By utilizing a dual architecture with two homogeneous discrete index generators, ED<sup>2</sup> synchronously generates indices for both users and items, assembling index generation and sequential recommendation into a unified end-to-end LLM pipeline. To facilitate the LLM comprehension of the untrained dynamic index tokens, the model incorporates a multi-grained token regulator that establishes alignment supervision between dynamic index tokens and corresponding natural language tokens. Additionally, ED<sup>2</sup> leverages customized instruction tuning tasks and associated user collection data to exploit implicit high-order user-item interaction patterns based on historical behaviors.

- **Augmentation & Substitution** replace popular head items in the user interaction sequence with their most similar tail items according to a pre-defined replacement probability  $p$ . This straightforward principle aims to mitigate popularity bias by artificially increasing the representation of tail items during training. Specifically, for a typical long-tailed distribution where head items account for 80% of total interactions and tail items account for 20%, we set the target replacement probability  $p = 0.375$ . This value is derived by solving the equilibrium equation  $0.8 - 0.8p = 0.2 + 0.8p$  to achieve a perfectly balanced ratio in the modified sequence. The target tail item for substitution is typically selected based on pre-calculated item similarity. Finally, the modified interaction sequence is integrated into the training process in two different ways. (i) The *Substitution* baseline directly overwrites the original sequence, achieving a strictly balanced head-to-tail ratio. (ii) Alternatively, the *Augmentation* baseline appends the modified sequence to the original training set. By combining the original data (0.8 head and 0.2 tail) with the modified data (0.5 head and 0.5 tail), *Augmentation* yields a final head-to-tail ratio of 13:7, approximately 1.857, which significantly alleviates the initial 4:1 imbalance while strictly preserving the original interaction contexts.

- **IFairLRS** [19] is an effective framework designed to enhance the item-side fairness of large language model-based recommendation systems (LRS). The framework addresses the item-side unfairness in LRS that primarily stems from two factors: the imbalanced distribution of historical user interactions and the inherent semantic biases present within LLMs. To

mitigate these issues without sacrificing recommendation accuracy, IFairLRS calibrates recommendations by deploying two specifically adapted strategies across the main stages of building an LRS. During the instruction finetuning stage, it employs a reweighting strategy that adjusts the weights of training samples based on the bias observed between the distribution of target items and historical interactions. In the post learning stage (i.e., inference), it utilizes a re-ranking strategy that incorporates a punishment term based on group unfairness (GU) to adjust the final top- $K$  recommendations.

### C. Metric

Following well-established benchmarks [5], [34], this study adopts Hit-Rate (HR) and normalized discounted cumulative gain (NDCG) to evaluate the recommendation performance. Formally, HR@ $K$  is measured as follows,

$$\text{HR}@K = \frac{1}{|U_{\text{test}}|} \sum_{u \in U_{\text{test}}} \mathbb{I}(\mathcal{V}_u \cap L_u^{(K)} \neq \emptyset), \quad (16)$$

where  $U_{\text{test}}$  denotes the user set for evaluation,  $L_u^{(K)}$  denotes the top- $K$  recommendation list for user  $u$ ,  $\mathcal{V}_u$  represents the set of ground-truth interacted items for user  $u$  in the test set, and  $\mathbb{I}(\cdot)$  is an indicator function that returns 1 if the condition is true and 0 otherwise.

Furthermore, NDCG@ $K$  is defined below,

$$\text{NDCG}@K = \frac{1}{|U_{\text{test}}|} \sum_{u \in U_{\text{test}}} \frac{\text{DCG}_u@K}{\text{IDCG}_u@K}, \quad (17)$$

$$\text{DCG}_u@K = \sum_{i=1}^K \frac{r_{u,i}}{\log_2(i+1)}. \quad (18)$$

Here,  $r_{u,i} \in \{0, 1\}$  indicates whether the  $i$ -th recommended item in  $L_u^{(K)}$  is relevant, i.e., exists in  $\mathcal{V}_u$ .  $\text{IDCG}_u@K$  represents the ideal DCG score obtained by perfectly ranking all relevant items in  $\mathcal{V}_u$  at the very top of the recommendation list.

This study adopts mean group unfairness (MGU) and average recommendation popularity (ARP) to quantify the fairness and the average popularity of the recommendation results [27]. Particularly, MGU is measured as follows,

$$\text{MGU} = (\text{GU}_{\text{head}} + \text{GU}_{\text{tail}})/2, \quad (19)$$

where  $\text{GU}_G$  stands for the unfairness of group  $G$  and  $G \in \{\text{head}, \text{tail}\}$ .

$\text{GU}_G$  is defined as  $\text{GU}_G = \text{GR}_G - \text{GH}_G$ , where  $\text{GR}_G$  and  $\text{GH}_G$  represent the popularity of group  $G$  in recommendation results and interaction history, respectively. Formally, let  $\mathbf{H}$  denote the set of all user interaction sequences in the history,  $\mathbf{L}$  denote the set of top- $K$  recommendations of all users at the inference phase, and  $\mathcal{G}$  denote the set of item groups. Given an item group  $G \in \mathcal{G}$ , we can measure the recommendation proportion of group  $G$  by

$$\text{GR}_G = \frac{\sum_{L \in \mathbf{L}} \sum_{v \in L} \mathbb{I}(v \in G)}{\sum_{G' \in \mathcal{G}} \sum_{L \in \mathbf{L}} \sum_{v \in L} \mathbb{I}(v \in G')}, \quad (20)$$

where  $\mathbb{I}(v \in G)$  is an identity function:

$$\mathbb{I}(v \in G) = \begin{cases} 1, & \text{item } v \text{ belongs to group } G \\ 0, & \text{otherwise} \end{cases}. \quad (21)$$

Intuitively,  $\text{GR}_G$  calculates the recommendation proportion of group  $G$  in the top- $K$  recommendations of all users. Accordingly, the interaction proportion of group  $G$  in the historical interaction sequences  $\mathbf{H}$  can be obtained by

$$\text{GH}_G = \frac{\sum_{H \in \mathbf{H}} \sum_{v \in H} \mathbb{I}(v \in G)}{\sum_{G' \in \mathcal{G}} \sum_{H \in \mathbf{H}} \sum_{v \in H} \mathbb{I}(v \in G')}. \quad (22)$$

Furthermore, ARP is defined below,

$$\text{ARP} = \frac{1}{|U_{\text{test}}|} \sum_{u \in U_{\text{test}}} \frac{\sum_{i \in L_u} \varphi(i)}{|L_u|}, \quad (23)$$

where  $U_{\text{test}}$  denotes the user set for evaluation,  $L_u$  denotes the recommendation list for user  $u$ , and  $\varphi(i)$  returns the popularity (i.e., the number of occurrences in the training set) of item  $i$ .

### D. Implementation Detail

This study employs the open source large language model *Qwen* [43]–[45] developed by Alibaba. In particular, the experiments in Section V, including the main result, the ablation study, and the SID length analysis, adopt *Qwen2.5-3B* as the backbone model for all the evaluated GR models (i.e, Ghost, standard GRs, and GRs popularity debiasing methods). For the scaling pattern analysis, we adopt the fancy *Qwen3 series*, and the backbone scale ranges from 0.6B to 8B. The hidden state dimension of *Qwen2.5-3B* is 2,048. The hidden state dimensions of the *Qwen3 series* are listed in Table VI. The original size of all the *Qwen* vocabulary is 151,936.

Within the *skeleton-founded item tokenization (SKT)*, a series of linear layers [4096, 2048, 1024, 512, 256, 128, 64] is adopted to gradually reduce the *Qwen* representation into a 32-dimensional space for RQ-Kmeans. Within each RQ-Kmeans iteration, the number of clustering centers is 256 for both head items and tail items. For head items and tail items, the lengths of SIDs generated by SKT are distinct. The length of the head item SIDs exactly equals  $L^h$ , while that of the tail item SIDs is  $L^h + L^t$ . Hence, the SIDs generated by SKT are essentially indices with variable lengths. Figure 9 presents a box plot illustrating the distribution of the number of head prefixes inherited by tail items across three datasets. The visualization highlights the central tendency, variance, and outliers within each domain. The data indicate that the *Ins* and *Arts* datasets share identical descriptive statistics, both demonstrating a median value of 4 and an interquartile range (IQR) spanning from 2 to 7. However, the *Games* dataset exhibits slightly lower overall values, featuring a median of 3 and an IQR ranging from 1 to 6. For the *asymmetric unlikelihood optimization (AUO)*, the default values of  $K_a, K_b$  controlling the candidate scale of the undesired item selection are set to 200 and 5. Therefore, for each tail item  $v'$ , the cardinality of the undesired collection  $\Omega_{v'}$  is 5. Regarding the Ghost training phase, we adopt the AdamW optimizer [46] with a learning rate of  $3 \times 10^{-5}$ . In Section VII, we conduct an analysis of the undesired collection size and the learning rate. All the experiments are completed on a machine with 8 *NVIDIA A100 Tensor Core 80GB* GPUs.

#### APPENDIX D EQUAL-SIZED GROUPING BASED ON ITEM POPULARITY

To provide a more fine-grained understanding of where the performance improvements originate, we analyze the recommendation results across different popularity segments. Figure 10 details the performance comparison on the *Ins* dataset, breaking down the item space into five equal-sized groups sorted by popularity (from the most popular 0-20 segment to the least popular 80-100 segment), alongside overall accuracy and debiasing metrics. In general, the fine-grained results confirm that Ghost successfully shifts exposure from over-recommended head items to underexposed mid-tail and tail items without sacrificing general utility. Specifically, we can draw the following three insights.

**First**, Ghost demonstrates an exceptional capability to excavate and accurately recommend items across the broad long-tail distribution. For the intermediate and long-tail segments (specifically the 20-40, 40-60, and 60-80 groups), Ghost consistently achieves the highest Hit@10 and NDCG@10 scores among all evaluated models. For example, in the 20-40 group, Ghost reaches a Hit@10 of 0.0273, substantially outperforming the standard LC-Rec baseline (0.0070) and the strongest fairness-aware baseline, IFair-RR (0.0226). This proves the model’s targeted effectiveness in surfacing relevant, less-mainstream content. **Second**, Ghost effectively suppresses the over-exposure of head items, leading to significantly enhanced recommendation fairness. This is clearly evidenced by the steep reductions in Average Recommendation Popularity (ARP) and Monopoly Gini of Users (MGU) metrics. Ghost lowers the ARP@10 to 212.73 and MGU@10 to 0.0462, representing a massive drop compared to LC-Rec’s 322.14 and 0.1079, respectively. While the heuristic “Replace” baseline achieves slightly lower ARP and MGU scores, it does so by indiscriminately swapping items, which severely damages recommendation quality.

#### APPENDIX E SUPPLEMENT HYPER-PARAMETER ANALYSIS

##### A. Learning rate

Here, we investigate the impact of the learning rate during Ghost optimization. Figure 11 illustrates the sensitivity of the Ghost model to varying learning rates on the *Ins* dataset, revealing its crucial role in balancing overall recommendation accuracy with popularity bias mitigation. The empirical trends identify  $3 \times 10^{-5}$  as the optimal learning rate configuration. At this

TABLE VI  
HIDDEN STATE DIMENSIONS OF QWEN3 SERIES.

Scale	Hidden Dimension
0.6B	1024
1.7B	2048
4B	2560
8B	4096

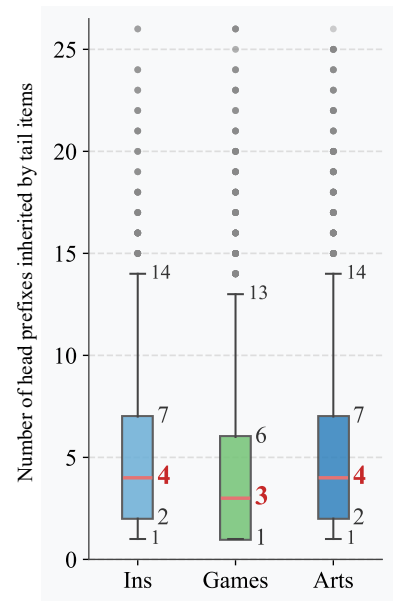


Fig. 9. Number of tail items that inherit SID prefix from the same head items.

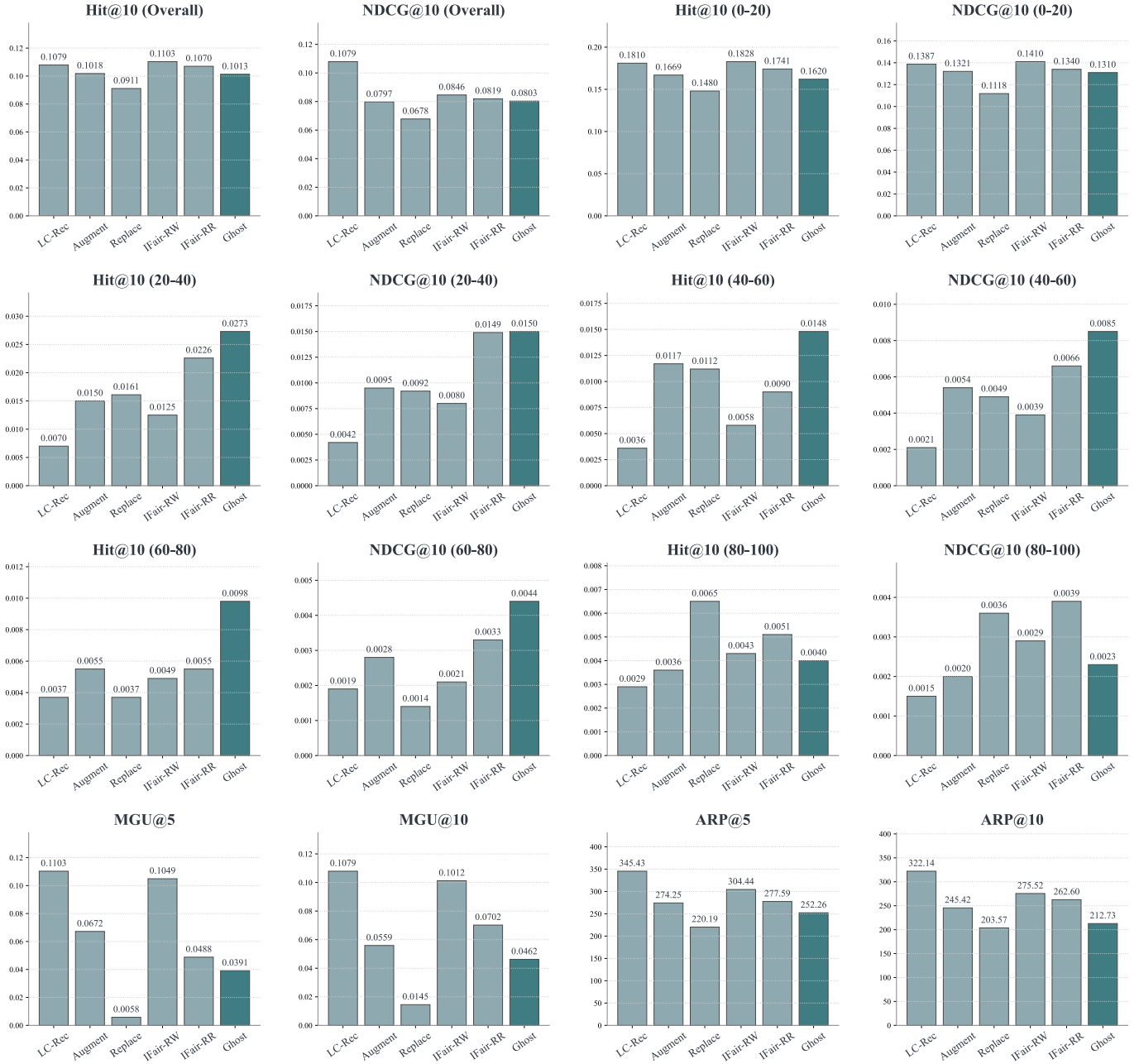


Fig. 10. Performance comparison of each equal-sized grouping on Ins dataset.

specific point, the model attains peak effectiveness in long-tail item excavation, evidenced by the highest Tail HR@10 and Tail NDCG@10, while simultaneously maximizing user coverage (MGU@10). Crucially, the ARP@10 reaches its global minimum at this setting, indicating a strong suppression of popularity bias without severely compromising the highly competitive overall accuracy metrics (Overall HR@10 and NDCG@10, which peak slightly earlier around  $1 \times 10^{-5}$ ). Conversely, employing an excessively large learning rate (e.g.,  $5 \times 10^{-4}$ ) yields detrimental effects across all dimensions: overall and tail-specific accuracies sharply degrade, user coverage diminishes, and ARP@10 spikes significantly. This demonstrates that an overly aggressive learning rate destabilizes the debiasing mechanism, causing the model to collapse back into disproportionately favoring mainstream items.

### B. Epoch

Here, we investigate the impact of the optimization epoch. Figure 12 visualizes the impact of the number of optimization epochs on the Ghost model performance using the Games dataset. The empirical trajectories indicate that determining an appropriate stopping point is crucial for balancing recommendation accuracy with popularity bias mitigation. Specifically, epoch 4 emerges as the optimal training duration, where both overall and tail-specific accuracy metrics (HR@10 and NDCG@10)

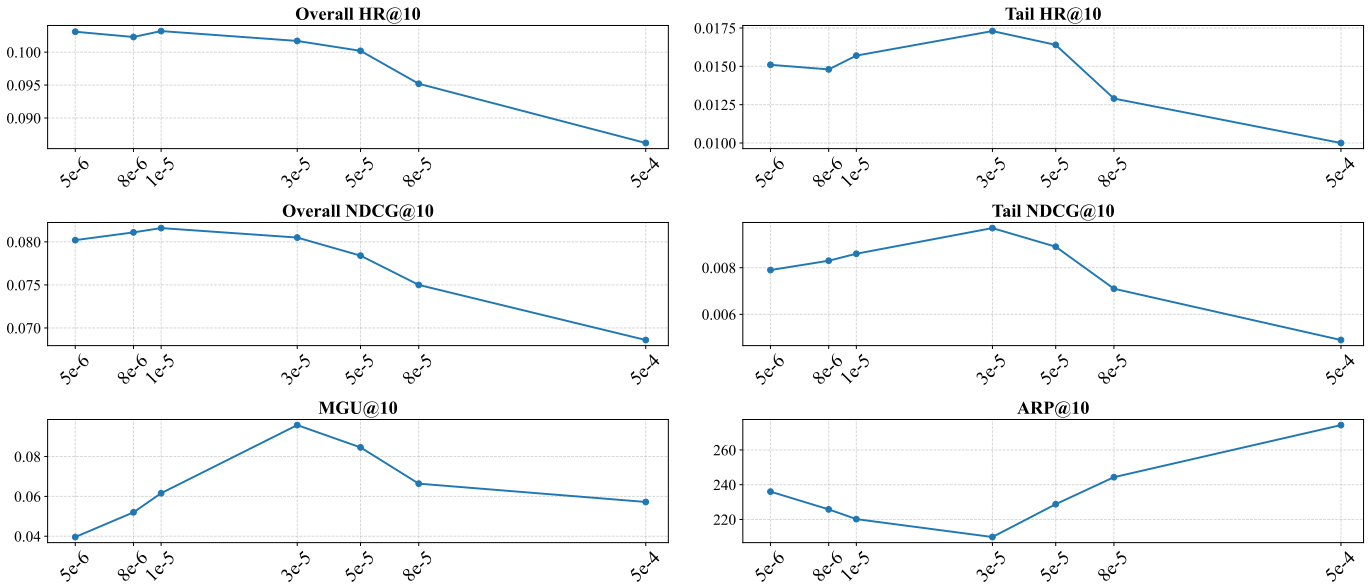


Fig. 11. Tendency of Ghost performance on Ins dataset, under different learning rates. The  $x$ -axis denotes learning rate, and the  $y$ -axis is the metric values.

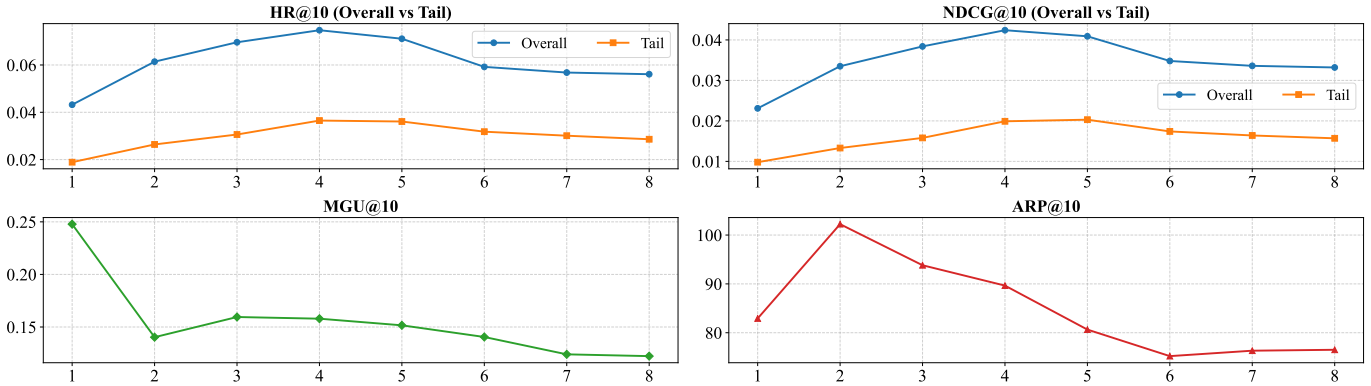


Fig. 12. Tendency of Ghost performance on Games dataset, under different optimization epochs. The  $x$ -axis denotes the epoch values, and the  $y$ -axis is the metric values.

simultaneously achieve their peak performance. During the initial training phase (epochs 1 and 2), the model exhibits suboptimal accuracy and a sharp spike in Average Recommendation Popularity (ARP@10), despite an initially high user coverage (MGU@10) that quickly drops. As optimization progresses towards epoch 4, the model effectively converges, significantly enhancing long-tail item retrieval while steadily driving down ARP@10. Conversely, prolonged training (epochs 6 through 8) evidently leads to overfitting; this over-optimization results in a consistent degradation across both overall and tail accuracies, a continued decline in MGU@10, and yields no further meaningful benefits for popularity debiasing. Therefore, maintaining a moderate number of optimization epochs is essential to maximize the framework’s capacity to deliver accurate, balanced, and diverse recommendations.

## APPENDIX F THEORETICAL DERIVATIONS AND JUSTIFICATIONS

This section provides the rigorous theoretical foundations and detailed mathematical proofs for the lemmas and corollaries presented in the main text regarding the Ghost model. Specifically, it elucidates the gradient starvation issue inherent in Maximum Likelihood Estimation (MLE), the localized bias amplification effect induced by undifferentiated tokenization, and the mathematical mechanisms through which Skeleton-Founded Tokenization (SKT) and Asymmetric Unlikelihood Optimization (AUO) structurally mitigate popularity bias.

### A. Prerequisites and Theoretical Foundations

To establish a rigorous theoretical framework for diagnosing popularity bias in Generative Recommenders (GRs), we first formalize three fundamental prerequisites regarding the data distribution and the architecture of SID-based GRs.

- **Prerequisite 1 (Long-tail Distribution).** In the sequential recommendation training dataset  $\mathcal{D}$ , the frequency of user-item interactions exhibits a heavily skewed, long-tailed distribution (e.g., a power-law distribution). Consequently, the occurrence probability of items in the head set  $\mathcal{V}_{\text{head}}$  (the top 20% most popular items) is vastly greater than that of items in the tail set  $\mathcal{V}_{\text{tail}}$  (the remaining 80%), i.e.,  $\mathbb{P}_{\mathcal{D}}(v \in \mathcal{V}_{\text{head}}) \gg \mathbb{P}_{\mathcal{D}}(v \in \mathcal{V}_{\text{tail}})$ .
- **Prerequisite 2 (Inner Product Mapping).** Current SOTA GR architectures formulate the next-token prediction task by computing generative logits via the inner product of the user historical behavior representation  $X_{h_u}$  and the candidate token embedding  $e_c$ . The generation probability is formalized via a Softmax operation across the token vocabulary,

$$\mathcal{P}_{\theta}(c|h_u) = \frac{\exp(\langle e_c, X_{h_u} \rangle)}{\sum_{c'} \exp(\langle e_{c'}, X_{h_u} \rangle)} \quad (24)$$

- **Assumption 1 (Identical Token Space under Undifferentiated Tokenization).** Current vector quantization techniques (e.g., RQ-VAE, RQ-KMeans) process head and tail items indiscriminately, meaning they share and compete within the identical representation space. Let  $c_{\text{tail}}$  denote a token structurally exclusive to tail items. Due to the heavy-tailed interaction distribution established in Prerequisite 1, the mathematical expectation of  $c_{\text{tail}}$  acting as a ground-truth target token within the training distribution  $\mathcal{D}$  asymptotically approaches zero.

### B. Detailed Derivations of LEMMA 1 and COROLLARY 1

#### Proof of LEMMA 1 (Gradient Starvation in MLE).

Basically, the model is optimized using standard Maximum Likelihood Estimation, specifically minimizing the Negative Log-Likelihood loss over the sequential tokens of the target item SID.

For an arbitrary generation step  $i$  and target item  $v$ , the instantaneous loss is

$$\mathcal{L}_{\text{NLL}}^{(i)} = -\log \mathcal{P}_{\theta}(c_v^{(i)}|h_u, c_v^{<i}).$$

According to the standard derivative properties of the Softmax function in Eq.(24), the partial derivative of the prediction probability with respect to an arbitrary candidate token embedding  $e_c$  yields a Jacobian formulation that separates into target and non-target cases. Consequently, the gradient of the NLL loss with respect to  $e_c$  is,

$$\frac{\partial \mathcal{L}_{\text{NLL}}^{(i)}}{\partial e_c} = \left( \mathcal{P}_{\theta}(c|h_u, c_v^{<i}) - \mathbb{I}\{c = c_v^{(i)}\} \right) X_{h_u} \quad (25)$$

In gradient descent, the parameter update direction opposes the computed gradient, denoted as  $\Delta e_c \propto -\frac{\partial \mathcal{L}}{\partial e_c}$ . Taking the mathematical expectation over the entire training distribution  $\mathcal{D}$ , we obtain,

$$\mathbb{E}_{\mathcal{D}}[\Delta e_c] \propto \mathbb{E}_{\mathcal{D}} \left[ \sum_i \left( \mathbb{I}\{c = c_v^{(i)}\} - \mathcal{P}_{\theta}(c|h_u, c_v^{<i}) \right) X_{h_u} \right] \quad (26)$$

Now, consider a tail-specific token  $c_{\text{tail}}$ . Based on Prerequisite 1 and Assumptions 1, its probability of being sampled as the ground-truth target token within any batch drawn from  $\mathcal{D}$  is negligible,  $\mathbb{P}_{\mathcal{D}}(c_{\text{tail}} = c_v^{(i)}) \approx 0$ . Consequently, the indicator function term  $\mathbb{I}\{c_{\text{tail}} = c_v^{(i)}\}$  reliably vanishes. By projecting the expected gradient update onto the user preference vector  $X_{h_u}$ , we isolate the token alignment trajectory, formulated as,

$$\mathbb{E}_{\mathcal{D}}[\langle \Delta e_{c_{\text{tail}}}, X_{h_u} \rangle] \approx -\mathbb{E}_{\mathcal{D}} \left[ \sum_i \mathcal{P}_{\theta}(c_{\text{tail}}|h_u, c_v^{<i}) \cdot \|X_{h_u}\|_2^2 \right] \leq 0 \quad (27)$$

**Conclusion.** Eq.(27) rigorously proves that tail tokens consistently act as trivial negative samples within the Softmax denominator during MLE optimization. Lacking positive reinforcement from the indicator function, their embeddings  $e_{c_{\text{tail}}}$  are pathologically pushed away from the user intent space, trapping them in a state of Gradient Starvation [31].

#### Proof of COROLLARY 1 (Head Token Dominance at Branching Point).

Building upon LEMMA 1, the parameter space is assumed to have converged to a skewed state where head tokens have received massive positive updates due to high interaction frequencies, i.e.,  $\mathbb{I}\{c_{\text{head}} = c_v^{(i)}\} = 1$  frequently, while tail tokens are starved.

This optimization skew manifests in the representation space as an extreme disparity in inner products,  $\langle e_{c_{\text{head}}}, X_{h_u} \rangle \gg \langle e_{c_{\text{tail}}}, X_{h_u} \rangle$ . Let  $i$  denote an arbitrary generative branching point where head and tail candidate tokens structurally compete. The predicted probability ratio for these tokens is given by,

$$\frac{\mathcal{P}_{\theta}(c_{\text{head}}^{(i)}|h_u, c^{<i})}{\mathcal{P}_{\theta}(c_{\text{tail}}^{(i)}|h_u, c^{<i})} = \frac{\exp(\langle e_{c_{\text{head}}}^{(i)}, X_{h_u} \rangle)}{\exp(\langle e_{c_{\text{tail}}}^{(i)}, X_{h_u} \rangle)} = \exp \left( \langle e_{c_{\text{head}}}^{(i)} - e_{c_{\text{tail}}}^{(i)}, X_{h_u} \rangle \right) \quad (28)$$

Because the gradient updates heavily favor the head token, the exponential term in Eq.(28) is pathologically amplified. This causes the predictive distribution  $\mathcal{P}_\theta$  to diverge significantly from the true underlying data distribution  $\mathcal{P}_d$ . We formally define this local divergence as the amplification factor  $\gamma_i$  as follows,

$$\gamma_i = \frac{\mathcal{P}_\theta(c_{\text{head}}^{(i)}|h_u, c^{<i})/\mathcal{P}_\theta(c_{\text{tail}}^{(i)}|h_u, c^{<i})}{\mathcal{P}_d(c_{\text{head}}^{(i)}|h_u, c^{<i})/\mathcal{P}_d(c_{\text{tail}}^{(i)}|h_u, c^{<i})} > 1 \quad (29)$$

Eq.(29) demonstrates the inevitable probability dominance of head tokens, establishing that the generative process becomes overconfident in predicting head tokens, regardless of context.

### C. Detailed Derivations of LEMMA 2 and LEMMA 3

#### Proof of LEMMA 2 (Bias Amplification via Undifferentiated Tokenization).

Basically, the current item tokenization method assigns Semantic Indices (SIDs) of length  $L$  indiscriminately, treating all items equally without accounting for popularity.

Under undifferentiated tokenization, a long-tail item  $v_{\text{tail}}$  shares prefixes of varying lengths with numerous head items. Consequently, generating  $v_{\text{tail}}$  requires decoding a sequence of tokens where it must repeatedly survive competition against popular items.

Suppose there exists a set  $\mathcal{Z}$  containing  $z$  unpredictable branching points ( $|\mathcal{Z}| = z \leq L$ ) where tail tokens structurally compete against head tokens. Based on the chain rule of autoregressive sequence generation and the dominance established in Corollary 1, each time the generation navigates a competitive branching point  $j \in \mathcal{Z}$ , the relative generation probability of the tail token is suppressed by at least a factor of  $\gamma_j$ ,

$$\mathcal{P}_\theta(c_{\text{tail}}^{(j)}|h_u, c^{<j}) \leq \gamma_j^{-1} \cdot \mathcal{P}_d(c_{\text{tail}}^{(j)}|h_u, c^{<j}) \quad (30)$$

Over the entire generation sequence of  $v_{\text{tail}}$ , we multiply the conditional probabilities. The suppressions at branching points accumulate geometrically,

$$\begin{aligned} \mathcal{P}_\theta(v_{\text{tail}}|h_u) &= \prod_{j=1}^L \mathcal{P}_\theta(c_{\text{tail}}^{(j)}|h_u, c^{<j}) \\ &\leq \left( \prod_{j \in \mathcal{Z}} \gamma_j^{-1} \right) \prod_{j=1}^L \mathcal{P}_d(c_{\text{tail}}^{(j)}|h_u, c^{<j}) \\ &\leq (\gamma_{\min})^{-z} \prod_{j=1}^L \mathcal{P}_d(c_{\text{tail}}^{(j)}|h_u, c^{<j}) \end{aligned} \quad (31)$$

where  $\gamma_{\min} = \min_{j \in \mathcal{Z}}(\gamma_j) > 1$ .

**Conclusion.** This derivation establishes that undifferentiated tokenization cascades localized, token-level gradient bias into a macroscopic, geometric  $\mathcal{O}(\gamma_{\min}^z)$  probability suppression for tail items, fully explaining their severe marginalization in recommendation lists.

#### Proof of LEMMA 3 (Mitigation of Bias Amplification via SKT).

The Skeleton-Founded Tokenization (SKT) mechanism asynchronously defines SIDs. The head items dictate the  $L^h$ -length skeleton of the SID space. A tail item  $v'$  is forced to explicitly inherit the first  $L^h$  tokens from its semantically closest head item  $v^*$ , and subsequently generates  $L^t$  specific tokens to characterize its distinctiveness.

During the initial generation steps  $j \in [1, L^h]$ , the tail item  $v'$  and the head item  $v^*$  share an identical prefix trajectory ( $c_{v'}^{(j)} = c_{v^*}^{(j)}$ ). Consequently, no probability divergence or chaotic token competition occurs between them. The unstructured branching points that comprised set  $\mathcal{Z}$  in LEMMA 2 are entirely eliminated. The genuine structural divergence is uniformly deferred to the single, predictable step ( $L^h + 1$ ). At this exact locus, the head item outputs an End-Of-Sequence (EOS) token, whereas the tail item generates the first token of its distinct semantic prefix. Because the competition is now strictly restricted to this single step, the cardinality of the branching set  $z$  is explicitly limited to 1.

The geometric suppression series detailed in Eq.(31) therefore collapses into a singular, localized deviation bounded by the head-dominance factor at the EOS step ( $\gamma_{EOS}$ ):

$$\mathcal{P}_\theta(v'|h_u) = \prod_{j=1}^{L^h+L^t} \mathcal{P}_\theta(c_{v'}^{(j)}|h_u, c^{<j}) \approx (\gamma_{EOS})^{-1} \prod_{j=1}^{L^h+L^t} \mathcal{P}_d(c_{v'}^{(j)}|h_u, c^{<j}) \quad (32)$$

**Conclusion.** This rigorously proves that SKT effectively halts the Markovian amplification chain of popularity bias. By establishing a unified structural branching point, it transforms a multi-step geometric suppression  $\mathcal{O}(\gamma_{\min}^z)$  into an insulated, single-step discrepancy  $\mathcal{O}(\gamma_{EOS})$ .

#### D. Detailed Derivations of LEMMA 4

##### Proof of LEMMA 4 (Gradient Rescue based on AUO).

The Ghost model incorporates Asymmetric Unlikelihood Optimization (AUO) to actively penalize a dynamically generated set of undesired tokens  $\bar{\Omega}$ . For a target tail item  $v'$ ,  $\bar{\Omega}$  consists of SIDs from head items that share high textual similarity with  $v'$  but possess divergent SID structures, acting as deceptive popular distractions.

The AUO loss function introduces an explicit penalty for generating tokens in  $\bar{\Omega}$ ,

$$\mathcal{L}_{\text{AUO}} = - \sum_{c \in \bar{\Omega}} \log(1 - \mathcal{P}_\theta(c)) \quad (33)$$

To derive the partial derivative with respect to an arbitrary candidate embedding  $e_k$ , let the logit be defined as  $z_k = \langle e_k, X_{h_u} \rangle$ . We first utilize the Jacobian matrix of the Softmax function for individual output probabilities,

- For the diagonal term ( $i = k$ ),

$$\frac{\partial \mathcal{P}_\theta(c_i)}{\partial z_k} = \mathcal{P}_\theta(c_i)(1 - \mathcal{P}_\theta(c_i))$$

- For the off-diagonal terms ( $i \neq k$ ),

$$\frac{\partial \mathcal{P}_\theta(c_i)}{\partial z_k} = -\mathcal{P}_\theta(c_i)\mathcal{P}_\theta(c_k)$$

Applying the chain rule to the AUO objective yields:

$$\frac{\partial \mathcal{L}_{\text{AUO}}}{\partial z_k} = \sum_{c_i \in \bar{\Omega}} \frac{1}{1 - \mathcal{P}_\theta(c_i)} \frac{\partial \mathcal{P}_\theta(c_i)}{\partial z_k} \quad (34)$$

##### Scenario 1: For a false positive head token $c_{\text{head}}^- \in \bar{\Omega}$ (Active Suppression).

Here, the embedding of interest  $k = c_{\text{head}}^-$  belongs to the actively penalized set. We separate the summation in Eq.(34) into terms where  $i = k$  and  $i \neq k$ ,

$$\frac{\partial \mathcal{L}_{\text{AUO}}}{\partial z_{\text{head}}^-} = \frac{\mathcal{P}_\theta(c_{\text{head}}^-)(1 - \mathcal{P}_\theta(c_{\text{head}}^-))}{1 - \mathcal{P}_\theta(c_{\text{head}}^-)} + \sum_{c_i \in \bar{\Omega} \setminus \{c_{\text{head}}^-\}} \frac{-\mathcal{P}_\theta(c_i)\mathcal{P}_\theta(c_{\text{head}}^-)}{1 - \mathcal{P}_\theta(c_i)} \quad (35)$$

Simplifying the first term, we get,

$$\frac{\partial \mathcal{L}_{\text{AUO}}}{\partial z_{\text{head}}^-} = \mathcal{P}_\theta(c_{\text{head}}^-) - \sum_{c_i \neq c_{\text{head}}^-} \frac{\mathcal{P}_\theta(c_i)\mathcal{P}_\theta(c_{\text{head}}^-)}{1 - \mathcal{P}_\theta(c_i)} \quad (36)$$

The overall objective is a weighted combination of MLE and AUO ( $\mathcal{L} = \mathcal{L}_{\text{NLL}} + \alpha \mathcal{L}_{\text{AUO}}$ ). Incorporating the standard MLE base gradient (which negatively updates  $c_{\text{head}}^-$  as it is a false negative, not the target), the overall parameter update direction becomes,

$$\Delta e_{c_{\text{head}}^-} \propto -(1 + \alpha)\mathcal{P}_\theta(c_{\text{head}}^-)X_{h_u} + \alpha \sum_{c_j \in \bar{\Omega} \setminus \{c_{\text{head}}^-\}} \frac{\mathcal{P}_\theta(c_j)\mathcal{P}_\theta(c_{\text{head}}^-)}{1 - \mathcal{P}_\theta(c_j)} X_{h_u} \quad (37)$$

This explicitly shows the direct penalty via the softmax derivative.

##### Scenario 2: For a target tail token $c_{\text{tail}} \notin \bar{\Omega}$ (Cross-Penalization Rescue).

Here, the embedding of interest  $k = c_{\text{tail}}$  corresponds to our ground-truth tail item, which is specifically excluded from the explicitly penalized head set  $\bar{\Omega}$ . Therefore, the condition  $i \neq k$  strictly holds for *all* terms in the summation of Eq.(34),

$$\frac{\partial \mathcal{L}_{\text{AUO}}}{\partial z_{\text{tail}}} = \sum_{c_j \in \bar{\Omega}} \frac{-\mathcal{P}_\theta(c_j)\mathcal{P}_\theta(c_{\text{tail}})}{1 - \mathcal{P}_\theta(c_j)} \quad (38)$$

In gradient descent, the specific parameter update contribution derived from AUO towards this tail token acts in the opposing direction ( $-\alpha \cdot \partial \mathcal{L}_{\text{AUO}} / \partial z_{\text{tail}}$ ),

$$\Delta e_{c_{\text{tail}}}^{\text{AUO}} \propto +\alpha \sum_{c_j \in \bar{\Omega}} \frac{\mathcal{P}_\theta(c_j)\mathcal{P}_\theta(c_{\text{tail}})}{1 - \mathcal{P}_\theta(c_j)} X_{h_u} \quad (39)$$

Merging this positive AUO rescue term with the suppressive MLE base term established in LEMMA 1 rigorously completes the derivation:

$$\Delta e_{c_{\text{tail}}} \propto -\mathcal{P}_\theta(c_{\text{tail}})X_{h_u} + \alpha \sum_{c_j \in \bar{\Omega}} \frac{\mathcal{P}_\theta(c_j)\mathcal{P}_\theta(c_{\text{tail}})}{1 - \mathcal{P}_\theta(c_j)} X_{h_u} \quad (40)$$

**Conclusion.** This Jacobian analysis illuminates how the AUO loss function systematically transfigures the penalization of notorious head tokens into a positive, structural rescue force for long-tail tokens. By redistributing probability mass, this mechanism actively counteracts and dismantles the gradient starvation trap established in LEMMA 1, ensuring tail tokens receive rational parameter updates.

## APPENDIX G LIMITATIONS AND FUTURE WORK

- **Supervised Finetuning.** While this study provides a comprehensive diagnosis and mitigation strategy for popularity bias in GRs, the scope of our investigation is currently constrained to the supervised fine-tuning (SFT) paradigm. Specifically, our theoretical analysis of the gradient starvation issue and the subsequent design of the AUO are fundamentally grounded in the MLE framework. Recent advancements in generative modeling have increasingly adopted reinforcement learning (RL) techniques to align models with complex objectives [47], [48]. The specific mechanisms by which popularity bias manifests in RL-based GRs, and whether our proposed skeleton-founded tokenization and unlikelihood penalties remain effective under reward-driven optimization, have not yet been analyzed and represent a critical direction for future research.

- **Pre-computed Undesired Collection.** The undesired collection utilized in the AUO module relies on a static, pre-computed approach. As detailed in our methodology, the undesired items are identified by retrieving head items with high initial semantic similarity but divergent SIDs relative to the target tail item. Currently, this collection is established prior to training and is not dynamically adjusted as the model’s parameters and internal representations evolve during the optimization process. While dynamically updating the undesired collection at each epoch could theoretically provide more precise and adaptive negative supervision, the static approach was deliberately adopted as a pragmatic trade-off to preserve computational efficiency and prevent prohibitive training overhead.