

Geodesic Expert Routing for Unbiased Knowledge Distillation in Recommendation

Xuan Zhang¹, Rongchuan Wei¹, Chunyu Wei^{2,*}, Hongxing Yuan^{1,*} and Yushun Fan¹

¹Department of Automation, Tsinghua University, Beijing National Research Center for Information Science and Technology, China

²School of Information, Renmin University of China, China

{xuan-zha23, wrc20, yuanhx24}@mails.tsinghua.edu.cn, weichunyu@ruc.edu.cn, fanyus@tsinghua.edu.cn

Abstract

Knowledge distillation has become a prevalent technique for deploying efficient recommender systems, enabling lightweight student models to approximate the performance of larger teachers. However, we identify a critical issue: distillation systematically amplifies popularity bias, as student models inherit and intensify the popularity-driven shortcuts encoded in teachers trained on interaction data dominated by popular items. To address this limitation, we propose GUIDE (Geodesic aware Unbiased Instructive Distillation with Experts), a collaborative distillation framework that incorporates domain-specific debiasing experts alongside the global teacher. GUIDE tackles two key challenges in this paradigm. First, for expert routing, we introduce Spherical Expert Alignment, which conducts expert-student matching on the spherical manifold with geodesic distance optimization, eliminating magnitude-induced bias and ensuring stable gradient flow. Second, for context fusion, a Meta-Debiasing Gate is designed to dynamically arbitrate teacher-expert influence using real-time context via end-to-end amortized meta-learning. Extensive experiments on multiple real-world datasets demonstrate that GUIDE significantly mitigates popularity bias while preserving recommendation accuracy, with state-of-the-art trade-offs among efficiency, accuracy, and fairness.

1 Introduction

Recommender systems have become indispensable infrastructure for modern online platforms, serving as the primary mechanism through which users discover products, content, and services aligned with their preferences [Deldjoo *et al.*, 2024]. By filtering the overwhelming volume of available options into personalized suggestions, these systems significantly enhance user experience across diverse domains such as e-commerce, streaming media, and social networks [Park

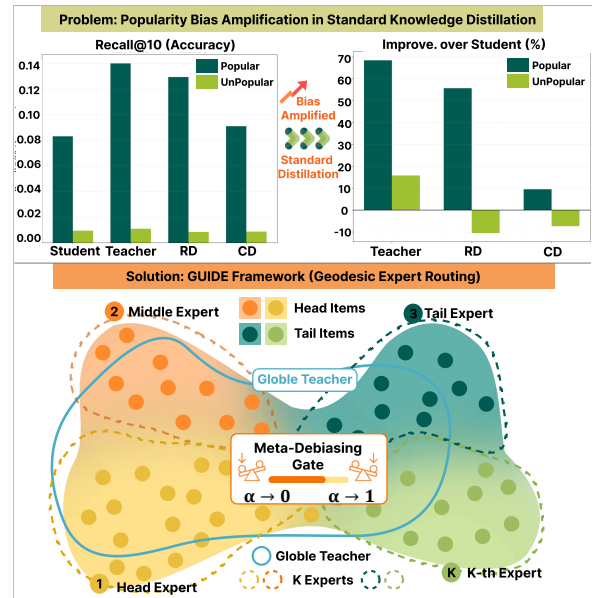


Figure 1: (Top) Recall@10 comparison on popular/unpopular items. Relative improvements over the base Student (directly learned from the data) reveal that standard distillation amplifies popularity bias. (Bottom) The proposed GUIDE dynamically fuses knowledge from global and expert teachers to balance accuracy and fairness.

et al., 2026; Gao *et al.*, 2023]. As data volume and item diversity continue to scale exponentially, system designers face an increasingly critical trade-off between recommendation accuracy and computational efficiency [Zhang *et al.*, 2019].

Knowledge Distillation (KD) has emerged as a widely adopted paradigm to address this efficiency challenge [Kang *et al.*, 2020]. By transferring knowledge from a computationally expensive teacher model to a lightweight student model, KD enables the deployment of compact recommenders that retain much of the predictive capability of their larger counterparts while satisfying stringent latency requirements [Lee *et al.*, 2019; Kweon *et al.*, 2021]. This approach has demonstrated remarkable success across various recommendation architectures, from collaborative filtering to sequential models [Sarwar *et al.*, 2001; Kang *et al.*, 2021]. However, we identify a critical yet underexplored issue: *knowledge distillation systematically amplifies popularity bias in the stu-*

*Corresponding author

The code and data are available at: <https://github.com/zx19971219/GUIDE>

55 *dent model*. Popularity bias, where models disproportion- 113
56 ately favor frequently-interacted items at the expense of long- 114
57 tail content, is already pervasive in recommender systems 115
58 trained on implicit feedback [Oestreicher-Singer and Sun- 116
59 dararajan, 2012; Abdollahpouri *et al.*, 2019]. Our empir- 117
60 ical investigation reveals that this bias is not merely pre- 118
61 served but substantially magnified during the distillation pro- 119
62 cess. We argue that the root cause lies in the reliance on 120
63 a single global teacher trained on complete interaction data 121
64 that is inherently dominated by popularity patterns. Figure 1 122
65 shows KD improvements are driven by popular items at the 123
66 expense of unpopular ones. By mimicking such a teacher, 124
67 the student internalizes and amplifies these popularity-driven 125
68 shortcuts, resulting in reduced exposure diversity and ex- 126
69 acerbated unfairness toward niche items [Lee *et al.*, 2019; 127
70 Tang and Wang, 2018].

71 To address this limitation, we propose **GUIDE** (Geodesic- 128
72 aware **U**nbiased **I**nstructive **D**istillation with **E**xperts), a novel 129
73 framework that augments conventional single-teacher distil- 130
74 lation with collaborative guidance from domain-specific de- 131
75 biasing experts. Our key insight is that while the global 132
76 teacher provides comprehensive knowledge essential for ac- 133
77 curacy, targeted experts trained on strategically partitioned 134
78 data subsets can serve as corrective signals that counteract 135
79 popularity bias. Specifically, we partition the original inter- 136
80 action data into multiple subsets based on item popularity 137
81 strata and train specialized expert models in parallel. Dur- 138
82 ing distillation, the student model learns not only from the 139
83 global teacher but also from these experts, enabling balanced 140
84 knowledge acquisition that promotes both accuracy and fair-
85 ness. However, transitioning from single-teacher distillation
86 to this collaborative teacher-expert paradigm introduces two
87 fundamental challenges:

- 88 • **Expert Routing**: How to select the most appropriate ex- 141
89 pert for each prediction context? A natural approach 142
90 measures expert-student similarity via cosine or inner 143
91 product in Euclidean space [Koren *et al.*, 2009]. How- 144
92 ever, embedding norms encode popularity-correlated in- 145
93 tensity information, as popular items naturally exhibit 146
94 larger magnitudes due to more frequent gradient up- 147
95 dates [Zhu *et al.*, 2021; Wang *et al.*, 2017]. This causes 148
96 the student to be attracted to high-norm experts rather 149
97 than semantically aligned ones. While post-hoc ℓ_2 nor- 150
98 malization offers an implicit remedy, it remains geomet- 151
99 rically imprecise and optimization-unfriendly.
- 100 • **Context Fusion**: How to dynamically balance knowl- 152
101 edge from the teacher and experts based on varying user- 153
102 item characteristics? Static fusion weights are funda- 154
103 mentally inadequate, since susceptibility to popularity 155
104 bias varies substantially across user-item contexts de- 156
105 pending on user engagement diversity, item popular- 157
106 ity characteristics, and interaction patterns [Han *et al.*, 158
107 2024; Abdollahpouri *et al.*, 2021]. Determining appro- 159
108 priate debiasing intensity requires real-time contextual 160
109 reasoning rather than a fixed arbitration strategy.

110 To tackle **Expert Routing**, we propose *Spherical Expert*
111 *Alignment*, which conducts expert-student matching directly
112 on the spherical manifold \mathcal{S}^{d-1} , ensuring that all compar-

isons reflect pure directional similarity free from magnitude
contamination. Crucially, we optimize geodesic distance
rather than cosine angle: cosine-based gradients scale with
 $\sin(\theta)$ and vanish as $\theta \rightarrow 0$, whereas geodesic gradients re-
main linear in θ , enabling stable optimization even when stu-
dent and expert representations are closely aligned.

To address **Context Fusion**, we design a *Meta-Debiasing Gate*, a lightweight gating network that ingests contextual features such as user history diversity, item popularity percentile, and interaction density. The gate outputs a dynamic fusion weight $\alpha \in [0, 1]$ that arbitrates the balance between teacher and expert influence. It is trained end-to-end via amortized meta-learning to explicitly optimize for both accuracy and debiasing effectiveness.

Our main contributions are summarized as follows:

- We identify the phenomenon of popularity bias amplification in knowledge distillation and propose GUIDE, a collaborative distillation framework that integrates domain-specific debiasing experts alongside the global teacher to enable accuracy-fairness balanced knowledge transfer.
- We introduce Spherical Expert Alignment for geometry-aware expert routing and Meta-Debiasing Gate for context-aware dynamic fusion, addressing the challenges of expert selection and adaptive knowledge arbitration respectively.
- Extensive experiments on real-world datasets demonstrate that GUIDE significantly mitigates popularity bias while maintaining competitive accuracy, achieving state-of-the-art trade-offs among efficiency, accuracy, and fairness.

2 RELATED WORK 141

2.1 Knowledge Distillation for Recommendation 142

Knowledge Distillation (KD) enhances recommender systems by transferring knowledge from intermediate layers [Romero *et al.*, 2015], privileged features [Xu *et al.*, 2020], or structured graphs [Zhang *et al.*, 2020]. Beyond efficiency [Zhu *et al.*, 2020] and collaborative learning [Gou *et al.*, 2024], KD is widely applied to mitigate biases. For instance, [Liu *et al.*, 2020] and UKDRec [Yang *et al.*, 2025] leverage counterfactual data and multi-teacher frameworks for debiasing. UNKD [Chen *et al.*, 2023a] and [Ding *et al.*, 2022] address popularity bias and model fusion, respectively. Research also optimizes distillation losses, ranging from ranking-based top-N selection [Tang and Wang, 2018] and instance construction [Lee *et al.*, 2019] to list-level losses [Kang *et al.*, 2020], extending to heterogeneous graph integration [Tao *et al.*, 2022; Wang *et al.*, 2021].

2.2 Bias in Recommendation 158

Recommendation bias arises from user behavior and feedback loops [Chen *et al.*, 2023b]. Mitigation strategies typically employ IPS [Joachims *et al.*, 2018; Li *et al.*, 2023; Schnabel *et al.*, 2016], Doubly Robust methods [Song *et al.*, 2023], or causal graphs [Liu *et al.*, 2021; Liu *et al.*, 2023; He *et al.*, 2022] on biased data [Ning *et al.*, 2024; Zhang and Shen, 2023], or leverage unbiased data via meta-learning [Chen *et al.*, 2021] and distillation [Liu *et al.*, 2020]. Beyond data, architectures [Lin *et al.*, 2019] and optimizers [Tang *et*

168 *al.*, 2020] contribute to bias. Notably, KD itself can aggravate
 169 popularity bias [Chen *et al.*, 2023a]: a teacher model trained
 170 on skewed data transfers inherent biases to the student, un-
 171 dermining debiasing efforts [Yang *et al.*, 2025].

172 3 Preliminaries

173 3.1 Problem Formulation

174 Let $\mathcal{U} = \{u_1, u_2, \dots, u_M\}$ denote the set of M users and
 175 $\mathcal{I} = \{i_1, i_2, \dots, i_N\}$ denote the set of N items. The user-
 176 item interaction matrix $\mathbf{R} \in \{0, 1\}^{M \times N}$ encodes implicit
 177 feedback, where $r_{ui} = 1$ indicates that user u has interacted
 178 with item i , and $r_{ui} = 0$ otherwise. For each user u , we de-
 179 note the set of interacted items as $\mathcal{I}_u^+ = \{i \in \mathcal{I} \mid r_{ui} = 1\}$.
 180 The primary objective of a recommender system is to learn a
 181 scoring function $f : \mathcal{U} \times \mathcal{I} \rightarrow \mathbb{R}$ that predicts user preferences
 182 over unobserved items.

183 3.2 Knowledge Distillation in Recommendation

184 Knowledge Distillation (KD) aims to transfer knowledge
 185 from a high-capacity teacher model \mathcal{T} to a lightweight stu-
 186 dent model \mathcal{S} . Given a user-item pair (u, i) , the teacher pro-
 187 duces a soft prediction $\hat{y}_{ui}^{\mathcal{T}}$ that encodes rich relational knowl-
 188 edge beyond binary labels. The student is trained to mimic
 189 these soft targets while maintaining computational efficiency.
 190 The distillation objective typically minimizes the divergence
 191 between teacher and student outputs:

$$\mathcal{L}_{KD} = \sum_{(u,i) \in \mathcal{D}} \ell(\hat{y}_{ui}^{\mathcal{S}}, \hat{y}_{ui}^{\mathcal{T}}), \quad (1)$$

192 where $\ell(\cdot, \cdot)$ denotes a divergence measure such as Mean
 193 Squared Error or KL divergence, and \mathcal{D} is the training set.

194 3.3 Spherical Manifold

195 The $(d-1)$ -dimensional unit spherical manifold \mathcal{S}^{d-1} is de-
 196 fined as the set of all d -dimensional vectors with unit ℓ_2 norm:

$$\mathcal{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_2 = 1\}. \quad (2)$$

197 The spherical manifold provides an embedding space where
 198 vector magnitudes are normalized, enabling the model to
 199 focus purely on directional relationships. For any vector
 200 $\mathbf{x} \in \mathbb{R}^d \setminus \{0\}$, the projection onto \mathcal{S}^{d-1} is given by:

$$\pi(\mathbf{x}) = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}. \quad (3)$$

201 The intrinsic distance on the spherical manifold is the
 202 geodesic distance, representing the length of the shortest arc
 203 connecting two points. For $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{d-1}$, the distance is:

$$d_{\text{geo}}(\mathbf{x}, \mathbf{y}) = \arccos(\mathbf{x}^\top \mathbf{y}). \quad (4)$$

204 4 Methodology

205 Figure 2 illustrates the overall framework of GUIDE. Given
 206 user-item interaction data, GUIDE first partitions the item set
 207 into K popularity-based subsets to construct domain-specific
 208 expert teachers. To address the **Expert Routing** challenge,
 209 we project both users and expert centroids onto a spherical
 210 manifold, where expert relevance is quantified via geodesic

distance. To address the **Context Fusion** challenge, a Meta-
 211 Debiasing Gate dynamically arbitrates between expert and
 212 global teacher knowledge based on real-time contextual fea-
 213 tures. The final distillation signal is derived as a context-
 214 aware ensemble tailored to each user-item interaction. We
 215 detail each component in the following subsections. 216

217 4.1 Domain-Specific Expert Construction

218 Relying on a single global teacher causes popularity bias am-
 219 plification in knowledge distillation. A global teacher trained
 220 on the complete interaction data inevitably inherits the long-
 221 tail distribution, leading to overfitting on popular items while
 222 under-representing niche content. To mitigate this limitation,
 223 we construct multiple domain-specific experts, each special-
 224 izing in a distinct popularity stratum.

225 Popularity-Based Data Partitioning

226 We quantify the popularity of each item i by its interaction
 227 frequency $d_i = |\{u \in \mathcal{U} \mid r_{ui} = 1\}|$. Based on these fre-
 228 quency values, we partition the item set \mathcal{I} into K disjoint sub-
 229 sets $\{\mathcal{I}^{(1)}, \mathcal{I}^{(2)}, \dots, \mathcal{I}^{(K)}\}$, where each subset corresponds to
 230 a distinct popularity tier. Specifically, we sort all items by d_i
 231 in descending order and divide them into K groups of ap-
 232 proximately equal size, yielding a spectrum from head (high-
 233 popularity) to tail (low-popularity) items. Each item subset
 234 $\mathcal{I}^{(k)}$ induces a corresponding interaction subset:

$$\mathcal{D}^{(k)} = \{(u, i) \mid i \in \mathcal{I}^{(k)}, r_{ui} = 1\}. \quad (5)$$

235 This partitioning creates isolated training environments that
 236 shield tail item learning from head item interference.

237 Expert Teacher Training

238 For each partition $\mathcal{D}^{(k)}$, we train a dedicated expert teacher
 239 $\mathcal{T}^{(k)}$ exclusively on interactions within $\mathcal{I}^{(k)}$. This constraint
 240 ensures that each expert captures the intrinsic collaborative
 241 patterns specific to its popularity domain without being dom-
 242 inated by global popularity shortcuts. Formally, the expert
 243 $\mathcal{T}^{(k)}$ is optimized via:

$$\mathcal{T}^{(k)} = \arg \min_{\theta} \sum_{(u,i,j) \in \mathcal{D}^{(k)}} \mathcal{L}_{\text{BPR}}(f_{\theta}(u, i), f_{\theta}(u, j)), \quad (6)$$

244 where j denotes a negative sample drawn from $\mathcal{I}^{(k)}$. In paral-
 245 lel, we train a global teacher $\mathcal{T}^{(g)}$ on the complete interaction
 246 data \mathcal{D} to preserve comprehensive collaborative signals.

247 **Remark.** The expert construction procedure is model-
 248 agnostic: any recommendation architecture (e.g., matrix fac-
 249 torization, neural collaborative filtering, or sequential mod-
 250 els) can serve as the backbone for both expert and student
 251 models. This flexibility ensures broad applicability across di-
 252 verse recommendation scenarios.

253 4.2 Spherical Expert Alignment

254 With K expert teachers constructed, the next challenge lies
 255 in dynamically routing each user to the most semantically
 256 aligned expert. Conventional approaches measure expert-
 257 student similarity via cosine or inner product in Euclidean
 258 space. However, as discussed in Section 1, embedding mag-
 259 nitudes encode popularity-correlated intensity information

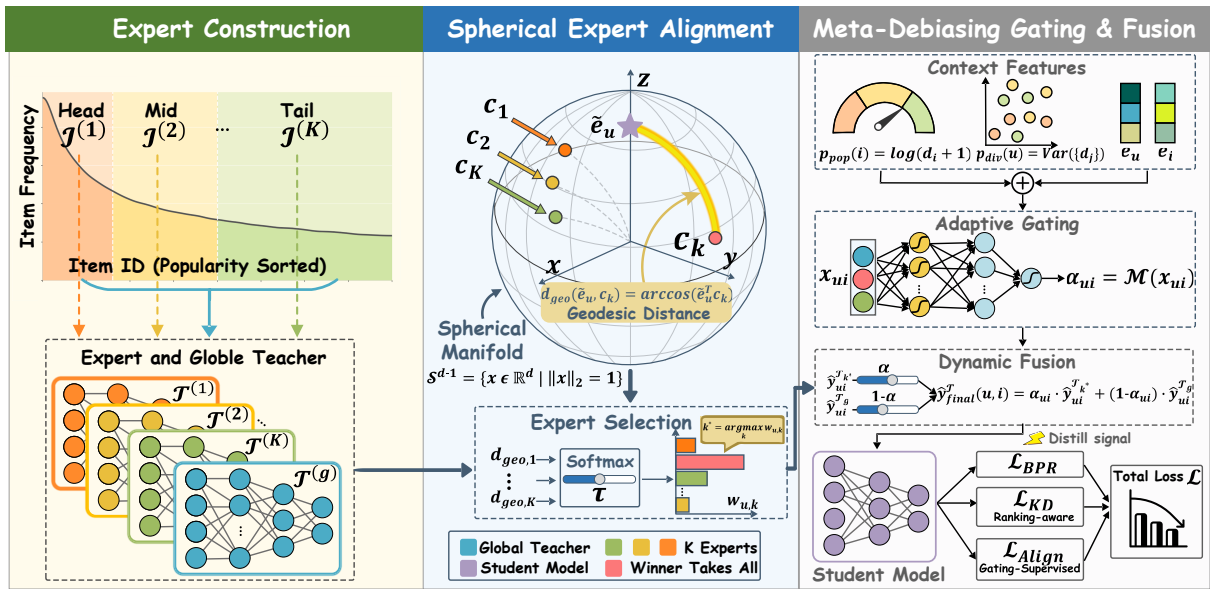


Figure 2: **Overall of GUIDE framework.** It consists of three key parts: (1) Expert Construction, which trains specialized teachers on popularity-based data partitions; (2) Spherical Expert Alignment, which routes users to experts via geodesic distance on a spherical manifold to eliminate magnitude bias; and (3) Meta-Debiasing Gating & Fusion, which adaptively fuses expert and global knowledge based on interaction context for personalized distillation.

260 that confounds directional similarity. We address this challenge through Spherical Expert Alignment, which conducts
 261 matching directly on the spherical manifold \mathcal{S}^{d-1} .
 262

263 Expert Centroid Representation

264 To enable efficient expert routing, we represent each expert $\mathcal{T}^{(k)}$ by a semantic centroid that summarizes its domain
 265 knowledge. Let $\mathbf{E}^{(k)} = \{\mathbf{e}_i \mid i \in \mathcal{I}^{(k)}\}$ denote the set of item embeddings learned by expert $\mathcal{T}^{(k)}$. We compute the expert
 266 centroid \mathbf{c}_k by aggregating these embeddings and projecting onto the spherical manifold:
 267
 268
 269

$$\mathbf{c}_k = \pi \left(\frac{1}{|\mathcal{I}^{(k)}|} \sum_{i \in \mathcal{I}^{(k)}} \mathbf{e}_i \right) = \frac{\sum_{i \in \mathcal{I}^{(k)}} \mathbf{e}_i}{\left\| \sum_{i \in \mathcal{I}^{(k)}} \mathbf{e}_i \right\|_2}, \quad (7)$$

270 where $\mathbf{c}_k \in \mathcal{S}^{d-1}$ serves as the representative prototype for the k -th expert domain.
 271

272 User Projection and Geodesic Matching

273 For any user u with embedding \mathbf{e}_u learned by the student model, we project it onto the same spherical manifold:
 274

$$\tilde{\mathbf{e}}_u = \pi(\mathbf{e}_u) = \frac{\mathbf{e}_u}{\|\mathbf{e}_u\|_2}. \quad (8)$$

275 The affinity between user u and expert k is then quantified by the geodesic distance on \mathcal{S}^{d-1} :
 276

$$d_{\text{geo}}(\tilde{\mathbf{e}}_u, \mathbf{c}_k) = \arccos(\tilde{\mathbf{e}}_u^\top \mathbf{c}_k). \quad (9)$$

277 This geodesic formulation offers two key advantages over cosine similarity. First, constraining representations to the unit
 278 sphere strictly decouples embedding magnitude—often correlated with item popularity—from semantic direction, ensuring
 279 pure content-based routing. Second, unlike cosine gradients which scale with $\sin(\theta)$ and vanish as $\theta \rightarrow 0$, the
 280
 281
 282

283 geodesic gradient $\nabla_v \mathcal{L}_{\text{geo}} \propto \frac{\partial \theta^2}{\partial \theta} \frac{\partial \theta}{\partial \cos \theta} \nabla_v \cos \theta$ introduces
 284 a $1/\sin(\theta)$ term that perfectly cancels the projection-induced
 285 $\sin(\theta)$ decay. This yields stable gradients linear in θ (i.e.,
 286 $\nabla \mathcal{L} \propto 2\theta$) even for closely aligned representations.

287 Expert Selection via Soft Routing

288 Based on geodesic distances, we compute a probability distribution over experts for each user u :
 289

$$w_{u,k} = \frac{\exp(-\tau \cdot d_{\text{geo}}(\tilde{\mathbf{e}}_u, \mathbf{c}_k))}{\sum_{j=1}^K \exp(-\tau \cdot d_{\text{geo}}(\tilde{\mathbf{e}}_u, \mathbf{c}_j))}, \quad (10)$$

290 where $\tau > 0$ is a temperature parameter controlling the sharpness of the distribution. To minimize noise from less relevant
 291 experts, we adopt a winner-takes-all strategy that selects the expert with the highest probability:
 292
 293

$$k^* = \arg \max_{k \in \{1, \dots, K\}} w_{u,k}. \quad (11)$$

294 The selected expert $\mathcal{T}^{(k^*)}$ provides the specialized soft label $\hat{y}_{ui}^{\mathcal{T}^{(k^*)}}$ for the current user-item pair.
 295

296 4.3 Meta-Debiasing Gate

297 While Spherical Expert Alignment identifies the optimal domain-specific expert for each user, relying solely on expert
 298 knowledge may be suboptimal. Expert teachers, trained on isolated popularity strata, may lack the broader collaborative
 299 signals captured by the global teacher. Conversely, the global teacher provides robust general knowledge but inherits popularity bias. To dynamically balance these complementary
 300 knowledge sources, we design a Meta-Debiasing Gate that adaptively arbitrates their influence based on real-time contextual features.
 301
 302
 303
 304
 305
 306

307 Context Feature Extraction

308 The optimal balance between expert and global teacher
309 knowledge depends on the specific characteristics of each
310 user-item interaction. We extract a context feature vector \mathbf{x}_{ui}
311 that captures two key factors:

312 **Item Popularity.** We quantify the global exposure of candi-
313 date item i using its log-normalized interaction frequency:

$$p_{\text{pop}}(i) = \log(d_i + 1). \quad (12)$$

314 This feature informs the gate about the item’s position in the
315 popularity spectrum.

316 **User Diversity.** We measure the diversity of user u ’s his-
317 torical interactions by computing the variance of popularity
318 values across interacted items:

$$p_{\text{div}}(u) = \text{Var}(\{d_j \mid j \in \mathcal{I}_u^+\}). \quad (13)$$

319 This feature indicates the user’s tendency to explore across
320 different popularity tiers. The context vector concatenates
321 these features with learned user and item embeddings, where
322 \oplus denotes concatenation:

$$\mathbf{x}_{ui} = [\mathbf{e}_u \oplus \mathbf{e}_i \oplus p_{\text{pop}}(i) \oplus p_{\text{div}}(u)]. \quad (14)$$

323 Adaptive Gating Mechanism

324 The Meta-Debiasing Gate \mathcal{M} uses a lightweight multi-layer
325 perceptron (MLP) with sigmoid activation, mapping the con-
326 text vector to a scalar fusion coefficient $\alpha_{ui} \in (0, 1)$:

$$\alpha_{ui} = \mathcal{M}(\mathbf{x}_{ui}) = \sigma(\mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \mathbf{x}_{ui} + \mathbf{b}_1) + \mathbf{b}_2), \quad (15)$$

327 where $\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2$ are learnable parameters, and $\sigma(\cdot)$
328 denotes the sigmoid function.

329 The coefficient α_{ui} quantifies the necessity of debiasing
330 for the current context: higher values indicate that the model
331 should prioritize expert knowledge (e.g., for long-tail items
332 or users with diverse histories), while lower values suggest
333 relying on the global teacher for stability.

334 Dynamic Knowledge Fusion

335 The final distillation target is synthesized as a convex combi-
336 nation of soft labels from the selected expert $\mathcal{T}^{(k^*)}$ and the
337 global teacher $\mathcal{T}^{(g)}$:

$$\hat{y}_{\text{final}}^{\mathcal{T}}(u, i) = \alpha_{ui} \cdot \hat{y}_{ui}^{\mathcal{T}^{(k^*)}} + (1 - \alpha_{ui}) \cdot \hat{y}_{ui}^{\mathcal{T}^{(g)}}. \quad (16)$$

338 This dynamic fusion ensures that the distillation signal is
339 adaptively tailored to each user-item interaction, effectively
340 resolving the trade-off between accuracy and fairness.

341 4.4 Optimization

342 We formulate a multi-objective loss function that jointly op-
343 timizes recommendation accuracy, knowledge transfer, and
344 gate calibration:

$$\mathcal{L} = \mathcal{L}_{\text{BPR}} + \lambda_{\text{KD}} \mathcal{L}_{\text{KD}} + \lambda_{\text{Align}} \mathcal{L}_{\text{Align}} + \lambda_{\text{Reg}} \|\Theta\|_2^2, \quad (17)$$

345 where $\lambda_{\text{KD}}, \lambda_{\text{Align}},$ and λ_{Reg} are hyperparameters balancing
346 each component.

347 **Recommendation Loss.** We employ the BPR loss to ensure
348 the student captures fundamental collaborative signals:

$$\mathcal{L}_{\text{BPR}} = \sum_{(u,i,j) \in \mathcal{D}} -\ln \sigma(\hat{y}_{ui}^{\mathcal{S}} - \hat{y}_{uj}^{\mathcal{S}}), \quad (18)$$

where (u, i, j) denotes a triplet with positive item $i \in \mathcal{I}_u^+$ and
negative item $j \notin \mathcal{I}_u^+$.

Ranking-Aware Distillation Loss. Standard pointwise distillation may fail to preserve the nuanced ranking relationships essential for recommendation. We propose a ranking-aware objective that aligns the preference margins between student and teacher:

$$\mathcal{L}_{\text{KD}} = \sum_{(u,i,j) \in \mathcal{D}} (\Delta p_{uij}^{\mathcal{S}} - \Delta p_{uij}^{\mathcal{T}})^2, \quad (19)$$

where $\Delta p_{uij}^{\mathcal{S}} = \sigma(\hat{y}_{ui}^{\mathcal{S}}) - \sigma(\hat{y}_{uj}^{\mathcal{S}})$ and $\Delta p_{uij}^{\mathcal{T}} = \sigma(\hat{y}_{\text{final}}^{\mathcal{T}}(u, i)) - \sigma(\hat{y}_{\text{final}}^{\mathcal{T}}(u, j))$ represent the preference margins from the student and fused teacher, respectively.

Gate Alignment Loss. To prevent the gate coefficient α_{ui} from collapsing to trivial solutions, we introduce an explicit supervisory signal. Let $t_i = \mathbb{I}(i \in \mathcal{I}_{\text{tail}})$ be an indicator for tail items. We employ binary cross-entropy to guide the gate:

$$\mathcal{L}_{\text{Align}} = - \sum_{(u,i) \in \mathcal{D}} [t_i \ln(\alpha_{ui}) + (1 - t_i) \ln(1 - \alpha_{ui})]. \quad (20)$$

This objective encourages the gate to prioritize expert knowledge ($\alpha_{ui} \rightarrow 1$) for tail items while retaining global knowledge ($\alpha_{ui} \rightarrow 0$) for popular items.

4.5 Complexity Analysis

Training Complexity. Given batch size B , embedding dimension d , and K experts, the primary computational costs include: (1) student forward pass: $O(Bd)$; (2) geodesic distance computation: $O(BKd)$; and (3) Meta-Debiasing Gate: $O(Bd^2)$. Since K is typically small (e.g., $K \leq 5$), and the winner-takes-all strategy activates only the selected expert and global teacher, the overall training complexity remains linear in batch size.

Inference Complexity. During inference, all auxiliary components (expert teachers, spherical projection, Meta-Debiasing Gate) are discarded. Only the lightweight student model is deployed, yielding inference complexity identical to the base architecture with no additional latency.

5 Experiments

We conduct comprehensive experiments to evaluate GUIDE, aiming to answer the following research questions:

- **RQ1:** How does GUIDE perform compared to state-of-the-art knowledge distillation methods for recommendation?
- **RQ2:** What are the individual contributions of each proposed component?
- **RQ3:** How sensitive is GUIDE to key hyperparameters?
- **RQ4:** How effective is GUIDE in handling popularity bias and enhancing item fairness?

5.1 Experimental Setup

Datasets. We evaluate our model on three public benchmarks: **CiteULike** (scholarly articles), **Amazon-Movie**, and **Amazon-Game** (e-commerce). Following standard preprocessing, we retain users with at least 20 interactions and split

Table 1: Dataset statistics after preprocessing.

Statistic	CiteULike	Amazon-Movie	Amazon-Game
#Users	5,219	123,960	24,303
#Items	25,181	50,052	10,672
#Interactions	125,580	1,697,533	231,780
Sparsity	99.90%	99.97%	99.91%

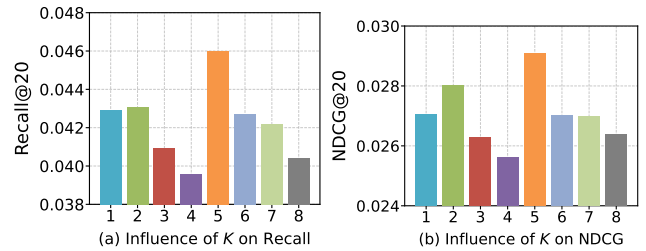


Figure 3: Effect of the number of experts K on Amazon-Game.

395 the data into training, validation, and test sets using an 8:1:1
396 ratio. Detailed statistics are reported in Table 1.

397 **Baselines.** We compare GUIDE against five state-of-the-
398 art distillation methods: (1) *Ranking-based*: **RD** [Tang
399 and Wang, 2018] and **CD** [Lee *et al.*, 2019] transfer rank-
400 ing knowledge via position-aware weighting and rank-aware
401 sampling; (2) *Structure-aware*: **DE-RRD** [Kang *et al.*, 2020]
402 and **HTD** [Kang *et al.*, 2021] distill latent knowledge and
403 topological structures; and (3) *Bias-aware*: **UNKD** [Chen *et*
404 *al.*, 2023a] mitigates popularity bias via stratified distillation.
405 **Evaluation Metrics.** We adopt two widely-used top- K
406 metrics: **Recall@ K** ($R@K$), which measures the proportion of
407 relevant items retrieved, and **NDCG@ K** ($N@K$), which ac-
408 counts for ranking position via discounted cumulative gain.
409 We report results for $K \in \{10, 20\}$.

410 **Implementation Details.** All hyperparameters are tuned
411 based on NDCG@20 on the validation set. The teacher
412 uses 3-layer LightGCN, for students, we evaluate 2-layer
413 GCN and BPRMF. The learning rate is searched over
414 $\{0.0002, 0.001, 0.005, 0.01\}$. The distillation coefficient λ_{KD}
415 is selected from $\{0.01, 0.02, \dots, 0.9, 0.99\}$, and the alignment
416 coefficient λ_{Align} from $\{0.0005, 0.001, \dots, 0.9, 0.99\}$. The
417 number of experts K is set to 5 by default. All experiments
418 are repeated 5 times with different random seeds, and we re-
419 port the mean and standard deviation.

420 5.2 Overall Performance (RQ1)

421 Table 2 presents the overall performance comparison. From
422 these results, we draw the following observations:

423 **GUIDE consistently outperforms all baselines across**
424 **datasets and metrics.** The improvements are particularly
425 pronounced on the sparser Amazon-Game dataset, demon-
426 strating GUIDE’s robustness in challenging scenarios where
427 popularity bias is most detrimental. Remarkably, GUIDE sur-
428 passes the teacher on several metrics, indicating that filtering
429 popularity noise yields superior generalization.

430 **Ranking-based methods exhibit limited effectiveness.** RD
431 yields the poorest performance among all methods, as simple
432 ranking alignment neglects the rich semantic relationships es-
433 sential for recommendation. This underscores the necessity
434 of more sophisticated knowledge transfer mechanisms.

435 **Structure-aware methods provide substantial improve-**
436 **ments over ranking-based approaches.** CD, HTD, and
437 UNKD significantly outperform RD by incorporating struc-
438 tural knowledge. Notably, CD achieves competitive results by
439 distilling collaborative similarity distributions, while UNKD
440 leverages unobserved feedback for enhanced supervision.

441 **The advantages of GUIDE stem from its two core inno-**
442 **vements.** Spherical Expert Alignment eliminates magnitude-
443 induced bias in expert routing, enabling semantically mean-
444 ingful knowledge selection. The Meta-Debiasing Gate fur-

ther adapts the fusion strategy to each context, ensuring debi- 445
asing preserves accuracy on well-represented items. 446

447 5.3 Ablation Study (RQ2)

To validate the contribution of each component, we compare 448
GUIDE against four ablated variants on Amazon-Game: 449

- 450 • **w/o Spherical:** Replaces spherical manifold matching with
451 Euclidean distance-based expert selection.
- 452 • **w/o Meta-Gate:** Substitutes the adaptive gate with a fixed
453 fusion weight ($\alpha = 0.5$).
- 454 • **w/o Ranking KD:** Replaces the ranking-aware distillation
455 loss with standard pointwise MSE.
- 456 • **w/o Alignment:** Removes the gate alignment loss \mathcal{L}_{Align} .

457 As shown in Table 3, the full GUIDE model consistently 458
outperforms all variants, confirming the necessity of each 459
component: (i) *w/o Spherical* shows that Euclidean embed- 460
dings are prone to magnitude-induced bias, emphasizing the 461
need for directional alignment; (ii) *w/o Meta-Gate* confirms 462
that static fusion is insufficient for handling user-item hetero- 463
geneity; (iii) *w/o Ranking KD* proves that pointwise distilla- 464
tion fails to capture nuanced preference margins; and (iv) *w/o* 465
Alignment indicates that explicit gate supervision is vital to 466
avoid trivial solutions.

467 5.4 Hyperparameter Sensitivity (RQ3)

468 We examine three key hyperparameters: (1) **Number of ex-** 468
erts K : As shown in Figure 3, performance peaks around 469
 $K = 5$; too few experts fail to capture popularity diversity, 470
while excessive K introduces noise due to data sparsity. (2) 471
Distillation coefficient λ_{KD} : Figure 4 (Left) shows optimal 472
result at 0.8, balancing knowledge transfer and student adap- 473
tation. (3) **Alignment coefficient λ_{Align} :** Figure 4 (Right) 474
indicates that $\lambda_{Align} = 0.01$ is optimal, while larger values 475
may conflict with the primary ranking objective. 476

477 5.5 Debiasing Effectiveness (RQ4)

478 To evaluate the effectiveness of GUIDE in addressing popu- 478
larity bias and enhancing fairness, we conduct three comple- 479
mentary analyses. 480

481 **Addressing Popularity Bias and Fairness.** We evaluate 481
GUIDE on Amazon-Game (stratified into Head/Tail groups) 482
and CiteULike (fairness metrics). Table 4 presents the re- 483
sults. GUIDE achieves substantial improvements on Tail 484
items (+6.21%) while simultaneously improving Head per- 485
formance (+2.75%). On CiteULike, GUIDE significantly 486

Table 2: Performance comparison on three datasets (Student: LightGCN). Best student results are in **bold**, second-best are underlined. “Improv.” denotes the relative improvement of GUIDE over the best baseline. Blue shading indicates teacher results.

Dataset	Metric	Teacher	RD	CD	DE-RRD	HTD	UNKD	GUIDE	Improv.
CiteULike	R@10	0.0964	0.062±.006	0.086±.005	0.062±.001	0.089±.002	0.093±.001	0.097±.001	+4.3%
	R@20	0.1456	0.085±.009	0.129±.006	0.092±.003	0.140±.001	0.132±.004	0.145±.002	+3.6%
	N@10	0.0788	0.051±.005	0.075±.002	0.050±.002	0.076±.001	0.078±.000	0.079±.000	+1.3%
	N@20	0.0937	0.058±.006	0.087±.003	0.060±.001	0.091±.001	0.089±.002	0.094±.001	+3.3%
Amazon-Movie	R@10	0.0284	0.020±.001	0.025±.003	0.022±.001	0.027±.002	0.027±.001	0.029±.001	+7.4%
	R@20	0.0487	0.033±.003	0.044±.004	0.037±.001	0.045±.002	0.046±.001	0.049±.000	+6.5%
	N@10	0.0313	0.023±.002	0.028±.002	0.024±.001	0.030±.001	0.030±.001	0.031±.002	+3.3%
	N@20	0.0394	0.028±.005	0.036±.003	0.030±.000	0.037±.001	0.037±.002	0.039±.001	+5.4%
Amazon-Game	R@10	0.0282	0.018±.001	0.025±.002	0.020±.002	0.026±.001	0.025±.001	0.029±.001	+11.5%
	R@20	0.0460	0.028±.003	0.043±.003	0.028±.001	0.040±.002	0.041±.002	0.046±.002	+6.9%
	N@10	0.0232	0.016±.001	0.020±.002	0.017±.000	0.021±.001	0.021±.001	0.023±.000	+9.5%
	N@20	0.0299	0.020±.001	0.027±.002	0.022±.000	0.026±.001	0.027±.001	0.029±.002	+7.4%

Table 3: Ablation study on Amazon-Game. Each row removes one component from the full GUIDE model.

Variant	R@10	R@20	N@10	N@20
w/o Spherical	0.0259	0.0422	0.0206	0.0266
w/o Meta-Gate	0.0259	0.0419	0.0209	0.0270
w/o Ranking KD	0.0258	0.0420	0.0206	0.0267
w/o Alignment	0.0261	0.0419	0.0210	0.0270
GUIDE (Full)	0.0288	0.0465	0.0230	0.0299

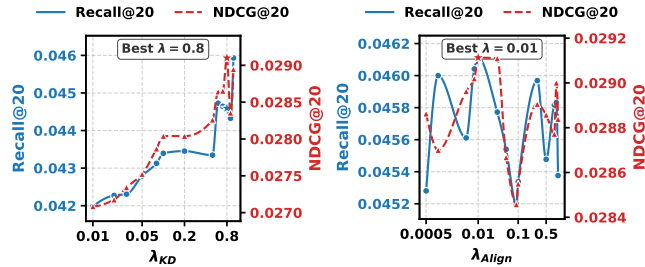


Figure 4: Sensitivity to λ_{KD} and λ_{Align} on Amazon-Game.

Table 4: Performance on Amazon-Game (Head/Tail) and fairness metrics on CiteULike.

Model	Amazon-Game			CiteULike			
	Head	Tail	Δ Tail	Gini ↓	Coverage ↑	Tail Ratio ↑	R@20
UNKD	0.0509	0.00177	–	0.9679	0.1234	0.0845	0.1331
GUIDE	0.0523	0.00188	+6.21%	0.9239	0.2521	0.0876	0.1447

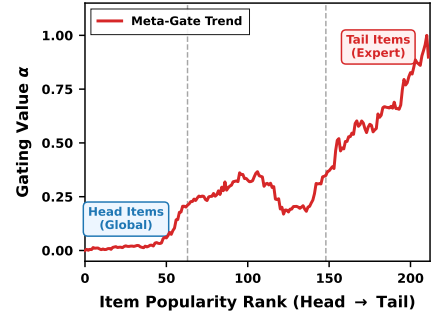


Figure 5: Distribution of gating coefficient α across item popularity. The gate adaptively shifts from global teacher ($\alpha \rightarrow 0$) for popular items to expert knowledge ($\alpha \rightarrow 1$) for tail items.

lowers the Gini index and doubles Item Coverage. These results across datasets prove that GUIDE effectively handles popularity bias and improves fairness while maintaining superior overall accuracy.

Visualization of Adaptive Gating. To verify that the Meta-Debiasing Gate learns meaningful context-aware strategies, we visualize the distribution of gating coefficients α_{ui} across item popularity on the Amazon-Game test set. As shown in Figure 5, a clear trend emerges: $\alpha \approx 0$ for popular (Head) items, indicating reliance on the global teacher, while $\alpha \approx 1$ for long-tail (Tail) items, shifting toward expert knowledge. This confirms that the gate learns a principled popularity-aware arbitration strategy rather than random assignment, validating the effectiveness of the alignment loss \mathcal{L}_{Align} .

Case Study: Cold-Start Item Recommendation. To evaluate performance under extreme sparsity, we analyze a cold-start interaction (User #8137, Item #10872) in Amazon-

Game. While the LightGCN teacher fails to recommend the item (Rank > 200) due to absent collaborative signals, GUIDE achieves Rank 1 with a high gating coefficient ($\alpha_{ui} = 0.84$). This confirms that GUIDE adaptively prioritizes specialized expert knowledge to compensate for data scarcity, effectively enabling zero-shot cold-start recommendations.

6 Conclusion

Existing KD approaches, using standard Euclidean teachers, often suffer from bias amplification. We propose GUIDE, projecting diverse expert knowledge onto a Spherical Manifold. This design eliminates magnitude-induced noise via distinct subspace experts, while our Meta-Debiasing Gate intelligently arbitrates between global and specialized supervision to adapt to dynamic user preferences. Experiments confirm that GUIDE consistently outperforms competing methods in both recommendation accuracy and bias mitigation.

520 Acknowledgments

521 This work was supported by the National Natural Science
522 Foundation of China (No.62173199 and No.62506366).

523 References

- 524 [Abdollahpouri *et al.*, 2019] Himan Abdollahpouri, Robin
525 Burke, and Bamshad Mobasher. Managing popularity bias
526 in recommender systems with personalized re-ranking. In
527 Roman Barták and Keith W. Brawner, editors, *Proceedings*
528 *of the Thirty-Second International Florida Artificial Intel-*
529 *ligence Research Society Conference, Sarasota, Florida,*
530 *USA, May 19-22 2019*, pages 413–418. AAAI Press, 2019.
- 531 [Abdollahpouri *et al.*, 2021] Himan Abdollahpouri, Masoud
532 Mansoury, Robin Burke, Bamshad Mobasher, and Ed-
533 ward C. Malthouse. User-centered evaluation of popu-
534 larity bias in recommender systems. In Judith Masthoff,
535 Eelco Herder, Nava Tintarev, and Marko Tkalcić, edi-
536 tors, *Proceedings of the 29th ACM Conference on User*
537 *Modeling, Adaptation and Personalization, UMAP 2021,*
538 *Utrecht, The Netherlands, June, 21-25, 2021*, pages 119–
539 129. ACM, 2021.
- 540 [Chen *et al.*, 2021] Jiawei Chen, Hande Dong, Yang Qiu, Xi-
541 angnan He, Xin Xin, Liang Chen, Guli Lin, and Keping
542 Yang. Autodebias: Learning to debias for recommenda-
543 tion. In *SIGIR*, pages 21–30. ACM, 2021.
- 544 [Chen *et al.*, 2023a] Gang Chen, Jiawei Chen, Fuli Feng,
545 Sheng Zhou, and Xiangnan He. Unbiased knowledge dis-
546 tillation for recommendation. In *WSDM*, pages 976–984.
547 ACM, 2023.
- 548 [Chen *et al.*, 2023b] Jiawei Chen, Hande Dong, Xiang
549 Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias
550 and debias in recommender system: A survey and future
551 directions. *ACM Trans. Inf. Syst.*, 41(3):67:1–67:39, 2023.
- 552 [Deldjoo *et al.*, 2024] Yashar Deldjoo, Zhankui He, Ju-
553 lian J. McAuley, Anton Korikov, Scott Sanner, Arnau
554 Ramisa, René Vidal, Maheswaran Sathiamoorthy, Atoosa
555 Kasirzadeh, and Silvia Milano. A review of modern rec-
556ommender systems using generative models (gen-recsys).
557 In Ricardo Baeza-Yates and Francesco Bonchi, editors,
558 *Proceedings of the 30th ACM SIGKDD Conference on*
559 *Knowledge Discovery and Data Mining, KDD 2024,*
560 *Barcelona, Spain, August 25-29, 2024*, pages 6448–6458.
561 ACM, 2024.
- 562 [Ding *et al.*, 2022] Sihao Ding, Fuli Feng, Xiangnan He, Jin-
563 qiu Jin, Wenjie Wang, Yong Liao, and Yongdong Zhang.
564 Interpolative distillation for unifying biased and debiased
565 recommendation. In Enrique Amigó, Pablo Castells,
566 Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and
567 Gabriella Kazai, editors, *SIGIR '22: The 45th Interna-*
568 *tional ACM SIGIR Conference on Research and Develop-*
569 *ment in Information Retrieval, Madrid, Spain, July 11 -*
570 *15, 2022*, pages 40–49. ACM, 2022.
- 571 [Gao *et al.*, 2023] Chen Gao, Yu Zheng, Nian Li, Yinfeng Li,
572 Yingrong Qin, Jinghua Piao, Yuhan Quan, Jianxin Chang,
573 Depeng Jin, Xiangnan He, and Yong Li. A survey of graph
574 neural networks for recommender systems: Challenges,
575 methods, and directions. *Trans. Recomm. Syst.*, 1(1):1–51,
576 2023.
- [Gou *et al.*, 2024] Jianping Gou, Liyuan Sun, Baosheng Yu,
577 Lan Du, Kotagiri Ramamohanarao, and Dacheng Tao.
578 Collaborative knowledge distillation via multiknowledge
579 transfer. *IEEE Trans. Neural Networks Learn. Syst.*,
580 35(5):6718–6730, 2024.
581
- [Han *et al.*, 2024] Yishan Han, Biao Xu, Yao Wang, and
582 Shanxing Gao. Towards popularity-aware recommenda-
583 tion: A multi-behavior enhanced framework with orthog-
584 onality constraint. *CoRR*, abs/2412.19172, 2024.
585
- [He *et al.*, 2022] Ming He, Xinlei Hu, Changshu Li, Xin
586 Chen, and Jiwen Wang. Mitigating confounding bias
587 for recommendation via counterfactual inference. In
588 *ECML/PKDD (1)*, volume 13713 of *Lecture Notes in Com-*
589 *puter Science*, pages 524–540. Springer, 2022.
590
- [Joachims *et al.*, 2018] Thorsten Joachims, Adith Swami-
591 nathan, and Tobias Schnabel. Unbiased learning-to-rank
592 with biased feedback. In Jérôme Lang, editor, *Proceed-*
593 *ings of the Twenty-Seventh International Joint Conference*
594 *on Artificial Intelligence, IJCAI 2018, July 13-19, 2018,*
595 *Stockholm, Sweden*, pages 5284–5288. ijcai.org, 2018.
596
- [Kang *et al.*, 2020] SeongKu Kang, Junyoung Hwang, Won-
597 bin Kweon, and Hwanjo Yu. DE-RRD: A knowledge dis-
598 tillation framework for recommender system. In *CIKM*,
599 pages 605–614. ACM, 2020.
600
- [Kang *et al.*, 2021] SeongKu Kang, Junyoung Hwang, Won-
601 bin Kweon, and Hwanjo Yu. Topology distillation for rec-
602ommender system. In *KDD*, pages 829–839. ACM, 2021.
603
- [Koren *et al.*, 2009] Yehuda Koren, Robert M. Bell, and
604 Chris Volinsky. Matrix factorization techniques for rec-
605ommender systems. *Computer*, 42(8):30–37, 2009.
606
- [Kweon *et al.*, 2021] Wonbin Kweon, SeongKu Kang, and
607 Hwanjo Yu. Bidirectional distillation for top-k recom-
608mender system. In Jure Leskovec, Marko Grobelnik, Marc
609 Najork, Jie Tang, and Leila Zia, editors, *WWW '21: The*
610 *Web Conference 2021, Virtual Event / Ljubljana, Slove-*
611 *nia, April 19-23, 2021*, pages 3861–3871. ACM / IW3C2,
612 2021.
613
- [Lee *et al.*, 2019] Jae-woong Lee, Minjin Choi, Jongwuk
614 Lee, and Hyunjung Shim. Collaborative distillation for
615 top-n recommendation. In *ICDM*, pages 369–378. IEEE,
616 2019.
617
- [Li *et al.*, 2023] Haoxuan Li, Yanghao Xiao, Chunyuan
618 Zheng, and Peng Wu. Balancing unobserved confounding
619 with a few unbiased ratings in debiased recommendations.
620 In *WWW*, pages 1305–1313. ACM, 2023.
621
- [Lin *et al.*, 2019] Kun Lin, Nasim Sonboli, Bamshad
622 Mobasher, and Robin Burke. Crank up the volume:
623 Preference bias amplification in collaborative recom-
624mendation. In *RMSE@RecSys*, volume 2440 of *CEUR*
625 *Workshop Proceedings*. CEUR-WS.org, 2019.
626
- [Liu *et al.*, 2020] Dugang Liu, Pengxiang Cheng, Zhenhua
627 Dong, Xiuqiang He, Weike Pan, and Zhong Ming. A gen-
628 eral knowledge distillation framework for counterfactual
629

- 630 recommendation via uniform data. In *SIGIR*, pages 831–
631 840. ACM, 2020.
- 632 [Liu *et al.*, 2021] Dugang Liu, Pengxiang Cheng, Hong Zhu,
633 Zhenhua Dong, Xiuqiang He, Weike Pan, and Zhong
634 Ming. Mitigating confounding bias in recommendation via
635 information bottleneck. In *RecSys*, pages 351–360. ACM,
636 2021.
- 637 [Liu *et al.*, 2023] Dugang Liu, Pengxiang Cheng, Hong Zhu,
638 Zhenhua Dong, Xiuqiang He, Weike Pan, and Zhong
639 Ming. Debaised representation learning in recommenda-
640 tion via information bottleneck. *Trans. Recomm. Syst.*,
641 1(1):1–27, 2023.
- 642 [Ning *et al.*, 2024] Wentao Ning, Reynold Cheng, Xiao Yan,
643 Ben Kao, Nan Huo, Nur Al Hasan Haldar, and Bo Tang.
644 Debiasing recommendation with personal popularity. In
645 *WWW*, pages 3400–3409. ACM, 2024.
- 646 [Oestreicher-Singer and Sundararajan, 2012] Gal
647 Oestreicher-Singer and Arun Sundararajan. Recom-
648 mendation networks and the long tail of electronic
649 commerce. *MIS Q.*, 36(1):65–83, 2012.
- 650 [Park *et al.*, 2026] Jongwon Park, Minku Kang, Wooseok
651 Sim, Soyung Lee, and Hogun Park. Federated recom-
652 mender system with data valuation for e-commerce plat-
653 form. *Expert Syst. Appl.*, 298:129695, 2026.
- 654 [Romero *et al.*, 2015] Adriana Romero, Nicolas Ballas,
655 Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta,
656 and Yoshua Bengio. Fitnets: Hints for thin deep nets. In
657 *ICLR (Poster)*, 2015.
- 658 [Sarwar *et al.*, 2001] Badrul Munir Sarwar, George Karypis,
659 Joseph A. Konstan, and John Riedl. Item-based collabora-
660 tive filtering recommendation algorithms. In Vincent Y.
661 Shen, Nobuo Saito, Michael R. Lyu, and Mary Ellen
662 Zurko, editors, *Proceedings of the Tenth International
663 World Wide Web Conference, WWW 10, Hong Kong,
664 China, May 1-5, 2001*, pages 285–295. ACM, 2001.
- 665 [Schnabel *et al.*, 2016] Tobias Schnabel, Adith Swami-
666 nathan, Ashudeep Singh, Navin Chandak, and Thorsten
667 Joachims. Recommendations as treatments: Debias-
668 ing learning and evaluation. In *ICML*, volume 48 of
669 *JMLR Workshop and Conference Proceedings*, pages
670 1670–1679. JMLR.org, 2016.
- 671 [Song *et al.*, 2023] Zijie Song, Jiawei Chen, Sheng Zhou,
672 Qihao Shi, Yan Feng, Chun Chen, and Can Wang. CDR:
673 conservative doubly robust learning for debaised recom-
674 mendation. In *CIKM*, pages 2321–2330. ACM, 2023.
- 675 [Tang and Wang, 2018] Jiayi Tang and Ke Wang. Rank-
676 ing distillation: Learning compact ranking models with
677 high performance for recommender system. In Yike Guo
678 and Faisal Farooq, editors, *Proceedings of the 24th ACM
679 SIGKDD International Conference on Knowledge Discov-
680 ery & Data Mining, KDD 2018, London, UK, August 19-
681 23, 2018*, pages 2289–2298. ACM, 2018.
- 682 [Tang *et al.*, 2020] Kaihua Tang, Jianqiang Huang, and Han-
683 wang Zhang. Long-tailed classification by keeping the
684 good and removing the bad momentum causal effect. In
685 *NeurIPS*, 2020.
- [Tao *et al.*, 2022] Ye Tao, Ying Li, Su Zhang, Zhirong Hou, 686
and Zhonghai Wu. Revisiting graph based social recom- 687
mendation: A distillation enhanced social graph network. 688
In *WWW*, pages 2830–2838. ACM, 2022. 689
- [Wang *et al.*, 2017] Feng Wang, Xiang Xiang, Jian Cheng, 690
and Alan Loddon Yuille. Normface: L_2 hypersphere em- 691
bedding for face verification. In Qiong Liu, Rainer Lien- 692
hart, Haohong Wang, Sheng-Wei ”Kuan-Ta” Chen, Su- 693
sanne Boll, Yi-Ping Phoebe Chen, Gerald Friedland, Jia 694
Li, and Shuicheng Yan, editors, *Proceedings of the 2017
ACM on Multimedia Conference, MM 2017, Mountain
View, CA, USA, October 23-27, 2017*, pages 1041–1049. 695
ACM, 2017. 696
697
698
- [Wang *et al.*, 2021] Shuai Wang, Kun Zhang, Le Wu, Haip- 699
ing Ma, Richang Hong, and Meng Wang. Privileged 700
graph distillation for cold start recommendation. In Fer- 701
nando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, 702
Rosie Jones, and Tetsuya Sakai, editors, *SIGIR ’21: The
44th International ACM SIGIR Conference on Research
and Development in Information Retrieval, Virtual Event,
Canada, July 11-15, 2021*, pages 1187–1196. ACM, 2021. 703
704
705
706
- [Xu *et al.*, 2020] Chen Xu, Quan Li, Junfeng Ge, Jinyang 707
Gao, Xiaoyong Yang, Changhua Pei, Fei Sun, Jian Wu, 708
Hanxiao Sun, and Wenwu Ou. Privileged features distil- 709
lation at taobao recommendations. In *KDD*, pages 2590– 710
2598. ACM, 2020. 711
- [Yang *et al.*, 2025] Xinxin Yang, Xinwei Li, Zhen Liu, Yafan 712
Yuan, and Yannan Wang. Multi-teacher knowledge distil- 713
lation for debiasing recommendation with uniform data. 714
Expert Syst. Appl., 273:126808, 2025. 715
- [Zhang and Shen, 2023] Fan Zhang and Qijie Shen. A 716
model-agnostic popularity debias training framework for 717
click-through rate prediction in recommender system. In 718
SIGIR, pages 1760–1764. ACM, 2023. 719
- [Zhang *et al.*, 2019] Shuai Zhang, Lina Yao, Aixin Sun, and 720
Yi Tay. Deep learning based recommender system: A sur- 721
vey and new perspectives. *ACM Comput. Surv.*, 52(1):5:1– 722
5:38, 2019. 723
- [Zhang *et al.*, 2020] Yuan Zhang, Xiaoran Xu, Hanning 724
Zhou, and Yan Zhang. Distilling structured knowledge 725
into embeddings for explainable and accurate recommen- 726
dation. In *WSDM*, pages 735–743. ACM, 2020. 727
- [Zhu *et al.*, 2020] Jieming Zhu, Jinyang Liu, Weiqi Li, Jincan 728
Lai, Xiuqiang He, Liang Chen, and Zibin Zheng. Ensem- 729
bled CTR prediction via knowledge distillation. In *CIKM*, 730
pages 2941–2958. ACM, 2020. 731
- [Zhu *et al.*, 2021] Ziwei Zhu, Yun He, Xing Zhao, Yin 732
Zhang, Jianling Wang, and James Caverlee. Popularity- 733
opportunity bias in collaborative filtering. In *WSDM ’21,
The Fourteenth ACM International Conference on Web
Search and Data Mining, Virtual Event, Israel, March 8-
12, 2021*, pages 85–93. ACM, 2021. 734
735
736
737