

SymCERE: Symmetric Contrastive Learning for Robust Review-Enhanced Recommendation

Toyotaro Suzumura

suzumura@acm.org

The University of Tokyo / Rakuten
Group, Inc.
Tokyo, Japan

Hisashi Ikari

hisashi@ikari.io

Independent Researcher
Tokyo, Japan

Hiroki Kanezashi

hkanezashi@acm.org

The University of Tokyo / Rakuten
Group, Inc.
Tokyo, Japan

Md Mostafizur Rahman

mdmostafizu.a.rahman@rakuten.com

Rakuten Group, Inc.

Tokyo, Japan

Yu Hirate

yu.hirate@rakuten.com

Rakuten Group, Inc.

Tokyo, Japan

ABSTRACT

Modern recommendation systems can achieve high performance by fusing user behavior graphs (via GNNs) and review texts (via LLMs). However, this fusion faces three significant issues: (1) **False Negatives** in contrastive learning can degrade the training signal by penalizing similar items; (2) **Popularity Bias**, often encoded as embedding magnitude, can distort similarity scores; and (3) **Signal Ambiguity**, which arises from the conflation of objective facts with subjective sentiment in reviews. These interconnected issues can prevent models from learning users' true preferences.

In this paper, we propose **SymCERE (Symmetric SINCERE)**, a contrastive learning method that addresses these three issues simultaneously through its structural design. First, we introduce a symmetric application of the SINCERE loss for cross-modal alignment, which is designed to eliminate false negatives in recommendation. Second, by integrating this with L2 normalisation under a "magnitude-as-noise" hypothesis, we aim to mitigate popularity bias by forcing the model to encode preferences primarily in the vector's direction.

Experiments on 15 datasets from three distinct platforms (e-commerce, local reviews, and travel) demonstrate that SymCERE outperforms several strong baselines, achieving a relative improvement of up to **43.6%** on NDCG@10. Furthermore, a detailed LIME analysis shows that the model learns to anchor alignment on objective, informative vocabulary (e.g., "OEM," "compatible," "gasket"), while placing less emphasis on generic sentiment (e.g., "good," "great"). This suggests that effective semantic alignment stems from understanding factual product attributes, offering a path toward more accurate recommendation systems. The code is available at: <https://anonymous.4open.science/r/ReviewGNN-2E1E>.

CCS CONCEPTS

- Information systems → Recommender systems;
- Computing methodologies → Machine learning.

KEYWORDS

Recommender Systems, Contrastive Learning, False Negatives, Popularity Bias, GNN, LLM

ACM Reference Format:

Toyotaro Suzumura, Hisashi Ikari, Hiroki Kanezashi, Md Mostafizur Rahman, and Yu Hirate. 2018. SymCERE: Symmetric Contrastive Learning for Robust Review-Enhanced Recommendation. In . ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Recent advances in recommender systems focus on multi-modal approaches that fuse user behavior, typically captured by Graph Neural Networks (GNNs), with side information from modalities like review texts, processed by Large Language Models (LLMs) [26, 36]. While this fusion can yield richer user and item representations [2, 23], its practical implementation faces a triad of interconnected challenges that can impair recommendation quality. These are: the **false negative problem** in contrastive learning, item **popularity bias**, and **signal ambiguity** in review texts. These issues often compound one another: popularity bias can exacerbate the false negative problem, while both can be masked by ambiguous textual signals, making it difficult to learn a user's true intent.

While various methods have been proposed to tackle these issues individually—using denoising modules [17, 19] or causal inference [35, 38]—we posit that these challenges can be systematically mitigated by modifying the core contrastive learning objective itself. We propose **SymCERE (Symmetric SINCERE)**, a unified methodology that addresses this triad with geometric mechanisms. SymCERE employs a symmetric, denoised contrastive loss based on SINCERE [3, 39] to eliminate intra-class repulsion and false negatives. Guided by our "magnitude-as-noise" hypothesis, it integrates this with L2 normalisation to mitigate popularity bias by forcing the model to encode preferences solely in the vector's direction.

Our experiments on 15 datasets from three distinct platforms, multiple domains, and two languages show that SymCERE significantly outperforms strong baselines, achieving up to a **43.6%**

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

relative improvement in NDCG@10. Furthermore, our analysis reveals that the model learns to anchor its alignment on objective, factual vocabulary (e.g., ‘OEM’, ‘compatible’) over generic sentiment. This mechanism, which we term ‘**semantic anchoring**’, suggests that robust alignment stems from understanding factual product attributes, presenting a path toward more accurate and interpretable recommender systems. Our main contributions are:

- **A Unified Contrastive Method (SymCERE):** We propose a method that structurally mitigates false negatives, popularity bias, and signal ambiguity by symmetrically applying a denoised SINCERE loss and integrating geometric debiasing via L2 normalisation.
- **Extensive Empirical Validation:** We demonstrate the effectiveness of SymCERE through comprehensive experiments on 15 diverse datasets, achieving significant improvements over strong baselines.
- **A Novel Interpretation of Multi-Modal Alignment:** We find, through in-depth LIME analysis, that our model learns to prioritize product-specific terminology over generic sentiment, a mechanism we term ‘semantic anchoring’.

2 RELATED WORK

Our work lies at the intersection of multi-modal recommendation, Graph Neural Networks (GNNs), and contrastive learning (CL). While GNNs like LightGCN [6] form the backbone of many modern recommenders, enhancing them with review texts via CL surfaces a triad of interconnected challenges.

2.1 A Triad of Interconnected Challenges

The False Negative Problem. The standard InfoNCE loss has a key limitation: by treating all non-positive samples as negatives, it may penalize “false negatives”—semantically similar items that are not the anchor’s direct positive pair [1, 10]. This “intra-class repulsion” can negatively affect the embedding space by pushing similar items apart and hindering the development of a coherent semantic structure. This is especially problematic in recommendation, where a user who bought one camera might also be interested in a compatible lens—an item that standard InfoNCE would treat as a hard negative.

Popularity Bias as a Geometric Artifact. Popularity bias, the over-recommendation of popular items, is a known challenge to fairness and diversity [14]. We approach this from a geometric perspective, positing that popularity is often encoded as the magnitude (L2-norm) of an item’s embedding vector [18, 20]. This magnitude can act as a noisy, confounding signal that obscures user preferences for specific attributes. For instance, a model may incorrectly learn to equate a large embedding magnitude with high item quality, when it merely reflects high interaction frequency. Therefore, we treat L2 normalisation not just as a technical step, but as a debiasing mechanism that can structurally remove the popularity signal, building on the theoretical properties of CL on a unit hypersphere [30].

Signal Ambiguity in Multi-Modal Fusion. Fusing collaborative signals with review texts can enrich representations but may introduce “signal ambiguity.” A review can conflate objective product attributes (e.g., “OEM part,” “compatible gasket”) with subjective,

often uninformative sentiment (e.g., “good,” “great”) [32]. A simple alignment may lead a model to learn spurious correlations from generic sentiment rather than from true semantic understanding, as it cannot distinguish the nuanced reasons for a user’s choice.

2.2 Our Approach: A Unified Solution via Principled CL

The literature often presents methods that address these issues individually [19, 21, 22]. We diverge from this trend, arguing that modifying the core CL objective itself offers a more direct and unified solution. To this end, we build upon **SINCERE** [3, 39], a loss that eliminates the intra-class repulsion found in other contrastive methods. By integrating a symmetric SINCERE loss with L2 normalisation, SymCERE is designed to simultaneously: (1) **Eliminate False Negatives** by definition, (2) **Mitigate Popularity Bias** structurally, and (3) **Resolve Signal Ambiguity** by focusing alignment on stable, objective information. By embedding these solutions directly into the learning objective, our approach avoids the need for complex, multi-stage pipelines or heuristic post-processing. We posit that this unified approach leads to a simpler, yet more powerful and interpretable recommender system.

3 METHODOLOGY

Our proposed model architecture is designed to create a unified representation space where embeddings from user-item interactions (graph modality) and review texts (text modality) are meaningfully aligned. This is achieved through a multi-component system, trained end-to-end, which consists of modality-specific encoders and a multi-task contrastive learning objective.

3.1 Modality Encoders

We employ two specialized encoders to process the distinct data types, capturing both collaborative filtering patterns and deep semantic information.

3.1.1 Graph Encoder for Collaborative Signal. To capture the collaborative filtering signals from the user-item interaction graph, we adopt the architecture of LightGCN [6] as our graph encoder, denoted $f_{\text{graph}}(\cdot)$. This choice is motivated by its effectiveness in distilling the core mechanism of GNNs for recommendation—neighborhood aggregation—while removing components like feature transformations and non-linearities, which have been shown to be less critical for this task [6].

The encoder begins with an initial embedding layer, which provides zero-th layer embeddings $\mathbf{e}_u^{(0)}$ for each user u and $\mathbf{e}_i^{(0)}$ for each item i . These are the only trainable parameters in the graph encoder. The embeddings are then refined over K layers using a linear propagation rule that aggregates information from neighbors [13]:

$$\mathbf{e}_u^{(k+1)} = \sum_{i \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u||\mathcal{N}_i|}} \mathbf{e}_i^{(k)}, \quad \mathbf{e}_i^{(k+1)} = \sum_{u \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i||\mathcal{N}_u|}} \mathbf{e}_u^{(k)} \quad (1)$$

where \mathcal{N}_u is the set of items user u has interacted with, and \mathcal{N}_i is the set of users who have interacted with item i . The term $1/\sqrt{|\mathcal{N}_u||\mathcal{N}_i|}$ is a symmetrically normalized coefficient that stabilizes the propagation process [13].

The final representation for a user or item is a weighted combination of the embeddings from all layers, which captures collaborative signals from different neighborhood depths:

$$\mathbf{g}_u = \sum_{k=0}^K \alpha_k \mathbf{e}_u^{(k)}, \quad \mathbf{g}_i = \sum_{k=0}^K \alpha_k \mathbf{e}_i^{(k)} \quad (2)$$

where α_k are hyperparameters weighting the importance of the k -th layer. Following common practice, we set $\alpha_k = 1/(K+1)$ to give equal weight to each layer [6]. To obtain a single representation for a user-item interaction (u, i) , we average their final embeddings: $\mathbf{g}_{u,i} = (\mathbf{g}_u + \mathbf{g}_i)/2$.

3.1.2 Text Encoder for Semantic Representation. To extract semantic information from user reviews, we employ a pre-trained Large Language Model (LLM) as our text encoder, $f_{\text{text}}(\cdot)$. We utilize Parameter-Efficient Fine-Tuning (PEFT) with LoRA [9] to adapt the LLM to the recommendation domain without the high computational cost of full fine-tuning. Given the raw review text t_i associated with an item i , the encoder produces a dense semantic embedding $\mathbf{t}_i = f_{\text{text}}(t_i)$.

3.2 Unified Contrastive Learning Method

Our methodology uses a multi-task contrastive method designed to fuse the collaborative and semantic signals. It is built upon two principles: geometric debiasing and a denoised alignment objective.

3.2.1 Geometric Debiasing via L2 Normalisation. A key step, applied before any loss calculation, is the L2 normalisation of all embeddings. This is central to our "magnitude-as-noise" hypothesis, which posits that the L2-norm of an embedding often encodes popularity bias rather than user preference [18, 20]. By projecting all embeddings onto a unit hypersphere, we aim to eliminate this magnitude component:

$$\hat{\mathbf{g}}_{u,i} = \frac{\mathbf{g}_{u,i}}{\|\mathbf{g}_{u,i}\|_2}, \quad \hat{\mathbf{t}}_i = \frac{\mathbf{t}_i}{\|\mathbf{t}_i\|_2} \quad (3)$$

This ensures that similarity is based on angular distance (cosine similarity). This is a theoretically-grounded debiasing mechanism that satisfies the geometric precondition for contrastive learning to optimize for the properties of **Alignment** and **Uniformity** [30].

3.2.2 Multi-Task Contrastive Objectives. Our method is trained with a combination of a cross-modal loss for alignment and an intra-modal loss for representation enhancement [33].

Cross-Modal Alignment with Symmetric SINCERE Loss ($\mathcal{L}_{\text{cross-modal}}$). To align the graph and text modalities, we address the "false negative" problem in standard contrastive losses. The InfoNCE loss [29] is defined for an anchor \mathbf{z}_i , a positive sample \mathbf{z}_i^+ , and a set of $N - 1$ negative samples $\{\mathbf{z}_j^-\}_{j=1}^{N-1}$ as:

$$\mathcal{L}_{\text{InfoNCE}} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_i^+)/\tau)}{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_i^+)/\tau) + \sum_{j=1}^{N-1} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j^-)/\tau)} \quad (4)$$

where $\text{sim}(\cdot, \cdot)$ is a similarity function (e.g., dot product) and τ is a temperature parameter. This formulation treats all samples that are not the designated positive as negatives, which can lead to false negatives.

The Supervised Contrastive (SupCon) loss [12] attempts to address this by leveraging class labels. For an anchor \mathbf{z}_i from class c , it treats all other samples from class c as positives. However, its formulation can lead to *intra-class repulsion* [3, 39]. The SupCon loss for an anchor \mathbf{z}_i and a positive \mathbf{z}_p from the same class is:

$$\mathcal{L}_{\text{SupCon}} = -\log \frac{\exp(\mathbf{z}_i^\top \mathbf{z}_p/\tau)}{\sum_{k \in P_i} \exp(\mathbf{z}_i^\top \mathbf{z}_k/\tau) + \sum_{j \in N_i} \exp(\mathbf{z}_i^\top \mathbf{z}_j/\tau)} \quad (5)$$

where P_i is the set of all positives for anchor i and N_i is the set of all negatives. The denominator includes other positive samples from the same class, which can create a repulsive gradient term [39].

To overcome this, we adopt the SINCERE loss [3, 39], which adheres to noise-contrastive estimation principles: only true negatives should be in the denominator [4]. For an anchor \mathbf{z}_i , a positive \mathbf{z}_p , and a set of true negatives N_i , the SINCERE loss is:

$$\mathcal{L}_{\text{SINCERE}} = -\log \frac{\exp(\mathbf{z}_i^\top \mathbf{z}_p/\tau)}{\exp(\mathbf{z}_i^\top \mathbf{z}_p/\tau) + \sum_{j \in N_i} \exp(\mathbf{z}_i^\top \mathbf{z}_j/\tau)} \quad (6)$$

This formulation is designed to eliminate intra-class repulsion [39]. We apply this in a symmetric, cross-modal fashion. In our setting, each sample in a mini-batch corresponds to a single user-item interaction (i.e., an edge in the user-item graph). The graph embedding $\hat{\mathbf{g}}_{u,i}$ is computed from both the user u and item i via the graph encoder, while the text embedding $\hat{\mathbf{t}}_i$ is derived from the review(s) of item i .

In our method, we process the data in mini-batches, where each sample in a mini-batch corresponds to a single user-item interaction (i.e., an edge in the user-item graph). For a batch of B interactions, the graph-to-text loss aligns each graph embedding $\hat{\mathbf{g}}_{u,i}$ with its corresponding text embedding $\hat{\mathbf{t}}_i$, treating all other text embeddings $\{\hat{\mathbf{t}}_j\}_{j \neq i}$ as true negatives (i.e., items in the batch that have different item IDs from the anchor, excluding those the user has previously interacted with):

$$\mathcal{L}_{G \rightarrow T} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\hat{\mathbf{g}}_{u,i}^\top \hat{\mathbf{t}}_i/\tau)}{\exp(\hat{\mathbf{g}}_{u,i}^\top \hat{\mathbf{t}}_i/\tau) + \sum_{j \in N_i} \exp(\hat{\mathbf{g}}_{u,i}^\top \hat{\mathbf{t}}_j/\tau)} \quad (7)$$

Symmetrically, the text-to-graph loss is:

$$\mathcal{L}_{T \rightarrow G} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\hat{\mathbf{t}}_i^\top \hat{\mathbf{g}}_{u,i}/\tau)}{\exp(\hat{\mathbf{t}}_i^\top \hat{\mathbf{g}}_{u,i}/\tau) + \sum_{j \in N_i} \exp(\hat{\mathbf{t}}_i^\top \hat{\mathbf{g}}_{u,j}/\tau)} \quad (8)$$

The final cross-modal loss is the average of these two, ensuring balanced alignment:

$$\mathcal{L}_{\text{cross-modal}} = \frac{1}{2} (\mathcal{L}_{G \rightarrow T} + \mathcal{L}_{T \rightarrow G}) \quad (9)$$

Intra-Modal Alignment with InfoNCE Loss ($\mathcal{L}_{\text{intra-modal}}$). To improve the robustness and discriminative power of the representations within each modality, we incorporate a self-supervised, intra-modal contrastive task [33]. For each modality, we create two distinct "views" through data augmentation and enforce consistency between their representations.

For the graph modality, we create an augmented view by applying random edge dropout. Let $\hat{\mathbf{g}}$ be the batch of embeddings from the original graph and $\hat{\mathbf{g}}'$ be from the augmented graph. For the text modality, we create an augmented view via random word masking.

Let $\hat{\mathbf{t}}$ be the embeddings of the original texts and $\hat{\mathbf{t}}'$ be those of the augmented texts.

We use the standard InfoNCE loss for this task, as within-modality augmentation does not have the same false negative issue as cross-modal alignment. The loss for a generic modality with original embeddings $\hat{\mathbf{z}}$ and augmented embeddings $\hat{\mathbf{z}}'$ is:

$$\mathcal{L}_{\text{InfoNCE}}(\hat{\mathbf{z}}, \hat{\mathbf{z}}') = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\hat{\mathbf{z}}_i^\top \hat{\mathbf{z}}'_i / \tau)}{\sum_{j=1}^B \exp(\hat{\mathbf{z}}_i^\top \hat{\mathbf{z}}'_j / \tau)} \quad (10)$$

The total intra-modal loss is the sum of the losses from both modalities:

$$\mathcal{L}_{\text{intra-modal}} = \mathcal{L}_{\text{InfoNCE}}(\hat{\mathbf{g}}, \hat{\mathbf{g}}') + \mathcal{L}_{\text{InfoNCE}}(\hat{\mathbf{t}}, \hat{\mathbf{t}}') \quad (11)$$

This formulation can resolve signal ambiguity because by removing potential noise from false negatives (via SINCERE) and popularity signals (via L2 norm), the optimization is encouraged to find stable, information-rich correlations. In many domains, these signals are objective, factual attributes (e.g., "OEM," "compatible") rather than subjective sentiments (e.g., "good," "love"), leading to a more grounded semantic alignment.

Intra-Modal Ranking Regularization. In addition to aligning modalities, we regularize the embedding space within the graph modality using the BPR loss. This complements our main objective by enforcing a well-ordered ranking structure based on collaborative signals. To maintain logical consistency with our "magnitude-as-noise" hypothesis, the BPR loss is critically applied to the L2-normalized embeddings ($\hat{\mathbf{g}}$):

$$\mathcal{L}_{\text{BPR}} = \sum_{(u,i,j)} -\log \sigma(\hat{\mathbf{g}}_u^\top \hat{\mathbf{g}}_i - \hat{\mathbf{g}}_u^\top \hat{\mathbf{g}}_j) \quad (12)$$

3.3 Final Objective Function

The final training objective is a weighted sum of the cross-modal loss, the intra-modal loss, and a complementary Bayesian Personalized Ranking (BPR) loss for intra-modal regularization. The total loss is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{cross-modal}} + \alpha \cdot \mathcal{L}_{\text{intra-modal}} + \beta \cdot \mathcal{L}_{\text{BPR}} \quad (13)$$

where α and β are hyperparameters that balance the contributions of the contrastive alignment tasks and the recommendation ranking task. In our experiments, we set $\alpha = 0.5$ and $\beta = 0.05$.

4 EXPERIMENTS

We conducted a comprehensive set of experiments to evaluate the performance of our proposed model, SymCERE. The primary research questions we aim to answer are:

- **RQ1:** Does our unified contrastive method outperform state-of-the-art GNN-based and self-supervised recommendation baselines?
- **RQ2:** How does our approach to mitigating false negatives and popularity bias contribute to the model's effectiveness? (Addressed via ablation studies in Section 6.1)
- **RQ3:** What is the underlying semantic mechanism by which our model aligns graph and text modalities? (Addressed via a case study in Section 6.2)

4.1 Experimental Setup

4.1.1 Datasets. To ensure a robust and comprehensive evaluation, we curated a collection of 15 datasets from three distinct sources. The first two sources are standard e-commerce benchmarks: 13 categories from the Amazon Reviews dataset [5, 24] and the Yelp Open Dataset [34]. These are suitable for their link between user-item interactions and English review texts [36, 40]. To assess the generalizability of our method beyond standard e-commerce benchmarks, we also used a semi-public travel site dataset. It contains public user reviews in Japanese [25] and private staying records. We extracted a subset filtered by user (posted at least two reviews) and the period corresponding to each posted review.

To ensure data quality, we followed a standard pre-processing protocol [6, 7]. For each dataset (except the semi-public travel site), we applied the 5-core setting, filtering out users and items with fewer than five interactions. This ensures that every user and item in the evaluation has a minimal history to learn from. The statistics of the processed datasets are summarized in Table 1.

Table 1: Basic statistics of the datasets used in our experiments.

Dataset Category	Users	Items	Reviews
Clothing Shoes and Jewelry	1,786	4,291	13,676
CDs and Vinyl	2,078	3,888	24,604
Tools and Home Improvement	1,290	4,571	18,456
Sports and Outdoors	1,372	4,190	18,660
Beauty and Personal Care	1,405	5,066	23,870
Automotive	1,187	2,833	12,106
Toys and Games	1,045	4,152	17,568
Pet Supplies	1,114	4,241	20,722
Health and Household	820	3,679	14,442
Video Games	870	2,812	17,572
Cell Phones and Accessories	401	1,943	6,872
Patio Lawn and Garden	282	1,373	3,944
Office Products	256	1,287	3,828
Yelp	802	1,274	11,879
Travel	1,275	9,809	20,365

4.1.2 Evaluation Protocol. For each user, we split their interaction history chronologically into training, validation, and test sets with a ratio of 80:10:10. This ensures that the model is evaluated on its ability to predict future interactions based on past behavior.

We adopt a top-K item ranking evaluation protocol. For each user in the test set, the model is tasked with ranking the ground-truth item against all other items in the catalogue that the user has not interacted with in the training set. This "all-ranking" protocol provides a comprehensive and unbiased evaluation of the model's performance, avoiding potential biases from negative sampling during evaluation [15, 27].

We evaluate model performance using two standard top-K ranking metrics [8, 11]:

- **HR@10 (Hit Ratio at 10):** Measures the proportion of users for whom the ground-truth item is ranked within the top 10.

- **NDCG@10 (Normalized Discounted Cumulative Gain at 10):** A position-aware metric that assigns higher scores to hits at higher ranks in the top-10 list.

4.1.3 Baseline Models. We compare our model against a suite of competitive baselines to establish a strong performance benchmark. These models represent a range of approaches, from classic neural models to state-of-the-art GNN and self-supervised methods:

- **NeuMF [7]:** A classic neural recommendation model combining matrix factorization and a multi-layer perceptron to capture both linear and non-linear user-item relationships.
- **SGL [33]:** A self-supervised learning approach that enhances GNNs by creating auxiliary training signals through graph augmentations and contrastive learning.
- **DGCLR [16]:** A self-supervised graph framework that disentangles user and item representations into different semantic factors using a channel-aware contrastive learning objective.
- **NGCF [31]:** A GNN-based model that explicitly encodes high-order connectivity in the user-item interaction graph through embedding propagation.
- **LightGCN [6]:** A simplified GNN that has become a standard, strong baseline, removing non-linearities and feature transformations for efficiency and effectiveness in collaborative filtering.

4.1.4 Implementation Details. To ensure fair and reproducible comparisons, the baseline models, except for DGCLR, were implemented within the RecBole framework [37]. The proposed SymCERE model, its related models, and DGCLR are implemented in non-RecBole implementations, which are available anonymously on Github. For the text encoder, we used the pre-trained LLM "sarashina2.2-0.5b-instruct" for the semi-public travel dataset and "blair-roberta-base" for other datasets, using LoRA. All experiments were performed on an NVIDIA H100 GPU (90GB).

5 RESULTS

This section presents the empirical results of our experiments, addressing the primary research question (RQ1) concerning the performance of our proposed SymCERE method against state-of-the-art baselines.

5.1 Overall Performance Comparison (RQ1)

The comprehensive performance comparison across all 15 datasets is presented in Table 2. The results show that our proposed multi-modal approach, particularly when implemented with an NGCF backbone, generally outperforms the suite of standard baselines across the majority of datasets and on both evaluation metrics.

5.1.1 Hit Ratio (HR@10) Performance. Across the HR@10 metric, our SymCERE method shows strong performance. Our model variants secure the top position in 13 out of the 15 datasets. For instance, on the **Yelp** dataset, our model with an NGCF backbone achieves an HR@10 of 0.3540, which represents a 34.4% relative improvement over the strongest baseline (DGCLR at 0.2560). Similarly, on the **Pet Supplies** dataset, our model achieves an HR@10 of 0.3860, a 38.7% relative improvement over the best-performing baseline, SGL. These results indicate that the integration of textual information via our contrastive method is effective at identifying relevant items for users.

5.1.2 Normalized Discounted Cumulative Gain (NDCG@10) Performance. A similar trend is observed for the NDCG@10 metric, which evaluates the ranking quality of the top-10 recommendations. Our multi-modal model with an NGCF backbone achieves the highest score in 13 out of 15 datasets. The improvements are substantial, highlighting the model's ability to not only find relevant items but also rank them correctly. On the **Yelp** dataset, our model achieves an NDCG@10 of 0.112, a 43.6% relative improvement over the best baseline (NeuMF at 0.0788). On **Cell Phones and Accessories**, our model yields an NDCG@10 of 0.1780, a 18.4% relative improvement over the baseline NGCF. The performance gains are consistent across nearly all datasets, although the margin of improvement varies depending on the domain characteristics, as discussed in Section 6.

In summary, the empirical results support a positive answer to RQ1. Our proposed SymCERE method demonstrates a significant performance advantage over established collaborative filtering and self-supervised baselines across two standard ranking metrics. This indicates that our unified approach to mitigating false negatives and popularity bias can lead to more accurate and effective recommendations.

6 DISCUSSION

The experiments detailed in the previous section show that our proposed method, SymCERE, outperforms a range of strong baselines. These performance gains raise a question: what are the underlying mechanisms that enable this performance? In this section, we deconstruct our model's behavior to answer this question. We first revisit our geometric hypotheses (RQ2) before conducting a series of qualitative case studies to understand the semantic nature of the learned alignment and its impact on performance across different domains (RQ3).

Fundamentally, the triad of challenges we identified—false negatives, popularity bias, and signal ambiguity—are not just concurrent but deeply intertwined. Popularity bias, by encoding popularity as vector magnitude, can severely distort the embedding space. This distortion exacerbates the false negative problem, as true semantic similarity (based on vector direction) is masked by the dominant, non-semantic popularity signal. A model struggling with this skewed space is then ill-equipped to resolve signal ambiguity; it may learn to rely on the strong but noisy popularity signal rather than discerning subtle, objective cues from text. This cascade effect underscores why a unified approach is critical. Addressing these issues in isolation is insufficient, as the unresolved problems will continue to undermine the effectiveness of partial solutions. Our work, therefore, provides a holistic solution by structurally tackling these interconnected challenges simultaneously.

6.1 A Geometric View of Debiasing and Performance Gains (RQ2)

Our methodology is built on the hypothesis that mitigating popularity bias and sampling noise can be achieved through geometric controls over the embedding space. In this section, we quantitatively demonstrate the validity of this hypothesis using empirical data.

Table 2: Unified performance comparison using HR@10 and NDCG@10. This table contrasts standard baselines with our multi-modal method (SymCERE) applied to different GNN backbones. For each metric (row), the best result is highlighted in bold and the second-best is underlined.

Dataset	Metric	Standard Baselines					Ours (SymCERE)		Max. Improv. (%)
		DGCLR	NeuMF	SGL	NGCF	LightGCN	+ NGCF	+ LightGCN	
Clothing, Shoes & Jewelry	HR@10	0.504	0.521	0.530	0.525	0.530	<u>0.538</u>	0.550	↑ 3.8%
	NDCG@10	0.357	0.415	0.412	<u>0.420</u>	0.411	0.425	0.398	↑ 1.2%
CDs & Vinyl	HR@10	0.335	0.453	0.456	0.450	0.419	<u>0.488</u>	0.495	↑ 8.5%
	NDCG@10	0.235	<u>0.324</u>	0.309	<u>0.324</u>	0.285	0.332	0.308	↑ 2.6%
Tools & Home Improvement	HR@10	0.248	0.330	0.336	0.323	0.334	<u>0.402</u>	0.412	↑ 22.5%
	NDCG@10	0.113	0.170	<u>0.175</u>	0.166	0.174	0.187	0.162	↑ 6.7%
Sports & Outdoors	HR@10	0.279	0.381	0.376	0.364	0.337	<u>0.424</u>	0.434	↑ 14.0%
	NDCG@10	0.135	<u>0.209</u>	0.206	0.202	0.191	0.226	0.195	↑ 8.1%
Beauty & Personal Care	HR@10	0.222	0.312	0.312	0.311	0.283	<u>0.386</u>	0.405	↑ 29.8%
	NDCG@10	0.085	<u>0.141</u>	0.135	0.138	0.124	0.162	0.135	↑ 15.1%
Automotive	HR@10	0.367	0.408	0.441	0.424	0.434	<u>0.446</u>	0.470	↑ 6.6%
	NDCG@10	0.210	0.260	<u>0.281</u>	0.278	0.276	0.285	0.265	↑ 1.4%
Toys & Games	HR@10	0.263	0.312	0.329	0.312	0.332	<u>0.415</u>	0.419	↑ 26.2%
	NDCG@10	0.125	0.170	0.174	<u>0.175</u>	0.173	0.198	0.173	↑ 13.0%
Pet Supplies	HR@10	0.192	0.268	0.279	0.268	0.234	<u>0.373</u>	0.387	↑ 38.7%
	NDCG@10	0.054	<u>0.104</u>	0.101	0.101	0.086	0.129	0.102	↑ 24.2%
Health & Household	HR@10	0.212	<u>0.256</u>	<u>0.256</u>	0.246	0.254	0.371	0.360	↑ 44.9%
	NDCG@10	0.064	0.106	<u>0.108</u>	0.102	0.106	0.132	0.107	↑ 22.2%
Video Games	HR@10	0.204	0.227	0.255	0.226	0.254	0.345	<u>0.343</u>	↑ 35.3%
	NDCG@10	0.074	0.099	<u>0.109</u>	0.098	0.104	0.132	<u>0.109</u>	↑ 21.1%
Cell Phones & Accessories	HR@10	0.289	0.254	<u>0.306</u>	0.296	0.299	0.411	0.411	↑ 34.5%
	NDCG@10	0.121	0.114	0.146	<u>0.152</u>	0.140	0.180	0.151	↑ 18.3%
Patio, Lawn & Garden	HR@10	0.340	0.216	0.358	0.351	0.354	<u>0.433</u>	0.450	↑ 25.8%
	NDCG@10	0.149	0.130	<u>0.195</u>	<u>0.195</u>	0.192	0.225	0.189	↑ 15.4%
Office Products	HR@10	0.316	0.195	<u>0.351</u>	0.347	<u>0.351</u>	0.430	0.430	↑ 22.4%
	NDCG@10	0.116	0.093	<u>0.183</u>	0.181	0.178	0.190	0.161	↑ 3.9%
Yelp	HR@10	0.256	0.238	0.193	0.182	0.146	0.344	<u>0.285</u>	↑ 34.4%
	NDCG@10	0.039	0.078	0.054	0.052	0.037	0.112	<u>0.082</u>	↑ 43.6%
Travel	HR@10	0.688	0.768	0.742	0.737	0.695	<u>0.768</u>	0.824	↑ 7.3%
	NDCG@10	0.266	0.268	<u>0.296</u>	0.269	0.279	0.269	0.380	↑ 28.4%

First, we posited a "magnitude-as-noise" hypothesis. Our ablation studies (Table 3) show that a model variant without L2 normalisation performs poorly. This is particularly evident with the non-linear NGCF backbone, where performance drops to near-zero levels. A likely explanation for this catastrophic failure is that without the constraint of normalisation, the embedding magnitudes can grow uncontrollably during training due to NGCF's non-linear transformations, leading to numerical instability and model divergence. This underscores the critical role of L2 normalisation in maintaining a stable training process, especially for more complex

GNN architectures. In contrast, this is likely because, without normalisation, the L2-norm of popular items can dominate similarity calculations, potentially masking the directional preference signal.

L2 normalisation addresses this issue and provides the foundation for improving the effectiveness of contrastive learning. We theorized that normalisation improves **Uniformity** by spreading embeddings more evenly across the unit hypersphere to maximize representational capacity. This claim is empirically supported by the data in Table 4. Across all datasets, applying L2 normalisation increased the standard deviation of cosine similarities, providing quantitative evidence that the embeddings are more uniformly distributed, thus forming a more discriminative representation space.

Table 3: Ablation study for different model backbones. ‘w/o Norm’ denotes the model without L2 normalisation, ‘MM w/ InfoNCE’ is the model variant using InfoNCE for cross-modal loss, and ‘Max. Drop (%)’ shows the performance degradation of the ‘MM w/ InfoNCE’ model compared to the ‘Full’ model.

Dataset	Metric	NGCF Backbone				LightGCN Backbone			
		Full	w/o Norm	MM w/ InfoNCE	Min. Drop (%)	Full	w/o Norm	MM w/ InfoNCE	Max. Drop (%)
Clothing, Shoes and Jewelry	HR@10	0.538	0.002	0.536	↓ 0.2%	0.550	0.548	0.547	↓ 0.6%
	NDCG@10	0.425	0.001	0.428	↑ 0.7%	0.398	0.385	0.384	↓ 3.6%
CDs and Vinyl	HR@10	0.488	0.007	0.490	↑ 0.4%	0.495	0.492	0.463	↓ 6.4%
	NDCG@10	0.332	0.002	0.338	↑ 1.8%	0.308	0.303	0.295	↓ 4.2%
Tools and Home Imp.	HR@10	0.402	0.005	0.397	↓ 1.3%	0.412	0.409	0.396	↓ 3.8%
	NDCG@10	0.187	0.003	0.186	↓ 0.3%	0.162	0.156	0.153	↓ 6.0%
Sports and Outdoors	HR@10	0.424	0.002	0.426	↑ 0.3%	0.434	0.433	0.422	↓ 2.9%
	NDCG@10	0.226	0.001	0.227	↑ 0.6%	0.195	0.187	0.184	↓ 5.6%
Beauty and Personal Care	HR@10	0.386	0.001	0.390	↑ 0.9%	0.405	0.401	0.374	↓ 7.7%
	NDCG@10	0.162	0.000	0.158	↓ 2.4%	0.135	0.128	0.122	↓ 9.6%
Automotive	HR@10	0.446	0.005	0.450	↑ 0.9%	0.470	0.469	0.459	↓ 2.3%
	NDCG@10	0.285	0.002	0.283	↓ 0.6%	0.265	0.254	0.252	↓ 4.9%
Toys and Games	HR@10	0.415	0.001	0.417	↑ 0.5%	0.419	0.416	0.407	↓ 3.0%
	NDCG@10	0.198	0.000	0.197	↓ 0.2%	0.173	0.164	0.162	↓ 6.3%
Pet Supplies	HR@10	0.373	0.004	0.371	↓ 0.7%	0.387	0.383	0.369	↓ 4.6%
	NDCG@10	0.129	0.001	0.124	↓ 3.9%	0.102	0.096	0.093	↓ 8.9%
Health and Household	HR@10	0.358	0.002	0.355	↓ 0.8%	0.357	0.343	0.291	↓ 18.5%
	NDCG@10	0.131	0.001	0.130	↓ 0.8%	0.112	0.102	0.098	↓ 12.5%
Video Games	HR@10	0.332	0.006	0.331	↓ 0.3%	0.326	0.280	0.222	↓ 31.9%
	NDCG@10	0.130	0.002	0.129	↓ 0.8%	0.116	0.099	0.089	↓ 23.3%
Cell Phones and Acc.	HR@10	0.408	0.005	0.397	↓ 2.7%	0.401	0.408	0.357	↓ 11.0%
	NDCG@10	0.178	0.001	0.170	↓ 4.5%	0.161	0.146	0.149	↓ 7.5%
Patio, Lawn and Garden	HR@10	0.433	0.007	0.440	↑ 1.6%	0.450	0.440	0.447	↓ 0.8%
	NDCG@10	0.225	0.002	0.217	↓ 3.7%	0.189	0.179	0.182	↓ 3.5%
Office Products	HR@10	0.430	0.004	0.430	↑ 0.0%	0.430	0.430	0.422	↓ 1.8%
	NDCG@10	0.190	0.002	0.203	↑ 6.9%	0.161	0.147	0.147	↓ 8.8%
Yelp	HR@10	0.344	0.017	0.350	↑ 1.7%	0.293	0.224	0.210	↓ 28.3%
	NDCG@10	0.112	0.003	0.113	↑ 0.9%	0.089	0.070	0.068	↓ 23.6%
Travel	HR@10	0.768	0.040	0.789	↑ 2.7%	0.824	0.829	0.695	↓ 15.7%
	NDCG@10	0.269	0.006	0.298	↑ 10.8%	0.380	0.390	0.228	↓ 40.0%

These geometric improvements—enhanced uniformity and dimensional utilization—are mechanisms that can suppress noise and bias, leading to gains in recommendation performance.

Interestingly, our ablation study (Table 3) also reveals that the InfoNCE-based model variant occasionally outperforms the SINCERE-based model when using the non-linear NGCF backbone. We hypothesize this is a result of implicit regularization; while InfoNCE’s treatment of false negatives is flawed, the resulting noisy gradient may prevent a complex model like NGCF from overfitting. In contrast, SINCERE’s theoretically cleaner signal might not provide the same regularizing effect on certain datasets. This suggests a nuanced trade-off between the theoretical correctness of a loss function and its practical behavior with different model architectures.

6.2 Qualitative Case Studies: From Semantics to Performance (RQ3)

The quantitative results and the geometric intuition above raise a deeper qualitative question: what semantic information does the model actually use to align embeddings, and why does its effectiveness vary across datasets? To answer this (RQ3), we performed in-depth case studies on three representative datasets exhibiting high, medium, and low performance gains, using LIME [28] to analyze influential review terms.

6.2.1 High-Improvement Case: Pet Supplies. The *Pet Supplies* dataset showed one of the largest performance improvements (e.g., +38.7% in HR@10). Our analysis reveals that the model’s success stems from its ability to prioritize objective, information-rich words. The nature of these words shifts from technical specifications to product functionalities and target subjects.

Table 4: Detailed statistical comparison of similarity distributions, including key distributional statistics. Metrics for models with and without L2 normalisation are presented for each dataset.

Dataset	Without Normalisation						With Normalisation					
	Mean	Std. Dev.	Min	25%	75%	Max	Mean	Std. Dev.	Min	25%	75%	Max
Pet Supplies	0.001	0.006	-0.037	-0.003	0.005	0.028	0.001	↑ 0.010	-0.040	-0.006	0.008	↑ 0.046
Office Products	0.000	0.007	-0.030	-0.003	0.006	0.026	0.000	↑ 0.011	-0.037	-0.007	0.008	↑ 0.040
Clothing, Shoes & Jewelry	0.001	0.008	-0.034	-0.004	0.007	0.035	0.001	↑ 0.011	-0.038	-0.006	0.008	↑ 0.048

As summarized in Table 5, the model identifies keywords related to pet type ('**puppy**', '**cat**'), product category ('**crate**', '**leash**', '**toy**'), and functional attributes ('**chew**', '**durable**', '**healthy**') as strong positive contributors to similarity.

Consider a review like: "My **puppy** loves this **durable chew toy**. Perfect size for teething." The model learns that objective terms like 'puppy', 'durable', and 'chew toy' provide a stronger alignment signal than the generic sentiment 'loves'. These keywords are information-dense and describe the specific context of use, allowing the model to distinguish between a toy for a teething puppy and a food bowl for a large dog. This semantic precision explains the substantial performance lift. The model resolves "Signal Ambiguity" by focusing on the reason for the user's interaction, not just the interaction itself.

6.2.2 Mid-Improvement Case: Office Products. The *Office Products* dataset represents a middle-ground case. A key characteristic of this domain is its blend of reviews containing objective, factual terms with those describing subjective user experience. For instance, in reviews concerning printer ink, objective terms that directly indicate product compatibility or malfunction—such as '**HP**', '**genuine**', '**clogged**', and '**dry**'—are correctly identified by the model as important signals. On the other hand, reviews for writing instruments like fountain pens are centered on subjective terms dependent on personal taste, such as '**nib**', '**writing feel**', and '**design**'.

We postulate that this duality is what leads to the moderate performance improvement. The model likely achieves high precision in fact-based sub-domains, such as printer-related issues, but is partially confounded by signal ambiguity in the highly subjective sub-domain of writing instruments. Consequently, the overall rate of improvement is not as pronounced as in a more uniformly objective dataset like Pet Supplies.

6.2.3 Low-Improvement Case: Clothing, Shoes and Jewelry. In contrast, the *Clothing, Shoes and Jewelry* dataset saw less improvement. Our analysis suggests this is due to a combination of factors that create high signal ambiguity.

Firstly, the domain is dominated by subjective aesthetic judgements ('**beautiful**', '**style**', '**elegant**') that lack an objective benchmark. Secondly, even seemingly factual terms like '**size**' and '**fit**' are personal and context-dependent. Thirdly, our LIME analysis revealed a presence of non-English (primarily Spanish) reviews (e.g., '**producto**', '**calidad**'). If the underlying language model is mainly English-trained, this can act as noise, affecting alignment accuracy. This combination of subjectivity, contextual dependency, and multilingual noise (summarized in Table 7) can make it difficult

for the model to extract a clear semantic signal, thereby limiting performance gains.

6.3 Summary of Findings

Our analysis suggests that robust cross-modal alignment is a form of semantic and contextual understanding. The case studies reveal that the performance of our method is correlated with the nature of the language in a given domain. The model excels when it can ground its understanding in objective, factual, and information-rich keywords that describe a product's function, attributes, and context of use. As the language becomes more subjective and ambiguous, the model's ability to disambiguate signals and create a well-structured geometric space may diminish, limiting performance gains. This answers RQ3 and highlights that a key aspect of our model's performance is its learned ability to identify and prioritize reliable and objective signals within user-generated text. This finding implies that future work on review-based systems may benefit from focusing on objective feature extraction rather than advanced sentiment analysis.

7 LIMITATIONS AND FUTURE WORK

Our study demonstrates that SymCERE yields substantial performance gains, particularly when employing a linear graph encoder like LightGCN. This success stems from the natural compatibility between LightGCN's linear propagation, which primarily encodes information in the direction of embeddings, and our methodology's reliance on L2 normalisation.

However, a potential limitation and a promising avenue for future research emerges when considering non-linear encoders such as NGCF. Our ablation study revealed that applying SymCERE to NGCF, while still effective, operates within a potential tension: NGCF's non-linear transformations may learn to encode meaningful information within the magnitude of the embedding vectors, whereas our L2 normalisation purposefully discards this magnitude to mitigate popularity bias. This suggests that the full potential of sophisticated non-linear architectures may not be realised due to this architectural conflict.

Future work could proceed in two main directions. First, one could design novel non-linear GNN architectures that are inherently compatible with hyperspherical embeddings, learning representations where directional information is paramount by design. Second, it remains an open question whether a more advanced normalisation technique could be developed—one that can disentangle and preserve potentially useful information encoded in vector magnitudes while still neutralising the detrimental effects of popularity

Table 5: LIME analysis summary of word contributions to similarity in the Pet Supplies dataset. The model prioritizes objective terms related to product function, type, and attributes.

Keyword Type	Example Words	Contribution	Reason
Product Function/Type	'chew', 'toy', 'crate', 'leash'	Strong Signal	Objective, high-impact functional term
Target/Attribute	'puppy', 'dog', 'durable', 'healthy'	Supporting Signal	Objective attribute of product/user
Generic Sentiment	'love', 'happy', 'great', 'nice'	Weak / Ambiguous	Subjective, lacks specific details
Stop-words	'I', 'my', 'the', 'a', 'it'	Noise / Detrimental	High frequency, no semantic value

Table 6: LIME analysis summary of word contributions to similarity in the Office Products dataset. The model leverages a mix of objective product specifications and subjective feedback.

Keyword Type	Example Words	Contribution	Reason
Product Spec/Function	'ink', 'printer', 'cartridge', 'genuine'	Strong Signal	Objective, high-impact spec/function
Subjective Experience	'nib', 'pen', 'smooth', 'paper'	Context-Dependent Signal	Subjective user experience/taste
Generic Sentiment	'great', 'good', 'nice', 'works'	Weak / Ambiguous	Subjective, lacks specific details
Stop-words	'I', 'my', 'the', 'a', 'it'	Noise / Detrimental	High frequency, no semantic value

Table 7: LIME analysis summary for the Clothing, Shoes and Jewelry dataset. The model may struggle with highly subjective, context-dependent, and multilingual terms.

Keyword Type	Example Words	Contribution	Reason
Subjective Aesthetics	'beautiful', 'elegant', 'style', 'looks'	Highly Ambiguous	Personal taste, no objective benchmark
Context-Dependent Fit	'size', 'fit', 'small', 'large'	Ambiguous	Depends on user's body/preference
Foreign Language	'producto', 'calidad', 'excelente'	Noise / Mismatch	Potential language model mismatch
Generic Sentiment	'love', 'nice', 'good', 'wear'	Weak / Ambiguous	Low information, lacks specific reason

Table 8: Examples of review texts providing Strong/Objective versus Weak/Subjective signals. Objective product attributes (underlined) serve as strong anchors for alignment, while subjective or anecdotal terms (underlined) provide weaker signals.

Dataset	Signal Type	Review Text with Highlighted Keywords
Pet Supplies	Strong / Objective	i rescued a toy poodle and she has <u>no teeth</u> . i feed her canned foods because of the <u>softness</u> and was feeding ...
	Weak / Subjective	This is a great container, and perfectly holds a <u>30 pound bag</u> of Taste of the Wild dog food. The wheels move ...
Office Products	Strong / Objective	We had this for about 3 months before our puppy realized he could <u>tear it apart</u> . We had this on the bottom ...
	Weak / Subjective	My one year old cat... took one of the packs of treats, opened it (<u>ripped it open</u>) and I caught her red-pawed ...
Clothing, Shoes & Jewelry	Strong / Objective	Stick with the original, generic inks will <u>damage your machine</u> . The <u>print quality</u> is better with HP and the ...
	Weak / Subjective	The <u>ink dried up</u> in the lines ruining my machine. Also, the print quality was not as good. Keep away.
	Strong / Objective	I ordered this pen based on its <u>looks</u> , the Parker <u>reputation</u> and its reasonable price. The nib that came with it ...
	Weak / Subjective	I have always come to expect <u>great quality</u> from Parker products. This is no exception. I am very satisfied ...
	Strong / Objective	I really do love this hat. It is <u>very warm</u> and completely <u>covers my ears</u> . I use Watch Caps to hike in cold ...
	Weak / Subjective	it wears like a verry expensive piece time piece it also feels <u>verry light</u> like a feather easy to set even a child ...
	Strong / Objective	The first set I ordered were <u>damaged</u> the material had <u>holes in it</u> . I returned the first set and they sent me a ...
	Weak / Subjective	I was impressed with this piece of jewelry. It was not only <u>beautiful</u> in appearance but also very well made ...

bias. Such research would further enhance the synergy between advanced graph encoders and robust contrastive learning methods.

And while our LIME-based analysis provides compelling evidence for "semantic anchoring," we acknowledge that LIME is a local surrogate model. Therefore, a potential limitation is that our interpretation is based on local explanations, and the extent to which this mechanism represents the model's global behavior remains an open question. Future work could employ global explanation techniques to validate whether this semantic anchoring phenomenon holds true across the entire representation space.

8 CONCLUSION

We addressed the interconnected challenges of false negatives, popularity bias, and signal ambiguity with SymCERE, a unified contrastive learning method. By integrating a symmetric SINCERE loss with hyperspherical projection (L2 normalisation), SymCERE significantly outperforms strong baselines on both NDCG@10 and HR@10 across 15 datasets (e.g., +43.6% NDCG@10). Our key finding, termed 'semantic anchoring', reveals the model learns to prioritize objective, information-rich, and domain-specific attributes over superficial sentiment, showing that robust alignment stems from factual understanding. This unified and principled approach of structurally mitigating biases thus creates a path toward recommender systems that are not only more accurate but also more robust, interpretable, and fair by promoting item diversity.

REFERENCES

- [1] Tsai-Shien Chen, Wei-Chih Hung, Hung-Yu Tseng, Shao-Yi Chien, and Ming-Hsuan Yang. 2022. Incremental False Negative Detection for Contrastive Learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. <https://openreview.net/forum?id=dDjSKKA5TP1>
- [2] Zhengyi Cheng, Ying-Chun Lin, Yassine Benajiba, Ori Levi, Hang Li, and Panayiotis Tsaparas. 2024. DynLLM: Dynamic and LLM-Assisted Scalable Architecture for Product Search. In *Companion Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*. 118–126. <https://doi.org/10.1145/3613904.3642828>
- [3] Kathleen M. Feeney, A. Stephen McGough, and Paula L. Smith. 2023. SINCERE: Supervised-learning-based and Noise-robust Contrastive Estimation for a Regression task. In *2023 International Joint Conference on Neural Networks, IJCNN 2023, Gold Coast, Australia, June 18-23, 2023*. 1–8. <https://doi.org/10.1109/IJCNN54540.2023.10191590>
- [4] Michael U. Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010 (JMLR Workshop and Conference Proceedings, Vol. 9)*. JMLR.org, 297–304. <http://proceedings.mlr.press/v9/gutmann10a.html>
- [5] Ruining He and Julian J. McAuley. 2016. Ups and Downs: Modeling the Visual Evolution of Fashion Trends with One-Class Collaborative Filtering. In *Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11 - 15, 2016*. 507–517. <https://doi.org/10.1145/2872427.2883037>
- [6] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation. In *SIGIR '20: The 43rd International ACM SIGIR conference on research and development in Information Retrieval, Virtual Event, China, July 25-30, 2020*. 639–648. <https://doi.org/10.1145/3397271.3401063>
- [7] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural Collaborative Filtering. In *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*. 173–182. <https://doi.org/10.1145/3038912.3052569>
- [8] Jonathan L. Herlocker, Joseph A. Konstan, Loren G. Terveen, and John T. Riedl. 2004. Evaluating collaborative filtering recommender systems. *ACM Trans. Inf. Syst.* 22, 1 (2004), 5–53. <https://doi.org/10.1145/963770.963772>
- [9] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. <https://openreview.net/forum?id=nZeVKeeFyf9>
- [10] T. Huynh, M. El-Khamy, J. Lee, and M. Le. 2022. Boosting Contrastive Learning with False Negative Cancellation. In *2022 IEEE Winter Conference on Applications of Computer Vision (WACV) (2022)*. 2794–2803. <https://doi.org/10.1109/WACV51458.2022.00285>
- [11] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf. Syst.* 20, 4 (2002), 422–446. <https://doi.org/10.1145/582415.582418>
- [12] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised Contrastive Learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. <https://proceedings.neurips.cc/paper/2020/hash/d89a66c7c80a29b1bdbab0f2a1a94af8-Abstract.html>
- [13] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. <https://openreview.net/forum?id=SLJuayYgl>
- [14] Anastasia Klimashewska, Dietmar Jannach, Mehdi Elahi, and Christoph Trattner. 2024. A Survey on Popularity Bias in Recommender Systems. *User Model. User Adapt. Interact.* 34, 2 (2024), 177–234. <https://doi.org/10.1007/s11257-024-09406-0>
- [15] Walid Krichene and Steffen Rendle. 2020. On Sampled Metrics for Item Recommendation. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*. 1748–1757. <https://doi.org/10.1145/3394486.3403226>
- [16] Weng-Jun Li, Chen Gao, Yu Zheng, Xiang Wang, Xiangnan He, and Depeng Jin. 2022. Disentangled Graph Contrastive Learning for Review-based Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (CIKM '22)*. 1269–1278.
- [17] Ying-Chun Lin, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, and Yassine Benajiba. 2021. Denoising implicit feedback for recommendation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining, WSDM 2021, Virtual Event, Israel, March 8-12, 2021*. 678–686. <https://doi.org/10.1145/3437963.3441793>
- [18] Ying-Chun Lin, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, and Yassine Benajiba. 2022. Neutralizing Popularity Bias in Recommendation via Causal Embeddings. *CoRR* abs/2208.01314 (2022). arXiv:2208.01314 <https://arxiv.org/abs/2208.01314>
- [19] Ying-Chun Lin, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, and Yassine Benajiba. 2023. Denoising Self-Attentive Sequential Recommendation. In *Proceedings of the 16th ACM International Conference on Web Search and Data Mining, WSDM 2023, Singapore, 27 February 2023 - 3 March 2023*. 709–717. <https://doi.org/10.1145/3539597.3570451>
- [20] Ying-Chun Lin, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, and Yassine Benajiba. 2023. Mitigating Popularity Bias in Recommendation via User-Centric Popularity-Disentangled Representations. *CoRR* abs/2309.00695 (2023). arXiv:2309.00695 <https://arxiv.org/abs/2309.00695>
- [21] Ying-Chun Lin, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, Yassine Benajiba, and Jian-Yun Nie. 2024. Triple-Aware Word-level Debating for Review-based Recommendation. *CoRR* abs/2402.13851 (2024). arXiv:2402.13851 <https://arxiv.org/abs/2402.13851>
- [22] Jing-Kai Lou, Yen-Wei Chen, Yen-Cheng Chang, Han-Shen Huang, Jiong-Kai Wang, Yassine Benajiba, and Jian-Yun Nie. 2024. Robust Recommender System with Dual-channel Negative Feedback. *CoRR* abs/2402.13854 (2024). arXiv:2402.13854 <https://arxiv.org/abs/2402.13854>
- [23] Zhe Luo, Jian-Yun Nie, and Yassine Benajiba. 2024. MoLAR: a Modular framework for LLM-as-Recommender. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining, WSDM 2024, Merida, Yucatan, Mexico, March 4-8, 2024*. 1269–1272. <https://doi.org/10.1145/3616855.3639578>
- [24] Julian J. McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. 2015. Image-Based Recommendations on Styles and Substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, Santiago, Chile, August 9-13, 2015*. 43–52. <https://doi.org/10.1145/2766462.2767755>
- [25] Rakuten Group, Inc. 2020. Rakuten Travel data. <https://doi.org/10.32130/idr.2.2> ([dataset](https://doi.org/10.32130/idr.2.2)).
- [26] Xuhui Ren, Jian-Yun Nie, and Panayiotis Tsaparas. 2024. A Survey on Large Language Models for Personalization. *CoRR* abs/2403.05924 (2024). arXiv:2403.05924 <https://arxiv.org/abs/2403.05924>
- [27] Steffen Rendle, Walid Krichene, Li Zhang, and John R. Anderson. 2020. Neural Collaborative Filtering vs. Matrix Factorization Revisited. In *Fourteenth ACM Conference on Recommender Systems, RecSys 2020, Virtual Event, Brazil, September 22-26, 2020*. 240–248. <https://doi.org/10.1145/3383313.3412488>
- [28] Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and*

- Data Mining, San Francisco, CA, USA, August 13-17, 2016.* 1135–1144. <https://doi.org/10.1145/2939672.2939778>
- [29] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning with Contrastive Predictive Coding. *CoRR* abs/1807.03748 (2018). arXiv:1807.03748 <http://arxiv.org/abs/1807.03748>
- [30] Tongzhou Wang and Phillip Isola. 2020. Understanding Contrastive Representation Learning through Alignment and Uniformity on the Hypersphere. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 9929–9939. <http://proceedings.mlr.press/v119/wang20k.html>
- [31] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '19, Paris, France, July 21-25, 2019.* 165–174. <https://doi.org/10.1145/3331184.3331267>
- [32] Chuhan Wu, Fangzhao Wu, Junxin Liu, and Yongfeng Huang. 2014. A Survey on User Modeling with Review Data. In *Artificial Intelligence: 21st ECAI 2014, Prague, Czech Republic, August 18-22, 2014 - Including 21st European Conference on Artificial Intelligence, PAIS 2014.* 1005–1006. <https://doi.org/10.3233/978-1-61499-419-0-1005>
- [33] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised Graph Learning for Recommendation. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021.* 726–735. <https://doi.org/10.1145/3404835.3462862>
- [34] Yelp. 2024. Yelp Open Dataset. <https://www.yelp.com/dataset>.
- [35] Yong-Feng Zhang and Xu Chen. 2021. Causal Inference in Recommender Systems: A Survey and Future Directions. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM 2021, Virtual Event, Queensland, Australia, November 1 - 5, 2021.* 4415–4424. <https://doi.org/10.1145/3459637.3481928>
- [36] Wayne Xin Zhao, Yupeng Hou, Junjie Zhang, Zihan Lin, Shuqing Bian, Xingyu Pan, Yushuo Chen, Hong-Jian Ai, Jiarui Qin, Yifei Li, and Ji-Rong Wen. 2023. A Survey of Large Language Models. *CoRR* abs/2303.18223 (2023). arXiv:2303.18223 <https://arxiv.org/abs/2303.18223>
- [37] Wayne Xin Zhao, Shanlei Mu, Yupeng Hou, Zihan Lin, Yushuo Chen, Xingyu Pan, Kaiyuan Li, Yujie Lu, Hui Wang, Changxin Tian, Yingqian Min, Zhichao Feng, Xinyan Fan, Xu Chen, Pengfei Wang, Wendi Ji, Yaliang Li, Xiaoling Wang, and Ji-Rong Wen. 2021. RecBole: Towards a Unified, Comprehensive and Efficient Framework for Recommender Systems. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM 2021, Virtual Event, Queensland, Australia, November 1 - 5, 2021.* 4653–4664. <https://doi.org/10.1145/3459637.3482051>
- [38] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Depeng Jin, and Yong Li. 2021. Disentangling User Interest and Conformity for Recommendation with Causal Embedding. In *The Web Conference 2021, Proceedings of the World Wide Web Conference WWW 2021, Ljubljana, Slovenia, April 19-23, 2021.* 2980–2991. <https://doi.org/10.1145/3442381.3449788>
- [39] Yue Zhong, Yong-Xin Zhuang, and Tai-Xiang Jiang. 2023. SINCERE: A New Noise-Robust Training Strategy for Various Supervised Learning Tasks. *CoRR* abs/2301.11262 (2023). arXiv:2301.11262 <https://arxiv.org/abs/2301.11262>
- [40] Chaoqun Zhou, Yuan-Wei Wang, Yuan-Fang Li, and Gholamreza Haffari. 2023. A Survey on Review-based Recommender Systems. *CoRR* abs/2301.01775 (2023). arXiv:2301.01775 <https://arxiv.org/abs/2301.01775>