# FITRep: Attention-Guided Item Representation via MLLMs

Guoxiao Zhang*
Ao Li*
zhangguoxiao@meituan.com
liao27@meituan.com
Meituan
Beijing, China

Tan Qu*
Meituan
Beijing, China
qutan@meituan.com

Qianlong Xie
Meituan
Beijing, China
xieqianlong@meituan.com

Xingxing Wang
Meituan
Beijing, China
wangxingxing04@meituan.com

## Abstract

Online platforms usually suffer from user experience degradation due to near-duplicate items with similar visuals and text. While Multimodal Large Language Models (MLLMs) enable multimodal embedding, existing methods treat representations as black boxes, ignoring structural relationships (e.g., primary vs. auxiliary elements), leading to **local structural collapse problem**. To address this, inspired by Feature Integration Theory (FIT), we propose **FITRep**, the first *attention-guided, white-box item representation* framework for fine-grained *item deduplication*. **FITRep** consists of: (1) **Concept-Hierarchical Information Extraction** (CHIE), using MLLMs to extract hierarchical semantic concepts; (2) **Structure-Preserving Dimensionality Reduction** (SPDR), an adaptive UMAP-based method for efficient information compression; and (3) **FAISS-Based Clustering** (FBC), a FAISS-based clustering that assigns each item a unique cluster id using FAISS. Deployed on Meituan's advertising system, **FITRep** achieves +3.60% CTR and +4.25% CPM gains in online A/B tests, demonstrating both effectiveness and real-world impact.

## CCS Concepts

• **Information systems → Information retrieval**.

## Keywords

Multimodal, Large Language Models, Dimensionality Reduction

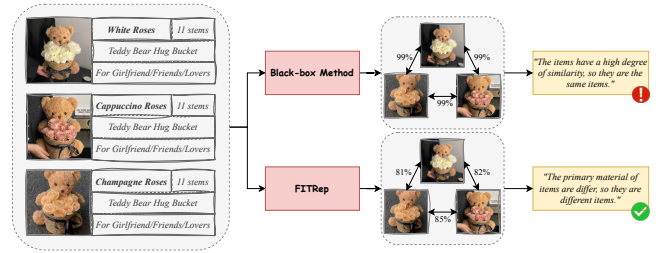*Both authors contributed equally to this research.

**Figure 1: An illustration of the Local structural collapse problem.**

## 1 Introduction

Most online platforms leverage multimodal information such as images and textual titles to enhance user engagement [1]. However, the repeated display of items with near-identical visual and textual content degrades user experience by introducing redundancy. To mitigate this, we propose a multimodal similarity metric for effective *item deduplication* task that identifies and removes redundant or near-duplicate items. Effective *item deduplication* critically depends on the quality of multimodal item representations.

Recently, Multimodal Large Language Models (MLLMs) have shown impressive capabilities in multimodal understanding and generating [2, 3]. Existing approaches for obtaining item embeddings via MLLMs broadly fall into two categories: (1) *prompt-based methods* [4, 5], where simple prompts are employed to generate image summaries which are transformed into latent vectors using encoding models like BERT [6]; and (2) *end-to-end joint learning*(NoteLLM-2 [7]), where lightweight MLLMs are fine-tuned exploiting co-occurrence items for better multimodal representation. We refer to such paradigms as *black-box coarse-grained item representation*, which suffers from **local structural collapse problem**: item images and texts inherently contain structural relationships, such as objective/subjective, form/content, primary/auxiliary materials, etc. Those black-box approaches, lacking explicit modeling of these structures, produce representations that may lose these structural relationships. As an example shown in Figure 1, due to the absence of attention mechanisms that differentiate primary from auxiliary materials, items with different primary materials

but highly similar auxiliary materials receive spuriously high multimodal similarity scores, leading to false-positive duplicates.

Inspired by Feature Integration Theory (FIT) [8], which suggests that the brain processes basic visual features automatically and separately in a pre-attentive stage, and then combines them accurately through focused attention during a subsequent attentive stage, we propose a holistic framework for *attention-guided, white-box item representation* that includes: (1) **Concept-Hierarchical Information Extraction** (CHIE) using MLLMs; (2) **Structure-Preserving Dimensionality Reduction** (SPDR) via adaptive UMAP, and (3) **FAISS-Based Clustering** (FBC) that assigns differentiated weights to distinct elements based on their attention coefficients, thereby generating the final item representation for FAISS based clustering.

This paper makes three key contributions: (1) We propose **FITRep**, the first *attention-guided, white-box item representation* framework that leverages MLLMs for interpretable and fine-grained perception. (2) **FITRep** comprises three core components: **CHIE** using MLLMs for high-efficiency information extraction, **SPDR** via adaptive UMAP for efficient information compression, and **FBC** for scalable duplicate detection. (3) We deploy **FITRep** on Meituan's advertising system for *item deduplication* and *CTR prediction*; online A/B tests demonstrate significant improvements of **+3.60% in CTR** and **+4.25% in CPM**, validating its real-world effectiveness.

## 2 Methodology

Our framework, illustrated in Figure 4 includes: (1) **CHIE** using MLLMs; (2) **SPDR** via adaptive UMAP; (3) **FBC** for scalable duplicate detection. Further details are provided in the following section.
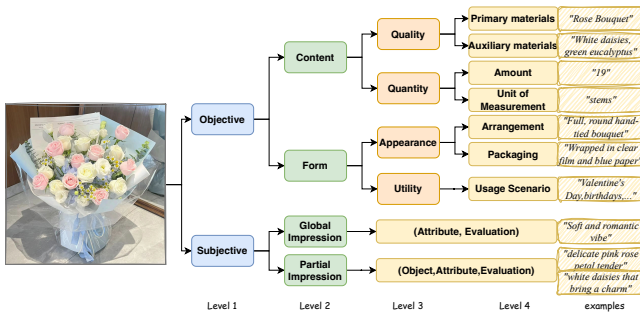


**Figure 2: An illustration of the *concept hierarchies*.**

### 2.1 Concept-Hierarchical Information Extraction (CHIE)

Given that item images and text inherently exhibit structural relationships, we begin by categorizing the multidimensional item information into four hierarchical concept levels (Levels 1–4), as depicted in Figure 2. Among these, Level 4 represents the finest granularity of item representation, consisting of $D$ distinct dimensions. In this study, we set $D = 8$, as we typically merge the *Amount* and *Unit of Measurement* dimensions into a single dimension *Quantity* for practical purposes.
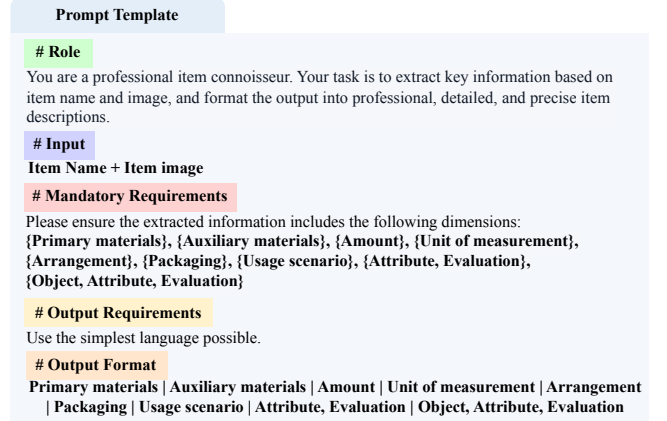


**Figure 3: Example of a structured prompt for FITRep.**

Concurrently, we design a structured prompt that integrates textual and visual inputs to guide the MLLM in extracting dimension-specific item representations (see Figure 3). Then we utilize a pretrained Text Encoder [9] to extract dense vector representations from the dimension-specific textual descriptions:

$$\mathbf{v}^k = \text{Encoder}\left(t^k\right), \quad k = 1, \ldots, D - 1. \tag{1}$$

Where $t^k$ represents MLLM generated textual descriptions in the $k$-th conceptual dimension (excluding quantity), we directly use the numerical quantity dimension without representation extraction.

### 2.2 Structure-Preserving Dimensionality Reduction (SPDR)

To reduce high-dimensional embeddings extracted from Equation (1) while preserving semantic structure, we apply an adaptive parameterized UMAP [10], which dynamically adjusts output dimensionality for different conceptual dimension:

$$\mathbf{e}^k = \text{UMAP}\left(k, \mathbf{v}^k\right), \quad k = 1, \ldots, D - 1. \tag{2}$$

To facilitate subsequent computations, we represent the embeddings $\mathbf{e}^k$ in their $L_2$-normalized form $\|\mathbf{e}^k\|_2$.

### 2.3 FAISS Based Clustering (FBC)

To enable efficient online deduplication, we assign each item a unique cluster id offline. Given the scale of over ten million items, as is shown in Figure 5, we employ FAISS [11] for efficient offline similarity search and clustering.

*2.3.1 Adapting to FAISS data processing.* We apply dimension-specific weights $w^k$ to the corresponding embeddings $\mathbf{e}^k$, producing weighted embeddings for the $D-1$ dimensions (excluding quantity).

The item quantity $q$ is treated separately and mapped nonlinearly to a point $(\cos\theta, \sin\theta)$ on the unit circle's first quadrant:

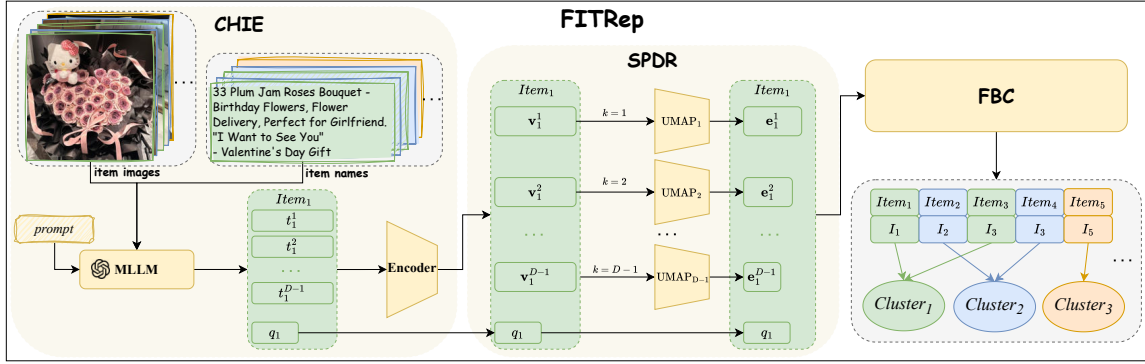$$\theta = \frac{\pi}{2} \times \left(1 - e^{-\alpha(q-1)/Q}\right) \tag{3}$$
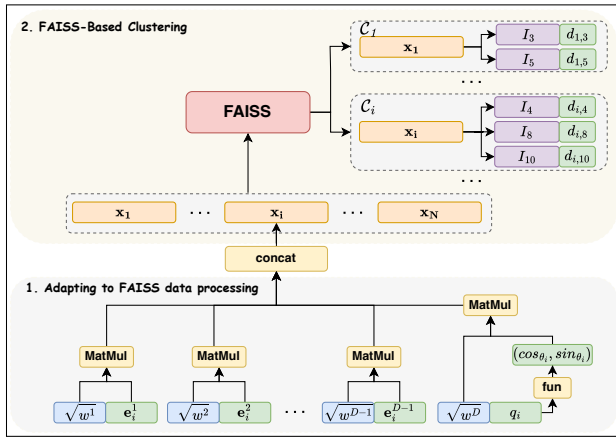
Figure 4: Overview of our proposed method.



Figure 5: Architecture of the FBC.

where $\alpha$ controls the nonlinearity of the mapping and $Q$ is the maximum possible value of $q$. This ensures that angular distance between points reflects semantic similarity in quantity.

Finally, we construct the unified item representation $\mathbf{x}$:

$$\mathbf{x} = \left[ \sqrt{w^1}\mathbf{e}^1, \ldots, \sqrt{w^{D-1}}\mathbf{e}^{D-1}, \sqrt{w^D}\cos\theta, \sqrt{w^D}\sin\theta \right]^\top. \quad (4)$$

where $w^D$ is the weight of item quantity $q$.

*2.3.2 FAISS-Based Clustering.* Based on the processed item embeddings, we perform FAISS-based retrieval for each item $I_i$, returning all items within an $L_2$ distance threshold $\tau$:

$$C_i = \{I_j \mid d_{i,j} < \tau, \ j \in \{1, \ldots, N\}\}$$
$$= \text{FAISS\_search}(I_i, \mathcal{I}, \tau), \quad \forall i \in \{1, \ldots, N\} \quad (5)$$

where $d_{i,j} = 2(1 - \mathbf{x}_i \cdot \mathbf{x}_j)$ denotes the $L_2$ distance between embeddings, and $N$ is the total number of items. Due to potential cluster overlaps, we apply a deduplication step to assign each item a unique cluster ID, resulting in $n$ final non-overlapping clusters.

## 3 EXPERIMENTS

To validate the effectiveness of **FITRep**, we compare it against a black-box method (**BBM**), which generate multimodal summaries through simple textual prompts, encode them using BERT, and apply dimensionality reduction. We exclude NoteLLM-2 [7] from our main comparison, as it is designed to incorporate cooperative signals through fine-tuning, whereas our goal is to learn intrinsic semantic representations independent of cooperative signals. We treat *item deduplication* (including *duplicate item identification* and *duplicate item removal*) as the primary task, with CTR prediction serving as an indirect validation of embedding quality, assuming that semantically meaningful representations improve item-user matching.

### 3.1 Duplicate Item Identification

For *duplicate item identification*, we evaluate on a manually annotated dataset, defining positive pairs as duplicates and negative pairs as non-duplicates.

Table 1: Offline Item deduplication performance

| Methods | Precise | Recall | F1 |
|---------|---------|--------|-----|
| **BBM** | 56.9% | **95.0%** | 71.2% |
| **FITRep** | **88.1%** | 87.5% | **87.8%** |

As shown in Table 1, **FITRep** significantly improves precision (88.1% vs. 56.9%) over **BBM**, with only a small drop in recall (87.5% vs. 95.0%), resulting in a much higher F1-score (87.8% vs. 71.2%). This demonstrates that attention-guided fine-grained embeddings enable more accurate duplicate detection.

### 3.2 Impact of Dimensionality Reduction Methods

On manually annotated dataset, we evaluate the impact of various dimensionality reduction methods. As shown in Table 2, **FITRep**, leveraging parameterized UMAP, surpasses both PCA-based and VAE-based methods, achieving the highest precision (88.1%), recall (87.5%), and F1-score (87.8%), demonstrating its superior ability to preserve structure during dimensionality reduction.

**Table 2: Comparison of Dimensionality Reduction Methods.**

| Methods | Precise | Recall | F1 |
|---|---|---|---|
| PCA-based | 75.3% | 81.4% | 78.2% |
| VAE-based | 81.2% | 85.7% | 83.4% |
| **FITRep** | **88.1%** | **87.5%** | **87.8%** |

## 3.3 CTR Prediction

We evaluate our framework on CTR prediction using a large-scale industrial dataset from Meituan's advertising system, which captures real user interactions in a massive local-services ecosystem. Key statistics are summarized in Table 3.

**Table 3: Dataset statistics.**

| Dataset | #Requests | #Users | #Items |
|---|---|---|---|
| Meituan | 98,362,548 | 24,215,750 | 5,365,286 |

Our proposed method, **CTR_FITRep**, integrates item embeddings for CTR prediction, and is compared against **CTR_BBM**, which uses item representations generated by **BBM** method.

**Table 4: Offline performance on the Meituan dataset.**

| Model | AUC ↑ | LogLoss ↓ |
|---|---|---|
| CTR_BBM | 0.6580 | 0.0234 |
| **CTR_FITRep** | **0.6640** | **0.0215** |

Following standard practice, we report **AUC** and **LogLoss**. As shown in Table 4, the gain over **CTR_BBM** (+0.6pp in AUC) demonstrates that replacing black-box prompt summaries with white-box, attention-guided representations better captures item semantics, leading to more accurate user interest modeling.

## 3.4 Online Deployment and A/B Test Results

We deploy **FITRep** on Meituan's advertising system to support two key tasks: *item deduplication* (primary) and *CTR prediction* (auxiliary).

For efficient online *item deduplication*, we leverage **FITRep** embeddings to cluster items during offline processing and store their cluster IDs in Redis. At serving time, for each request, we retrieve the cluster IDs of all candidate items, group those sharing the same cluster ID as duplicates, and retain only the highest-ranked item per cluster, yielding the **FITRep_ID** strategy.

To further enhance CTR prediction, we integrate **FITRep**-based semantic representations into the ranking model (**CTR_FITRep_ID**). For comparison, we also evaluate a variant that uses FITRep embeddings without *item deduplication* (**CTR_FITRep**).

We conducted an online A/B test from October 6–14, 2025. As shown in Table 5, **FITRep_ID** improves user experience by eliminating redundant ads, achieving +2.15% CTR and +2.40% CPM over the Baseline with negligible latency overhead (21.5 ms vs. 21.1 ms). When combined with the CTR model, **CTR_FITRep_ID** delivers

substantial gains of +3.60% in CTR and +4.25% in CPM, significantly outperforming both the Baseline and the non-deduplicated variant (**CTR_FITRep**).

**Table 5: Online A/B testing results on Meituan's advertising system.**

| Method | CTR Gain | CPM Gain | Latency |
|---|---|---|---|
| Baseline | — | — | 21.1 ms |
| **FITRep_ID** | +2.15% | 2.40% | 21.5 ms |
| CTR_FITRep | +1.50% | +1.68% | 21.3 ms |
| **CTR_FITRep_ID** | **+3.60%** | **+4.25%** | 21.8 ms |

## 4 Conclusion

We identify the *local structural collapse* problem in existing black-box multimodal representations, which leads to false-positive duplicates and degraded user experience. To address this, we propose **FITRep**, a white-box, attention-guided framework inspired by FIT that enables fine-grained, interpretable item representation through concept-hierarchical extraction, structure-preserving dimensionality reduction, and scalable indexing. Deployed on Meituan's advertising system, **FITRep** significantly improves user engagement, yielding +3.60% CTR and +4.25% CPM in online A/B tests, demonstrating its practical value in large-scale recommendation systems.

## References

[1] Chuhan Wu, Fangzhao Wu, Tao Qi, Chao Zhang, Yongfeng Huang, and Tong Xu. 2022. Mm-rec: Visiolinguistic model empowered multimodal news recommendation. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval.* 2560–2564.

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923* (2025).

[3] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[4] Yuyang Ye, Zhi Zheng, Yishan Shen, Tianshu Wang, Hengruo Zhang, Peijun Zhu, Runlong Yu, Kai Zhang, and Hui Xiong. 2025. Harnessing multimodal large language models for multimodal sequential recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. 13069–13077.

[5] Peilin Zhou, Chao Liu, Jing Ren, Xinfeng Zhou, Yueqi Xie, Meng Cao, Zhongtao Rao, You-Liang Huang, Dading Chong, Junling Liu, et al. 2025. When Large Vision Language Models Meet Multimodal Sequential Recommendation: An Empirical Study. In *Proceedings of the ACM on Web Conference 2025.* 275–292.

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers).* 4171–4186.

[7] Chao Zhang, Haoxin Zhang, Shiwei Wu, Di Wu, Tong Xu, Xiangyu Zhao, Yan Gao, Yao Hu, and Enhong Chen. 2025. Notellm-2: Multimodal large representation models for recommendation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 1.* 2815–2826.

[8] Anne M Treisman and Garry Gelade. 1980. A feature-integration theory of attention. *Cognitive psychology* 12, 1 (1980), 97–136.

[9] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards General Text Embeddings with Multi-stage Contrastive Learning. arXiv:2308.03281 [cs.CL]

[10] Tim Sainburg, Leland McInnes, and Timothy Q Gentner. 2021. Parametric UMAP embeddings for representation and semisupervised learning. *Neural Computation* 33, 11 (2021), 2881–2907.

[11] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The Faiss library. (2024). arXiv:2401.08281 [cs.LG]