

# Advancing Re-Ranking with Multimodal Fusion and Target-Oriented Auxiliary Tasks in E-Commerce Search

Enqiang Xu  
JD.com, Inc.  
Beijing, China  
xuenqiang@jd.com

Xinhui Li  
JD.com, Inc.  
Beijing, China  
lixinhui9@jd.com

Zhigong Zhou  
JD.com, Inc.  
Beijing, China  
zhouzhigong1@jd.com

Jiahao Ji\*  
JD.com, Inc.  
Beijing, China  
jiahaoji@buaa.edu.cn

Jinyuan Zhao†  
JD.com, Inc.  
Beijing, China  
zhaojinyuan1@jd.com

Dadong Miao  
JD.com, Inc.  
Beijing, China  
miaodadong@jd.com

Songlin Wang  
JD.com, Inc.  
Beijing, China  
wangsonglin3@jd.com

Lin Liu  
JD.com, Inc.  
Beijing, China  
liulin1@jd.com

Sulong Xu  
JD.com, Inc.  
Beijing, China  
xusulong@jd.com

## ABSTRACT

In the rapidly evolving field of e-commerce, the effectiveness of search re-ranking models is crucial for enhancing user experience and driving conversion rates. Despite significant advancements in feature representation and model architecture, the integration of multimodal information remains underexplored. This study addresses this gap by investigating the computation and fusion of textual and visual information in the context of re-ranking. We propose **Advancing Re-Ranking with Multimodal Fusion and Target-Oriented Auxiliary Tasks (ARMMT)**, which integrates an attention-based multimodal fusion technique and an auxiliary ranking-aligned task to enhance item representation and improve targeting capabilities. This method not only enriches the understanding of product attributes but also enables more precise and personalized recommendations. Experimental evaluations on JD.com's search platform demonstrate that ARMMT achieves state-of-the-art performance in multimodal information integration, evidenced by a 0.22% increase in the Conversion Rate (CVR), significantly contributing to Gross Merchandise Volume (GMV). This pioneering approach has the potential to revolutionize e-commerce re-ranking, leading to elevated user satisfaction and business growth.

## CCS CONCEPTS

• **Information systems** → Information retrieval.

\*Corresponding author

†Corresponding author

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](mailto:permissions@acm.org).

CIKM '24, October 21–25, 2024, Boise, ID, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0436-9/24/10

<https://doi.org/10.1145/3627673.3680063>

## KEYWORDS

Multimodal Fusion, Neural Network, Information Retrieval

### ACM Reference Format:

Enqiang Xu, Xinhui Li, Zhigong Zhou, Jiahao Ji, Jinyuan Zhao, Dadong Miao, Songlin Wang, Lin Liu, and Sulong Xu. 2024. Advancing Re-Ranking with Multimodal Fusion and Target-Oriented Auxiliary Tasks in E-Commerce Search. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management (CIKM '24)*, October 21–25, 2024, Boise, ID, USA. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3627673.3680063>

## 1 INTRODUCTION

In the dynamic landscape of e-commerce, search ranking models are fundamental to enhancing user experience and optimizing conversion rates. These models play a crucial role in organizing the vast array of products that match a user's query, significantly impacting user satisfaction. Leading e-commerce platforms typically employ a cascade structure, which includes stages such as matching, pre-ranking, ranking, and re-ranking. The re-ranking stage, as the final step before a user's purchase decision, focuses on refining the scores of the top-k items by utilizing detailed user and item-specific information, thus providing a more personalized and relevant product list.

Traditional re-ranking methods, often based on deep learning, primarily focus on model architecture and enhanced feature representation. These methods can be categorized into two types [11]. The first type is step-greedy re-ranking strategies [11, 13, 36], which determine the display result for each position sequentially, considering only the information of the previous item. This approach often falls short of the optimal outcome because it neglects the information of subsequent items. In contrast, contextual re-ranking strategies [6, 25, 30] capture the interdependencies among items using a contextual evaluation model, refining the click-through rate predictions for each item. For instance, PRM [25] takes the initial ranking list as input and generates the optimal permutation based on the contextual model's predictions. Most existing re-ranking

frameworks build upon this by adopting a two-stage architecture [6, 30], involving item sequence generation and ranking evaluation. The sequence generation task [16] refines the initial ranking results based on user queries and personalized information, while the ranking evaluation [28] ensures the quality of the generated results to meet user needs and provide a satisfactory user experience.

These re-ranking models heavily rely on unique IDs and categorical features for user-item matching [16, 34]. However, these methods primarily deal with sparse ID features and may not be adequately trained when the IDs appear infrequently in the data. In contrast, images provide intrinsic visual descriptions that can enhance model generalization. Given that users directly interact with item images, these images can offer additional visual information about user interests. Previous innovative works have introduced image features into recommendation systems [9, 23], focusing on representing ads with image features in click-through rate prediction. For instance, AMS [12] explored a new method based on a visual model to analyze user behavior using advanced model servers. However, most of these works concentrate on the acquisition of image features, with less research on the integration of multimodal features in ranking models. While progress has been made in extracting multimodal features [27], there remains a gap in research on effectively incorporating these different types of features into re-ranking models.

This study bridges the gap by investigating the computation and fusion of multimodal information within the realm of re-ranking models, using JD.com’s re-ranking system as a case study. The inclusion of multimodal cues, encompassing both textual and visual information, aims to mitigate the limitations of traditional textual and ID-based features, which lack the richness of visual information. To achieve this, we developed several key advancements. Firstly, we introduce an attention-based fusion mechanism, specifically the Context-Aware Fusion Unit (CAFU), which synergistically integrates textual and visual information. This method is designed to enhance item representation by incorporating visual information which might be neglected by traditional methods. Secondly, we propose the use of Multi-Perspective Self-Attention, which integrates multimodal information with other fields. This mechanism enhances the model’s ability to capture complex interactions between different types of information, thereby improving the precision of the re-ranking process. Thirdly, to maximize the utility of multimodal information, we introduce an auxiliary task aligned with the ranking objective, providing additional supervision for the learning of multimodal representations. The integration of multimodal information into re-ranking models not only enhances our understanding of user preferences but also improves the precision and personalization of search results.

Our contributions can be summarized as follows:

**1. Attention-based Fusion Mechanism:** We introduce the Context-Aware Fusion Unit (CAFU) and Multi-Perspective Self-Attention, which integrate textual and visual information for re-ranking models. This approach addresses the limitations of traditional textual and ID-based features by incorporating rich visual cues, enhancing item representations.

**2. Auxiliary Task for Supervision:** We design an auxiliary task aligned with the ranking objective to supervise the learning

process of multimodal representations, improving their quality and relevance for the ranking task.

**3. Empirical Validation:** We validate our approach through rigorous experimentation. Our results demonstrate the effectiveness of the CAFU, Multi-Perspective Self-Attention, and the auxiliary task in enhancing re-ranking models.

These contributions are significant for e-commerce. They introduce novel methods for optimizing re-ranking models, which have been successfully deployed, potentially improving user satisfaction, conversion rates, and the shopping experience.

## 2 RELATED WORK

In this section, we provide a concise overview of the latest developments in ranking models, with a focus on re-ranking algorithms and multimodal fusion within ranking frameworks. These areas are highly pertinent to the scope of our research.

### 2.1 Re-ranking Models

In the dynamic and competitive landscape of e-commerce, ranking models play a pivotal role in curating and presenting items in a manner that aligns with user preferences and needs. Traditional ranking models [4, 9, 14, 18, 22, 26, 35] primarily focus on scoring individual items for click-through rate (CTR) estimation, aiming to optimize performance at a single point. Re-ranking models distinguish themselves from traditional ranking models by their unique ability to model the contextual relationships within a sequence of candidate items. These models are typically categorized into two distinct approaches: step-greedy strategies and context-wise strategies [11].

**Step-greedy Re-ranking Strategies:** Step-greedy approaches utilize a sequential decision-making process for each position in the display results. These methods often employ recurrent neural networks or approximation solutions to determine the order of items. For instance, DPP [7] identifies the most relevant and diverse subset of candidates by calculating the determinant of a kernel matrix. In contrast, MMR [5] relies on pairwise similarity. Seq2Slate [2] uses a pointer network, while MIRNN [1] employs a gated recurrent unit to sequentially establish the order of items. However, this category of methods tends to overlook subsequent information in the ranking sequence, often resulting in suboptimal outcomes.

**Context-wise Re-ranking Strategies:** Context-wise methods aim to capture the mutual influence among items by using evaluation models that reassess the CTR or conversion rate (CVR) for each item. Methods like PRM [25] and DLCM [1] process the initial ranking list with RNNs or self-attention mechanisms to model context signals and predict values for each item. To circumvent the evaluation-before-reranking dilemma [30], some researchers adopt a two-stage re-ranking framework, comprising permutation generation followed by permutation evaluation, such as GRN [11] and PRS [10]. PIER [29] employs a fine-grained permutation selection module to choose the top-K candidates from the entire permutation space, along with a context-aware prediction module that predicts the list-wise CTR for each item.

In our approach, we adopt the two-stage architecture that emphasizes the seamless integration of multimodal information within a

framework that captures reciprocal contextual interactions, thereby facilitating the identification of the most optimal sequence.

## 2.2 Multimodal Fusion

In the context of e-commerce recommendations, visual signals play a significant role as intuitive factors influencing user purchase decisions. Enhancing item representation with comprehensive visual information can enrich the description of the current item and elevate the model’s level of personalization. Multimodal Recommender Systems [19] stand out from ID-based recommendation models due to their superior generalization capabilities [33]. Additionally, their proficiency in processing information across diverse modalities makes them particularly advantageous for multimedia services. Among the various feature fusion techniques, the attention mechanism is notably effective [19], significantly enhancing the system’s ability to understand and cater to user preferences. For instance, UVCAN [20] utilizes user-side ID features to generate fusion weights for item-side multimodal information through self-attention, thereby facilitating micro-video recommendations.

In the realm of multimodal information fusion, several studies have made significant advancements. During its pre-training phase, MCPTR [21] leverages self-supervised multimodal contrastive learning to acquire fusion weights across different modalities, simultaneously deriving multimodal user and item representations. CMBF [8] adopts a cross-modal fusion approach to comprehensively integrate multimodal features, learning the interplay between different modalities. MML [24] and MARIO [17] employ an attention network to assess the impact of each modality on the interactions between users and items, preserving modality-specific attributes to obtain personalized embeddings for items relative to users. VLSNR [15] initially processes images and titles through a CLIP [27] encoder, followed by a series of attention layers to derive multimodal representations of news, and ultimately employs a GRU network to learn users’ temporal interests.

Building on these studies, our approach focuses on the application of multimodal information fusion in ranking models. By integrating signals from different modalities, we aim to enhance the representation of items.

## 3 PRELIMINARIES

In this section, we initiate our discussion with the formal definitions of symbols that delineate the foundational framework of our re-ranking model and its associated tasks. Subsequently, we elucidate on the most prevalent model architectures in re-ranking, such as ID-based Deep Interest Network (DIN [35]) and the context-based Transformer structure, which form the cornerstone of what we term as the base model.

### 3.1 Background

In the architecture of conventional industrial search engine systems, a modular approach is typically employed, encompassing stages such as recall, pre-ranking, ranking, and re-ranking. Notably, the re-ranking phase serves to refine the output of the ranking module, optimizing the final presentation of search results to the end-user.

Mathematically, given a user  $U$  entering a query  $Q$  into the search bar, a ranking list of candidate items  $I = \{item_i\}_{i=1}^N$  is generated,

where  $N$  represents the top  $N$  candidate items output by the ranking module. The task of the re-ranking stage is to learn a strategy  $\mathcal{F}$  such that  $I^* = \mathcal{F}(I, U, Q)$ , which selects and rearranges items from  $I$ , presenting a final ranking list  $I^*$  to users with the aim of maximizing the conversion rate (CVR).

## 3.2 Base Model

**3.2.1 ID-based Deep Interest Network.** To capture user preferences from historical interactions, we utilize a target attention mechanism, specifically the Deep Interest Network (DIN). For each candidate item, the query is formed by concatenating the representations of the item ID, shop ID, and brand ID. Similarly, the key, representing the user’s historical behavior sequence, is formed in the same way. By applying the attention mechanism, we derive the personalized item representation  $P_{id}$  for each candidate item.

**3.2.2 Context-based Transformer Encoder.** We concatenate  $P_{id}$  with item features to create the representation of candidate items. To capture the interactions among these candidate items, specifically as inspired by the Personalized Ranking Model (PRM), we linearly transform this representation to obtain  $Q$ ,  $K$ , and  $V$ , and then utilize a transformer encoder to compute  $A$ , as shown in Equation 1.

$$A = \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where,  $d_k$  is the dimensionality of the key vectors  $K$ .

Subsequently, we pass  $A$  through a fully connected layer to reduce its dimensionality. To effectively rank the candidate items, we leverage the listwise modeling capability by employing the softmax activation function, which outputs the predicted scores  $\hat{y}$  for each item. We minimize the cross-entropy loss  $\mathcal{L}$  between the conversion labels  $y$  (binary, 0 or 1) and the predicted scores  $\hat{y}$ , as shown in Equation 2.

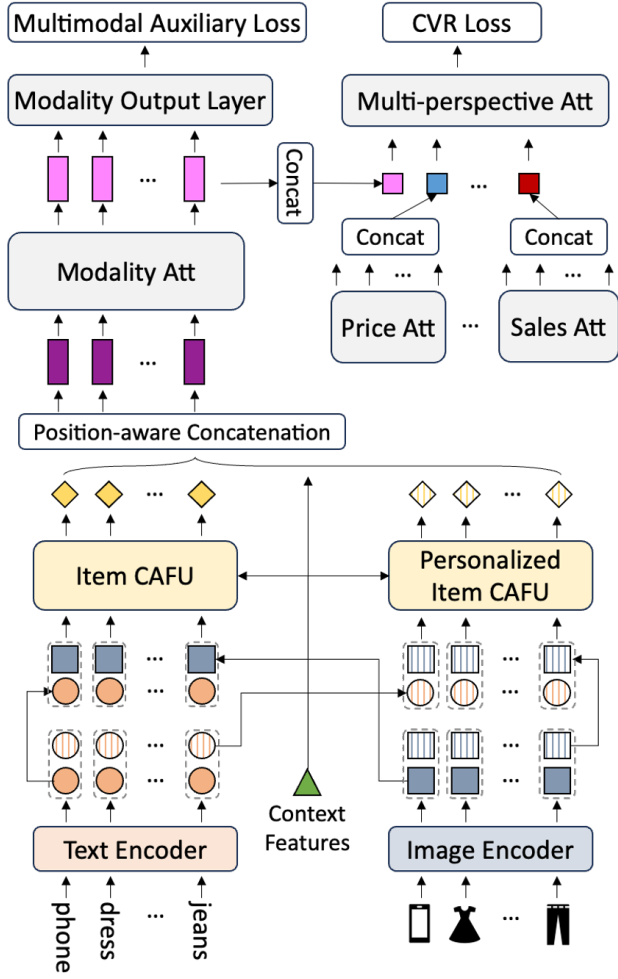
$$\mathcal{L} = - \sum_{i=1}^N y_i \log(\hat{y}_i) \quad (2)$$

## 4 METHOD

In this section, we present the overall architecture of our proposed ARMMT framework, as illustrated in Figure 1. We will explore its key components: the generation of multimodal representations for items and personalized items (derived from the interaction between items and user historical interests), the hierarchical multimodal fusion process using the Context-Aware Fusion Unit (CAFU), the integration of multimodal information with other domains through Multi-Perspective Self-Attention, and the critical multimodal auxiliary losses that enhance our model’s performance.

### 4.1 Multimodal representations

Figure 2 illustrates the encoding processes for textual and visual features in our method, with the left and right parts depicting these processes, respectively. In the context of our sequence analysis, we derive pre-trained embeddings for items by indexing the vocabulary. Subsequently, we employ the target item as a query to apply

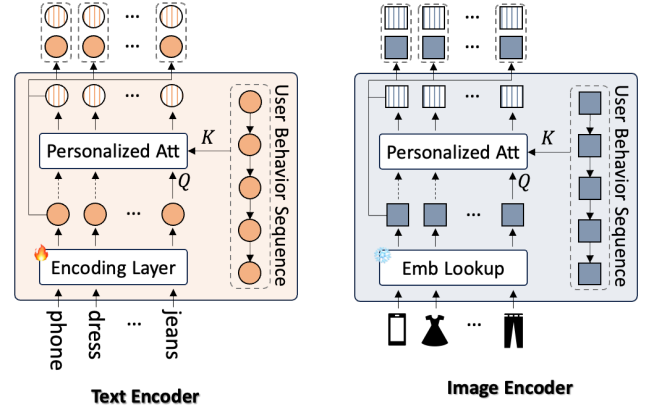


**Figure 1: The framework of Advancing Re-Ranking with Multimodal Fusion and Target-Oriented Auxiliary Tasks (AR-MMT).**

target attention mechanisms, which effectively model the interactions between the target item and the other items in the sequence, culminating in the derivation of the final feature embeddings.

**4.1.1 Multimodal Representation of Item.** To optimize training efficiency and avoid increasing online serving pressure, we adopted a two-stage integration strategy for multimodal representations. In contrast, previous methods often fuse image and text modalities too early during the pre-training stage, which can lead to inconsistencies with user preferences and suboptimal search results. For example, when users search for dresses, they may prioritize different attributes such as color or length, which early fusion methods struggle to dynamically adjust to.

To address this, we first obtain separate representations for each modality. Specifically, for the image modality, we use frozen image embeddings,  $I_{img}$ , extracted by fine-tuning the ResNet [32] convolutional network on JD.com product images after object detection with YOLO4 [3]. For the text modality, we use the



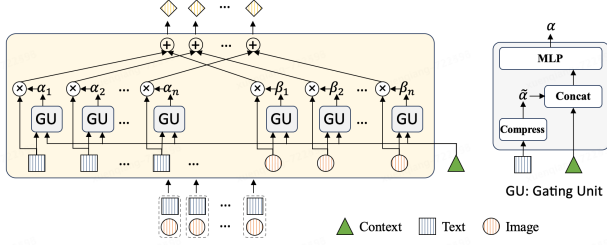
**Figure 2: The encoding process of textual and image information. Effective information from user behavior sequences is extracted through multi-head attention.**

title features directly, representing them as  $I_{text}$ . This separation allows us to maintain the integrity of each modality’s features before they are combined in the ranking model.

**4.1.2 Multimodal Representation of Personalized Item.** We implemented a two-stage behavioral sequence modeling approach to accurately capture the influence of user preferences on purchasing decisions, particularly within similar product categories. For instance, a user’s preference for the appearance of clothing significantly influences their decisions when purchasing similar apparel but has minimal impact on unrelated categories, such as computers. In the first stage, we use the category of the query, predicted by a large language model (LLM), to filter the user’s historical behavior sequence  $H$ , retaining only those behaviors that match the query’s category, resulting in a relevant sub-sequence  $H^*$ . In the second stage, we use a target attention mechanism to model interactions between these behaviors and the target item.

To capture the fine-grained features of both modalities, we model personalized item representations separately for text and images. For the image modality, we first leverage the image representations of items and those from the user’s historical behaviors. Building on this foundation, we incorporate user features such as age and gender to enhance generalization across similar user groups. Furthermore, each historical behavior is detailed, including the type of interaction (e.g., click, order), frequency of actions, and the recency of the actions. These features are integral in calculating the attention weights, which are then used to create personalized item image representations  $P_{img}$ , thereby revealing the user’s latent preferences in decision-making. As illustrated in Figure 2 (right side), the process of generating personalized product image representations involves obtaining image features through embedding lookup and integrating these features with the user’s behavioral sequence using a personalized attention mechanism.

Similarly, for the text modality, we begin by extracting text representations from item titles, attributes, and other textual features using an encoding layer. These representations are generated for



**Figure 3: The diagram of the Context-Aware Fusion UNIT. In this diagram, triangles, rectangles, and circles represent context, text, and image features, respectively.**

both the candidate items and the items in the user’s historical behavior sequence. To further personalize these text representations  $P_{\text{text}}$ , we incorporate user demographics such as age and gender, as well as detailed historical behavior information. This includes the type of interactions (e.g., clicks, orders), the frequency of these actions, and their recency. This information is crucial for calculating the attention weights of items within the user’s behavioral sequence. Figure 2 (left side) illustrates this process, highlighting how textual features are encoded and integrated with the user’s behavioral sequence through a personalized attention mechanism.

## 4.2 Hierarchical Multimodal Fusion

**4.2.1 Context-Aware Fusion UNIT.** We introduce a Context-Aware Fusion Unit (CAFU) designed to seamlessly integrate the representations of items and personalized items from both image and text modalities. The input to this module consists of a list of features  $X \in \mathbb{R}^{N \times D \times M}$ , along with contextual information  $C \in \mathbb{R}^J$ , such as query features and user features, where  $N$  represents the number of item candidates,  $D$  denotes the item feature dimension,  $M$  indicates the number of modalities (e.g., image and text), and  $J$  indicates the contextual feature dimension.

Initially, we compress the representations using mean pooling to obtain an initial weighting vector  $z_n = [z_{n1}, z_{n2}, \dots, z_{nM}] \in \mathbb{R}^M$  for each item  $x_n \in \mathbb{R}^{D \times M}$ , as shown in Equation 3.

$$z_{nm} = \frac{1}{D} \sum_{d=1}^D x_{ndm} \quad (3)$$

Next, we concatenate the initial weights  $z_n$  with the contextual information  $C$ , thereby incorporating the context into our weighting scheme, as shown in Equation 4.

$$z_n^c = \text{concat}(z_n, C) \in \mathbb{R}^{M+J} \quad (4)$$

To calculate the weights for each modal representation, we employ a Multi-Layer Perceptron (MLP). Specifically, the weights  $s_n$  are computed as shown in Equation 5, where  $\alpha_n$  and  $\beta_n$  in Figure 3 represent the weights for different modalities to facilitate understanding.

$$s_n = \text{softmax}(W_2 \delta(W_1 z_n^c)) \in \mathbb{R}^M \quad (5)$$

where  $\delta$  denotes the ReLU activation function,  $W_1 \in \mathbb{R}^{\frac{M+J}{r} \times (M+J)}$  and  $W_2 \in \mathbb{R}^{M \times \frac{M+J}{r}}$  are the weight matrices, and  $r$  is the reduction ratio that controls the size of the hidden layer.

In the final stage, these weights are multiplied by their corresponding representations, and a sum pooling operation is performed to obtain the fused representation  $B \in \mathbb{R}^{N \times D}$ , which serves as the output of this module, as shown in Equation 6 and illustrated in Figure 3.

$$b_n = \sum_{m=1}^M s_{nm} \cdot x_{nm} \quad (6)$$

This method offers a degree of interpretability, allowing for easy extraction and visual analysis of the fusion weights between different modalities. Moreover, this approach demonstrates a high degree of robustness to missing values. Utilizing this methodology, we successfully fuse the image and text representations of items into a cohesive multimodal item representation  $I_{\text{mo}}$ , as shown in Equation 7.

$$I_{\text{mo}} = \text{CAFU}([I_{\text{img}}, I_{\text{text}}], C) \quad (7)$$

Similarly, a personalized multimodal item representation  $P_{\text{mo}}$  is derived through the same process, as shown in Equation 8.

$$P_{\text{mo}} = \text{CAFU}([P_{\text{img}}, P_{\text{text}}], C) \quad (8)$$

**4.2.2 Multi-Perspective Self-Attention.** Building on the integrated representations achieved through the CAFU, we apply a Multi-Perspective Self-Attention mechanism to deeply merge the multimodal field with other critical fields, such as price and sales, thereby achieving a comprehensive global feature fusion. Specifically, we combine the multimodal representations of items with personalized multimodal item representations and contextual information through concatenation, which is then processed via an MLP to obtain a unified multimodal representation, denoted as  $M_{\text{mo}}$ , as shown in Equation 9.

$$M_{\text{mo}} = \text{MLP}(\text{concat}(I_{\text{mo}}, P_{\text{mo}}, C)) \quad (9)$$

In conventional approaches, multimodal representations  $M_{\text{mo}}$  are typically incorporated directly into the input layer, where they are combined with other features through concatenation. However, in re-ranking models, there is a risk that during self-attention computations, the model may favor strong fields, thereby overshadowing the multimodal information. To address this, we explicitly model the context by considering the interactions and variations among the multimodal information of items to be ranked. Specifically, we employ a transformer encoder architecture to capture the interactive representation, resulting in an enriched  $A_{\text{mo}}$ , as shown in Equation 10.

$$A_{\text{mo}} = \text{Attention}(Q_{\text{mo}}, K_{\text{mo}}, V_{\text{mo}}) \quad (10)$$

where  $Q_{\text{mo}}$ ,  $K_{\text{mo}}$ , and  $V_{\text{mo}}$  represent the query, key, and value, respectively, all of which are linearly transformed from  $M_{\text{mo}}$ .

We observe a similar pattern in the main task, where users vary in their sensitivity to details such as price, promotions, sales, and quality. To address this, we process each field individually through self-attention mechanisms, yielding representations such as  $A_{\text{price}}$

and  $A_{\text{sales}}$ . We then merge these fields with the multimodal field, resulting in a combined representation denoted as  $M_{\text{main}}$ , as shown in Equation 11.

$$M_{\text{main}} = \text{MLP}(\text{concat}(A_{\text{mo}}, A_{\text{price}}, A_{\text{sales}})) \quad (11)$$

Subsequently, we apply global attention on the combined representation to extract higher-order semantic features. This process helps to better align with user behavior patterns, as shown in Equation 12.

$$A_{\text{main}} = \text{Attention}(Q_{\text{main}}, K_{\text{main}}, V_{\text{main}}) \quad (12)$$

where,  $Q_{\text{main}}$ ,  $K_{\text{main}}$ , and  $V_{\text{main}}$  represent the query, key, and value, respectively, all of which are linearly transformed from  $M_{\text{main}}$ .

Finally, we score each item in the candidate list using the softmax function to obtain the model’s predicted scores  $\hat{y}$ , as shown in Equation 13.

$$\hat{y} = \text{softmax}(\text{MLP}(A_{\text{main}})) \quad (13)$$

As illustrated in Figure 1, the price attention ( $A_{\text{price}}$ ), sales attention ( $A_{\text{sales}}$ ), and multimodal attention ( $A_{\text{mo}}$ ) are concatenated and then processed through Multi-Perspective Self-Attention to produce the final representation, which is used to compute the CVR loss.

### 4.3 Multimodal Auxiliary Tasks

In our study, we use supervised learning specifically designed for multimodal auxiliary tasks due to the significant impact of multimodal information on user click behavior. We incorporate click labels  $y^{ctr}$  as direct user feedback. Each item in the candidate list is scored using a sigmoid function to predict the click probability  $\hat{y}_i^{ctr}$ , as shown in Equation 14.

$$\hat{y}_i^{ctr} = \text{sigmoid}(\text{MLP}(A_{\text{mo}})) \quad (14)$$

The auxiliary task loss  $\mathcal{L}_{\text{aux}}$  is computed using cross-entropy, as shown in Equation 15.

$$\mathcal{L}_{\text{aux}} = -\frac{1}{N} \sum_{i=1}^N y_i^{ctr} \log \hat{y}_i^{ctr} + (1 - y_i^{ctr}) \log(1 - \hat{y}_i^{ctr}) \quad (15)$$

### 4.4 Model Training

The main task loss  $\mathcal{L}_{\text{main}}$ , which is used to compute the CVR loss, is the cross-entropy between conversion labels  $y$  and predicted scores  $\hat{y}$  (Equation 13), as shown in Equation 16.

$$\mathcal{L}_{\text{main}} = -\sum_{i=1}^N y_i \log \hat{y}_i \quad (16)$$

The final optimization objective combines both the main task loss and the auxiliary task loss, with  $\lambda$  as a hyper-parameter, as shown in Equation 17.

$$\mathcal{L} = \mathcal{L}_{\text{main}} + \lambda \mathcal{L}_{\text{aux}} \quad (17)$$

## 5 EXPERIMENTS

### 5.1 Dataset and Evaluation Metric

In this paper, we conduct our experiments on an internal dataset provided by JD.com. This dataset is extracted from user behavioral logs specifically related to search activities and includes only those user sessions that resulted in a conversion. The dataset is organized by user search sessions, with each row representing a different session. Each session contains 30 item samples along with their corresponding labels indicating clicks and conversions. For the training set, data was collected over a consecutive 21-day period, resulting in a substantial dataset comprising hundreds of millions of sessions and billions of products. The data from the 22nd day was reserved for the test set, allowing us to evaluate the model’s performance on unseen data.

For offline evaluation, we predict the likelihood of a user’s conversion for each item and employ the widely used AUC (Area Under the ROC Curve) as our performance metric. This metric is well-established in search systems, advertising, and recommendation systems, as supported by the references in [9, 23], and represents the cumulative area under the characteristic curve, effectively quantifying the predictive accuracy of our model.

### 5.2 Experimental Settings

**5.2.1 Baseline Models.** To ensure a fair comparison, we selected some of the most representative models from the industry, such as PRM [25] and PIER [29], as benchmarks for our model. Our model is an optimization directly based on the PIER model. Specifically:

- **PRM:** This stands for the context-sensitive sequence awareness model, which is a commonly used re-ranking model in the industry.
- **PIER:** This model utilizes a two-stage re-ranking architecture comprising a Sequence Generator and a Sequence Evaluator. It serves as our baseline model that has been deployed online.

**5.2.2 Implementation Details.** All our experiments, including ARMMT and the comparative baseline experiments, were implemented using TensorFlow 1.15 and run on multiple NVIDIA V100 GPUs with CUDA 10. All experiments were conducted using the same offline dataset introduced in Section 5.1.

For the ARMMT method, we used the AdaGrad optimizer with a learning rate of 0.07 and trained the model for 20 epochs with a batch size of 128. The pre-trained ID and image vocabulary contained 8 million items, with an embedding dimension set to 32. The Transformer encoder consisted of two layers, each employing multi-head self-attention with 6 heads, and each head having a dimension of 32. The weight of the auxiliary loss,  $\lambda$ , was set to 1. This configuration was chosen to facilitate efficient convergence and ensure effective learning from the training data.

### 5.3 Performance Comparison

In this subsection, we conduct a comparative study of our proposed ARMMT model, which is benchmarked against PIER [29] and PRM [25], as presented in Table 1. PIER, characterized by its two-stage architecture and the introduction of an evaluator for scoring sequence generation, shows a notable enhancement in the AUC metric compared to the traditional PRM. Extending PIER, our

ARMMT model incorporates visual cues into the re-ranking process, achieving alignment between the modal information of the model’s inputs and the user’s browsing behavior. This enhancement yields the highest performance, with an offline AUC score of 0.9647, marking a 0.0005 increase over PIER. Given that the re-ranking module focuses on scoring only the top 30 items from the ranking model, this AUC improvement of 0.0005 is statistically meaningful, thereby validating the efficacy of our ARMMT model.

**Table 1: Performance comparison between the baseline models and our proposed ARMMT method. Bold indicates the best result, underlined indicates the second-best, and the AUC gain shows the improvement achieved by ARMMT over the second-best result.**

Methods	AUC	AUC Gain
PRM [25]	0.9636	-
PIER [29]	<u>0.9642</u>	-
ARMMT	<b>0.9647</b>	<b>0.0005</b>

## 5.4 Ablation Studies

In this subsection, we perform an extensive ablation analysis to quantitatively assess the impact of various architectural components on the model’s performance. The analysis includes the following experimental setups: (1) **PIER [29]**: A benchmark model relying solely on textual and ID-embedding features, excluding image embeddings. (2) **w/o CAFU and Auxiliary Tasks**: A model that directly incorporates image embeddings, bypassing the Context-Aware Fusion Unit and auxiliary losses. (3) **w/o Auxiliary Tasks**: A model using CAFU for multimodal fusion, excluding auxiliary losses. (4) **ARMMT**: The complete framework integrating both CAFU and auxiliary tasks.

**Table 2: Ablation study of ARMMT.**

Methods	AUC	AUC gain
PIER [29]	0.9642	-
w/o CAFU and Auxiliary Tasks	0.9643	0.0001
w/o Auxiliary Tasks	0.9645	0.0003
ARMMT	<b>0.9647</b>	<b>0.0005</b>

As shown in Table 2, our ablation experiments on the PIER [29] baseline model yielded varying degrees of improvement in the AUC metric. Directly incorporating image embeddings into the re-ranking model resulted in a 0.0001 increase in AUC compared to PIER [29], validating that incorporating image embeddings can enhance the model’s predictive accuracy. The introduction of CAFU and auxiliary tasks brought about additional AUC improvements of 0.0002 and 0.0002, respectively. These results suggest that the strategic design of CAFU and the incorporation of auxiliary tasks

effectively enhance multimodal feature fusion, thereby improving prediction precision. Our comprehensive ARMMT model achieves the best performance, with an offline AUC of 0.9647.

## 5.5 Online Results

On JD.com’s platform, which boasts tens of millions of daily active users, we conducted a 7-day A/B test, allocating 10% of the user traffic to each experimental variant. The control variant (Variant A) used the PIER [29] model, whereas the experimental variant (Variant B) implemented the novel ARMMT model proposed in our research.

As indicated in Table 3, the ARMMT model demonstrated superior performance across all pivotal metrics. Notably, it yielded a statistically significant 0.22% rise in CVR (p-value = 0.04) [31], alongside a 0.03% uplift in CTR and a substantial 0.49% increase in GMV. These improvements have led to a marked escalation in business revenue. The experimental results are consistent with our prior offline analyses, confirming the ARMMT model’s effectiveness in integrating visual information into the re-ranking process. This seamless integration significantly improves sorting efficiency, ensures a more precise alignment with user preferences, and notably enhances the platform’s conversion rate.

In light of these positive outcomes, the ARMMT model was deployed on the platform in early 2024, providing substantial benefits to hundreds of millions of users.

**Table 3: Online A/B test of ARMMT. The improvements are averaged over 7 days in April 2024.**

Online Metrics	Improvement	p-value
CVR	<b>0.22%</b>	<b>0.043</b>
GMV	0.49%	0.121
CTR	0.03%	0.668

## 6 CONCLUSION

This study introduces an innovative attention-based multimodal fusion method for re-ranking in e-commerce search. By integrating textual and visual information, it overcomes the limitations of traditional re-ranking models that rely solely on unimodal data, providing a more intuitive representation of the content users interact with. Additionally, an auxiliary task designed to predict click-through rates aligns closely with the ranking task, ensuring the effective utilization of multimodal information. Empirical results demonstrate that these combined approaches enhance user satisfaction and conversion rates, highlighting significant practical implications for the e-commerce industry and providing a competitive edge in the market. Furthermore, this research opens new avenues for incorporating additional modalities and dynamic ranking objectives, advancing re-ranking technologies and offering practical guidance for e-commerce platforms to enhance personalized shopping experiences.

## REFERENCES

- [1] Qingyao Ai, Keping Bi, Jiafeng Guo, and W Bruce Croft. 2018. Learning a deep listwise context model for ranking refinement. In *The 41st international ACM SIGIR conference on research & development in information retrieval*. 135–144.
- [2] Irwan Bello, Sayali Kulkarni, Sagar Jain, Craig Boutilier, Ed Chi, Elad Eban, Xiyang Luo, Alan Mackey, and Ofer Meshi. 2018. Seq2Slate: Re-ranking and slate optimization with RNNs. *arXiv preprint arXiv:1810.02019* (2018).
- [3] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020).
- [4] Yue Cao, Xiaojiang Zhou, Jiaqi Feng, Peihao Huang, Yao Xiao, Dayao Chen, and Sheng Chen. 2022. Sampling is all you need on modeling long-term user behaviors for CTR prediction. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2974–2983.
- [5] Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*. 335–336.
- [6] Chi Chen, Hui Chen, Kangzhi Zhao, Junsheng Zhou, Li He, Hongbo Deng, Jian Xu, Bo Zheng, Yong Zhang, and Chunxiao Xing. 2022. Extr: click-through rate prediction with externalities in e-commerce sponsored search. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2732–2740.
- [7] Laming Chen, Guoxin Zhang, and Eric Zhou. 2018. Fast greedy map inference for determinantal point process to improve recommendation diversity. *Advances in Neural Information Processing Systems* 31 (2018).
- [8] Xi Chen, Yangsiyi Lu, Yuehai Wang, and Jianyi Yang. 2021. CMBF: Cross-modal-based fusion recommendation algorithm. *Sensors* 21, 16 (2021), 5275.
- [9] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Isipir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [10] Yufei Feng, Yu Gong, Fei Sun, Junfeng Ge, and Wenwu Ou. 2021. Revisit recommender system in the permutation prospective. *arXiv preprint arXiv:2102.12057* (2021).
- [11] Yufei Feng, Binbin Hu, Yu Gong, Fei Sun, Qingwen Liu, and Wenwu Ou. 2021. GRN: Generative Rerank Network for Context-wise Recommendation. *arXiv preprint arXiv:2104.00860* (2021).
- [12] Tiezheng Ge, Liqin Zhao, Guorui Zhou, Keyu Chen, Shuying Liu, Huimin Yi, Zelin Hu, Bochao Liu, Peng Sun, Haoyu Liu, et al. 2018. Image matters: Visually modeling user behaviors using advanced model server. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 2087–2095.
- [13] Xudong Gong, Qinlin Feng, Yuan Zhang, Jiangling Qin, Weijie Ding, Biao Li, Peng Jiang, and Kun Gai. 2022. Real-time short video recommendation on mobile devices. In *Proceedings of the 31st ACM international conference on information & knowledge management*. 3103–3112.
- [14] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 1725–1731.
- [15] Songhao Han, Wei Huang, and Xiaotian Luan. 2022. VLSNR: Vision-Linguistics Coordination Time Sequence-aware News Recommendation. *arXiv preprint arXiv:2210.02946* (2022).
- [16] Balázs Hidasi and Alexandros Karatzoglou. 2018. Recurrent neural networks with top-k gains for session-based recommendations. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 843–852.
- [17] Taeri Kim, Yeon-Chang Lee, Kijung Shin, and Sang-Wook Kim. 2022. MARIO: modality-aware attention and modality-preserving decoders for multimedia recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 993–1002.
- [18] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong Sun. 2018. xdeepfm: Combining explicit and implicit feature interactions for recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1754–1763.
- [19] Qidong Liu, Jiayi Hu, Yutian Xiao, Jingtong Gao, and Xiangyu Zhao. 2023. Multi-modal recommender systems: A survey. *arXiv preprint arXiv:2302.03883* (2023).
- [20] Shang Liu, Zhenzhong Chen, Hongyi Liu, and Xinghai Hu. 2019. User-video co-attention network for personalized micro-video recommendation. In *The World Wide Web Conference*. 3020–3026.
- [21] Zhuang Liu, Yunpu Ma, Matthias Schubert, Yuanxin Ouyang, and Zhang Xiong. 2022. Multi-modal contrastive pre-training for recommendation. In *Proceedings of the 2022 International Conference on Multimedia Retrieval*. 99–108.
- [22] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. 2018. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1930–1939.
- [23] Kaixiang Mo, Bo Liu, Lei Xiao, Yong Li, and Jie Jiang. 2015. Image Feature Learning for Cold Start Problem in Display Advertising. In *IJCAI International Joint Conference on Artificial Intelligence*, Vol. 2015. 3728.
- [24] Xingyu Pan, Yushuo Chen, Changxin Tian, Zihan Lin, Jinpeng Wang, He Hu, and Wayne Xin Zhao. 2022. Multimodal meta-learning for cold-start sequential recommendation. In *Proceedings of the 31st ACM international conference on information & knowledge management*. 3421–3430.
- [25] Changhua Pei, Yi Zhang, Yongfeng Zhang, Fei Sun, Xiao Lin, Hanxiao Sun, Jian Wu, Peng Jiang, Junfeng Ge, Wenwu Ou, et al. 2019. Personalized re-ranking for recommendation. In *Proceedings of the 13th ACM conference on recommender systems*. 3–11.
- [26] Qi Pi, Guorui Zhou, Yujing Zhang, Zhe Wang, Lejian Ren, Ying Fan, Xiaoqiang Zhu, and Kun Gai. 2020. Search-based user interest modeling with lifelong sequential behavior data for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2685–2692.
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [28] Tetsuya Sakai, Jin Young Kim, and Inho Kang. 2023. A versatile framework for evaluating ranked lists in terms of group fairness and relevance. *ACM Transactions on Information Systems* 42, 1 (2023), 1–36.
- [29] Xiaowen Shi, Fan Yang, Ze Wang, Xiaoxu Wu, Muzhi Guan, Guogang Liao, Wang Yongkang, Xingxing Wang, and Dong Wang. 2023. PIER: Permutation-Level Interest-Based End-to-End Re-ranking Framework in E-commerce. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4823–4831.
- [30] Yunjia Xi, Weiwen Liu, Xinyi Dai, Ruiming Tang, Weinan Zhang, Qing Liu, Xiuqiang He, and Yong Yu. 2021. Context-aware reranking with utility maximization for recommendation. *arXiv preprint arXiv:2110.09059* (2021).
- [31] Enqiang Xu, Yiming Qiu, Junyang Bai, Ping Zhang, Dadong Miao, Songlin Wang, Guoyu Tang, Lin Liu, and Mingming Li. 2024. Optimizing E-commerce Search: Toward a Generalizable and Rank-Consistent Pre-Ranking Model. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2875–2879.
- [32] Jing Xu, Yu Pan, Xinglin Pan, Steven Hoi, Zhang Yi, and Zenglin Xu. 2022. RegNet: self-regulated network for image classification. *IEEE Transactions on Neural Networks and Learning Systems* (2022).
- [33] Zheng Yuan, Fajie Yuan, Yu Song, Youhua Li, Junchen Fu, Fei Yang, Yunzhu Pan, and Yongxin Ni. 2023. Where to go next for recommender systems? id-vs. modality-based recommender models revisited. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2639–2649.
- [34] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5941–5948.
- [35] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1059–1068.
- [36] Tao Zhuang, Wenwu Ou, and Zhirong Wang. 2018. Globally optimized mutual influence aware ranking in e-commerce search. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*. 3725–3731.