

Re-M: Adapting Multi-Modal Large Language Models for Zero-Shot Cross-Modal Hybrid Retrieval and Reranking

Anonymous ACL submission

Abstract

Multi-Modal Large Language Models (MLLMs) demonstrate remarkable capabilities in vision-language understanding. Leveraging MLLMs to extract features for modern retrieval pipelines, including sparse retrieval, dense retrieval, and reranking, offers a promising direction that eliminates the need for expensive training data. In this paper, we investigate the feasibility of extracting high-quality sparse representations from MLLMs and propose a multi-perspective prompting method to enhance representational expressivity. Furthermore, we identify a significant performance disparity between image-to-text and text-to-image tasks during the reranking phase, indicating the necessity for distinct strategies. Building on these insights, we introduce Re-M, a two-stage zero-shot cross-modal retrieval framework. By integrating sparse-dense hybrid retrieval with asymmetric reranking, Re-M achieves performance that rivals or even surpasses retrieval-oriented dense and sparse baselines in zero-shot settings.

1 Introduction

Cross-modal retrieval (Wang et al., 2016) plays a pivotal role in diverse real-world applications, such as recommendation systems and search engines. With the advent of Multi-modal Large Language Models (MLLMs), researchers have begun adapting these models for retrieval tasks to leverage their superior comprehension capabilities. However, existing works (Faysse et al., 2025; Ma et al., 2024; Zhang et al., 2024) primarily utilize MLLMs as backbones for supervised contrastive fine-tuning, a process that necessitates a vast amount of high-quality labeled data. In contrast, MLLMs inherently possess extensive cross-modal knowledge derived from pre-training, achieving remarkable performance in understanding tasks like zero-shot image classification (Yan et al., 2024; Hong et al., 2025; Atabuzzaman et al., 2025). Consequently,

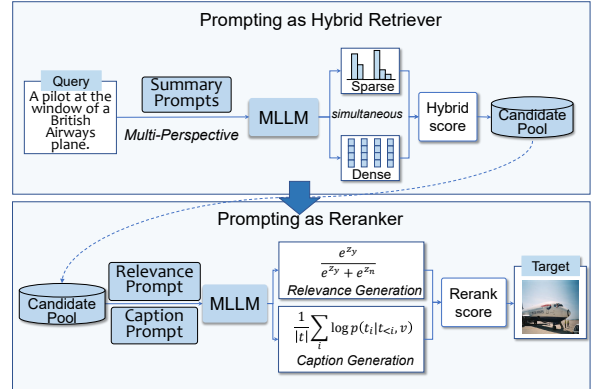


Figure 1: The proposed two-stage zero-shot retrieval framework, Re-M. The framework achieves hybrid retrieval with multi-perspectives prompting in the first stage and reranks the candidates via two different generation paradigms in the second stage.

harnessing the capabilities of MLLMs across various stages of modern retrieval pipelines represents a highly promising avenue for research.

Existing research on zero-shot retrieval has primarily focused on text modality, which demonstrate that prompting strategies enable LLMs to function as both retrievers and rerankers. For retrieval, standard approaches extract dense representations from the last hidden states using contrastive learning (Xu et al., 2024; Li et al., 2024b; Jiang et al., 2024c). While recent methods, such as PromptEOL (Jiang et al., 2024a), enhance discriminative power via Explicit One-word Limitation (EOL), they require resource-intensive training. To address this, PromptReps (Zhuang et al., 2024b) introduces a training-free hybrid retriever that generates dense and sparse representations in a single pass. Regarding reranking, LLM-based point-wise rerankers (Liu et al., 2024b; Long et al., 2025), pair-wise rerankers (Li et al., 2025; Qin et al., 2024), and list-wise rerankers (Ma et al., 2023; Chen et al., 2025) have been extensively ex-

065 plored. In comparison, exploration of the zero-shot
 066 setting in cross-modal retrieval and reranking has
 067 been quite limited. Only a few successful experi-
 068 ences (Jiang et al., 2024b) from text modality have
 069 been adapted to cross-modal domains and demon-
 070 strated their effectiveness.

071 To address the aforementioned research gap, this
 072 paper explores adapting MLLMs as hybrid retriev-
 073 ers and rerankers. Specifically, our research pro-
 074 ceeds in two stages:

075 First, regarding hybrid retrieval, we evaluate
 076 two effective, training-free sparse representation
 077 methods from text-modality studies. Empirically,
 078 neither achieves satisfactory sparse and hybrid re-
 079 trieval performance in cross-modal settings, primar-
 080 ily because they fail to incorporate prior dataset
 081 knowledge from diverse perspectives. To ad-
 082 dress this, we propose Multi-Perspective Prompt-
 083 ing (MPP) sparsification strategy. MPP guides the
 084 MLLM to summarize information from dataset-
 085 specific angles and aggregates weights from these
 086 perspectives, yielding sparse representations with
 087 enhanced expressive power.

088 Second, for reranking, we evaluate two point-
 089 wise paradigms suitable for MLLM’s context lim-
 090 its: Relevance Generation Paradigm (RGP) and
 091 Caption Generation Paradigm (CGP). RGP pre-
 092 dicts text-image relevance directly, while CGP
 093 predicts candidate conditional probability. Inter-
 094 estingly, CGP excels in text-to-image reranking,
 095 whereas RGP performs better in image-to-text
 096 tasks. We attribute CGP’s shortcomings in image-
 097 to-text reranking to deviations in the prior distri-
 098 bution of text candidates. Consequently, to avoid
 099 the high computational cost of debiasing, we im-
 100 plement distinct strategies for each direction, si-
 101 multaneously optimizing both image-to-text and
 102 text-to-image retrieval.

103 Based on our analysis, we propose Re-M, a
 104 two-stage zero-shot cross-modal retrieval frame-
 105 work (Figure 1). First, Re-M enhances dense rep-
 106 resentations with MPP-derived sparse features to
 107 achieve robust hybrid retrieval. Second, it em-
 108 ploys RCP and CGP reranking to optimize both
 109 text-to-image and image-to-text tasks. Extensive
 110 experiments show that Re-M rivals or surpasses
 111 specialized dense and sparse baselines in zero-shot
 112 settings. Ultimately, Re-M helps the community
 113 better understand the retrieval potential of MLLMs,
 114 enabling fine-tuning improvements with a more
 115 model-adaptive approach or utilizing it as a higher-
 116 quality distillation supervision signal.

2 Background 117

2.1 MLLM for Zero-Shot Dense Retrieval 118

119 Previous works, such as PromptEOL (Jiang et al.,
 120 2024a) and E5-V (Jiang et al., 2024b), show that
 121 Explicit One-word Limitation (EOL) embedding
 122 prompts can produce discriminative representa-
 123 tions. Specifically, the MLLM F can be divided
 124 into the language modeling head g and other parts
 125 f containing the vision encoder, the projector, and
 126 the LLM transformer. which can be expressed
 127 as $F = g \circ f$. To obtain dense representations,
 128 we prepare prompt templates $\langle sent \rangle \backslash n$ *Summary*
 129 *above sentence in one word:* for texts and $\langle im-$
 130 *age \rangle \backslash n* *Summary above image in one word:* for
 131 images. For image v , the MLLM image encoder
 132 is used to project v into the input embedding
 133 space, and the projected patch embeddings are
 134 filled into the template to replace $\langle image \rangle$. For
 135 text t , it is directly replaced $\langle sent \rangle$ in the tem-
 136 plate. We uniformly define the filled sequence as
 137 $s = \{s_1, s_2, \dots, s_l\}$, where l is the sequence length,
 138 then, the d -dimensional hidden states from the last
 139 layer can be expressed as

$$\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_l] = f(s) \quad (1) \quad 140$$

141 where $\mathbf{H} \in \mathbb{R}^{d \times l}$ and $\mathbf{h}_i \in \mathbb{R}^{d \times 1}$. Guided by
 142 the templates, the semantics of images or text is
 143 aggregated at the final position \mathbf{h}_l , and thus can be
 144 regarded as representations for dense retrieval. For
 145 convenience, we define the text representation as
 146 $\mathbf{h}^{(t)}$ and the image representation as $\mathbf{h}^{(v)}$.

2.2 LLM for Zero-Shot Hybrid Retrieval 147

Sparse Representation 148 With the language mod-
 149 eling head g , the hidden state can be mapped into
 150 logit $\mathbf{w} = g(\mathbf{h}_l) \in \mathbb{R}^{|\mathcal{V}|}$, where $|\mathcal{V}|$ denotes the vo-
 151 cabulary size. PromptReps (Zhuang et al., 2024b)
 152 shows \mathbf{w} can be used for sparse representation
 153 through post-processing. Specifically, following
 154 the recipe of SPLADE (Formal et al., 2021), the
 155 initial sparse representation can be calculated as

$$\hat{\mathbf{w}} = \text{Round}(100 \times \log(1 + \text{ReLU}(\mathbf{w}))) \quad (2) \quad 156$$

157 where the ReLU function avoids negative values,
 158 log-saturation prevents the domination of certain to-
 159 kens, and the Round function quantizes logit values
 160 to the nearest integers. To build sparser representa-
 161 tions, two potential methods from previous works
 162 are available:

- **Source-guided Sparse Representing (SSR)** retains the tokens that appear in the original text while zero the other positions. (Zhuang et al., 2024b)

$$\hat{\mathbf{w}}_{\text{ssr}} = \hat{\mathbf{w}} \odot (\mathbb{1}(\mathbb{V}_j \in t))_{j=1 \dots |\mathbb{V}|}, \quad (3)$$

where t is the original text and \mathbb{V}_j is the token corresponding to j^{th} indice in $\hat{\mathbf{w}}$.

- **Prediction-based Top-k Truncation (PTT)** retains the tokens with the top- k largest logit values. (Nie et al., 2025)

$$\hat{\mathbf{w}}_{\text{ptt}} = \hat{\mathbf{w}} \odot (\mathbb{1}(\hat{\mathbf{w}}_j \in \text{Topk}(\hat{\mathbf{w}})))_{j=1 \dots |\mathbb{V}|}, \quad (4)$$

where $\hat{\mathbf{w}}_j$ is the token weight in j^{th} indice.

Hybrid Scoring For hybrid retrieval, a linear interpolation with appropriate weight α is introduced to fusion the dense and sparse scores score_d and score_s to generate hybrid scores score_h . Given the query-document dense representation $\mathbf{h}^{(q)}$ and $\mathbf{h}^{(d)}$ and sparse representation $\hat{\mathbf{w}}^{(q)}$ and $\hat{\mathbf{w}}^{(d)}$, the hybrid score is calculated as

$$\begin{aligned} \text{score}_d &= \mathbf{h}^{(q)} \cdot \mathbf{h}^{(d)} & \text{score}_s &= \hat{\mathbf{w}}^{(q)} \cdot \hat{\mathbf{w}}^{(d)} \\ \text{score}_h &= \alpha \Gamma(\text{score}_d) + (1 - \alpha) \Gamma(\text{score}_s) \end{aligned} \quad (5)$$

where α is the weight hyper-parameter, Γ is the min-max normalization. Since hidden states and logits can be obtained in a single forward pass, we integrate three distinct basic paradigms for zero-shot cross-modal retrieval within a hybrid retriever without double-counting.

2.3 LLM for Zero-Shot Reranking

LLMs for zero-shot reranking have been extensively studied, including point-wise (Liu et al., 2024b; Long et al., 2025), pair-wise (Li et al., 2025; Qin et al., 2024), and list-wise (Ma et al., 2023; Chen et al., 2025). Considering the context length limitations of MLLMs, we primarily focus on point-wise strategies for reranking. Point-wise requires a query and a candidate as input, prompts the LLM to measure the relevance, and outputs a ranking score. Point-wise methods include two paradigms: query generation and relevance generation. Query generation calculates the average log-likelihood of the query conditioned on the candidate as the ranking score (Sachan et al., 2022; Zhuang et al., 2023). Relevance generation applies the softmax function on the logit of token 'yes' and 'no' or other relevance levels (Zhuang et al., 2024a; Liu et al., 2024b) to generate ranking scores.

3 Methodology

3.1 Acquisition of Sparse Representation

Earlier, we detailed constructing hybrid retrievers via two sparse methods, previously limited to document retrieval. We first explore their cross-modal feasibility. Since images lack original sentences, only SSR is applicable for the image modality. For the text modality, we build corresponding hybrid retrievers to compare both strategies, using E5-V prompts and setting PTT with $k=128$. Results on Flickr30K are shown in Figure 2. We report average recall@1 on Flickr30K (Lin et al., 2014) in Figure 2. Notably, similar conclusions hold for other datasets (e.g., MSCOCO (Plummer et al., 2015)), with details in Appendix C.

As shown in Figure 2, there is a contradiction between sparse and hybrid retrieval with SSR and PTT, which indicates that neither of them can perform well simultaneously. We conjecture that, compared to SSR, PTT’s sparse representations are simply the spatial projection of dense representations; therefore, they share more identical tokens (blue tokens in Table 1) when their corresponding dense representations are semantically similar. This property makes PTT’s sparse representations suitable for token matching, leading to respectable sparse retrieval performance. However, this results in dense and sparse representations encoding homogeneous information. For hybrid retrieval to outperform dense-only or sparse-only approaches, the two representations should provide heterogeneous content. Relying solely on PTT for both modalities thus limits the potential gains in hybrid retrieval. Therefore, while each sparsification method has its strengths, further research is needed to develop new sparse approaches better suited for cross-modal retrieval and improve both sparse and hybrid retrieval performance simultaneously.

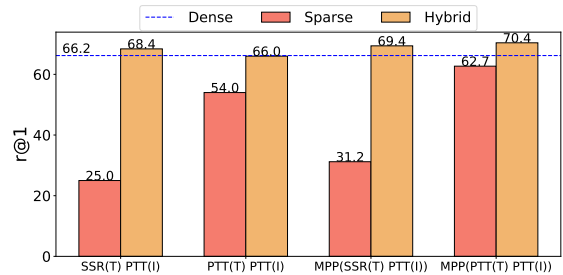


Figure 2: Average recall@1 on Flickr30K of text-to-image and image-to-text on LLaVA-Next-7B. The letter in parentheses denotes the modality, where “I” means image modality and “T” means text modality.


Image & Text	Method	Modality	Perspective	Top 10 Predicted Tokens
 <p>A brown dog is shown standing in the water near a muddy beach.</p>	SSR	text	coarse	_Dog, _Beach, _Mud, _Brown , _Standing , _Water, _Near, _Shown, _Dy
	PTT	text	coarse	_Dog , _Beach , _Brown , _Be , _W , _A , _Sw , _Water , _Yes , _The
		image	coarse	_F, _Dog , _Mist , _M , _Rain , _Beach , _W , _Cloud , _H , _D
	MPP	text	person or objects	_Dog , _Beach , _M , _Brown , _ , _Water , _Sand , _Be , _One , _No
			relations	_Dog , _Beach , _M , _Water , _Sand , _None , _Own , _Near , _No , _Brown
			environment, weather or places	_Beach , _Water , _Sh , _M , _Sw , _Ocean , _W , _River , _Be , _Coast
		image	actions or movements	_Stand , _Dog , _Sw , _St , _W , _Beach , _M , _Wait , _Walk , _Dr
			appearance	_Dog , _Brown , _M , _Color , _Water , _The , _ , _Beach , _A , _Be
			person or objects	_Dog , _F , _ , _None , _No , _One , _Water , _Nothing , _Cloud , _Sand
	image	relations	_Dog , _None , _No , _ , _F , _Un , _Run , _Water , _Sep , _One	
		environment, weather or places	_F , _Mist , _M , _Beach , _Cloud , _Sw , _W , _Rain , _D , _H	
		actions or movements	_Walk , _Run , _Running , _Dog , _Stand , _J , _Sw , _W , _Play , _S	
			appearance	_Dog , _Brown , _F , _M , _W , _Water , _D , _Gray , _Mist , _Sand

Table 1: The case study of the top 10 predicted tokens with SSR, PTT, and MPP. The 'coarse' indicates the ordinary EOL templates used in the Background Section.

3.2 Multi-Perspectives Prompting

Motivation When analyzing the sparse weights from existing methods, we find that MLLM struggles to grasp the aspects that summarize information in the dataset and focuses on core tokens via coarse-grained summary prompt templates. For example, when summarizing the image, predicted tokens with the top-10 largest weights exclude token '_Standing', which reflects the action of the dog, and token '_Brown' (red tokens in Table 1) reflecting the color of the dog. Meanwhile, the empirical success of test-time scaling in other areas, such as semantic text similarity (Lei et al., 2024) and image classification (Xiao, 2024), motivates us to consider predefined dataset-specific perspectives in prompt templates, enabling the MLLM to summarize information from multiple angles and thereby supplement the model with multifaceted prior knowledge about the dataset. Thus, we propose **Multi-Perspectives Prompting (MPP)** to improve sparse representations.

Prompting Design Through analysis of the image-text pairs in datasets, we realize that although data annotators describe in diverse forms, they always annotate in a constrained number of dataset-specific perspectives. For example, images and texts in MSCOCO and Flickr30K can be described from five angles: person or object, relations, environments, actions, and appearances. Based on this discovery, we collect a set of perspectives for each dataset¹. We revise prompt templates to *<image>/<sent>\n Summary the <angle> in the*

¹In this paper, we adopt a manual approach for perspective construction, reserving the method of iterative optimization via LLMs for future work.

above image/sentence in one word:. Then, we replace the *<angle>* token with a perspective from the set and obtain a group of prompt templates \mathcal{T} . We denote $L = |\mathcal{T}|$ as the sum of perspectives to guide MLLM in summarizing information from various perspectives. Concrete modified prompt templates for different benchmarks are shown in the Appendix B. We quantize token weights for each prompt template and accumulate them to produce sparse representations \hat{w}_{mpp} :

$$\hat{w}_{mpp} = \sum_{a \in \mathcal{T}} \hat{w}_a \quad (6)$$

where a is a prompt template with one perspective in \mathcal{T} , \hat{w}_a is the sparse representation produced by the prompt template a . We apply PTT to both modalities for sparsity, enabling the MLLM to explicitly capture the salient perspectives in the dataset content while maintaining high token matching accuracy.

Efficiency Optimization During retrieval, we reserve the prompt templates used by basic retrievers for dense representations to embed global semantics and use MPP for sparse representations. To address inefficient multiple forward passes, we introduce prompt-wise attention masking (Kim and Angelova, 2025; Li et al., 2023) and concatenate all prompt templates into a single sequence, as shown in Figure 3. We adjust the attention mask to prevent subsequent prompts from attending to earlier ones and extract representations at the end of each prompt to simulate the last pooling operation. In this way, MPP reduces latency while retaining the advantage of acquiring multiple retrieval paradigms within one forward pass and achieves simultaneous

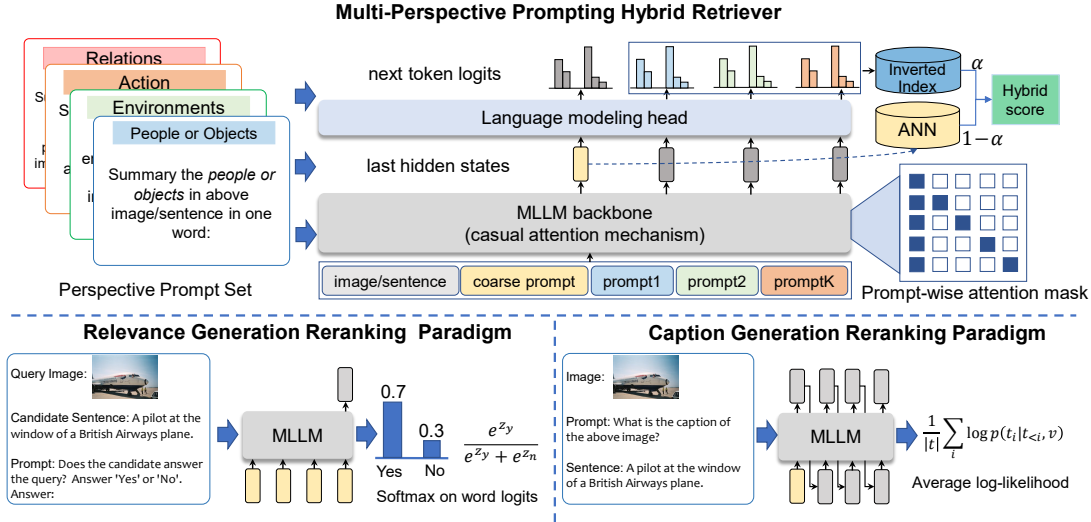


Figure 3: The illustration of our two-stage zero-shot cross-model framework, Re-M. In the retrieval phase, we propose multi-perspective prompting to enhance the performance of sparse retrieval and hybrid retrieval. In the reranking phase, we employed relevance generation for text reranking and caption generation for image reranking to improve performance on both sides simultaneously.

progress in sparse and hybrid retrieval.

Analysis of MPP From Table 1, we find MPP effectively addresses the issue of SSR and PTT with a coarse summary prompt. MLLM understands how to summarize information from various angles and predict the key tokens (highlighted green) with a top-weight ranking, such as the tokens ‘_Brown’ and ‘_Stand’. Finally, it is intuitive that some key tokens appear in various predictions repeatedly, and the weight accumulation mechanism emphasizes the key tokens (highlighted orange) with a larger weight to enable the sparse representations to be more discriminative, avoiding interference from similar negative candidates. Additionally, we report the performance of MPP based on SSR and PTT in Figure 2, denoted as “MPP(SSR(T), PTT(I))” and “MPP(PTT(T), PTT(I))”. It is observed that when achieving similar hybrid retrieval results, MPP based on PTT gains remarkable sparse retrieval performance compared to MPP based on SSR. Therefore, we employ construct MPP upon PTT by default to support hybrid retrieval.

3.3 Prompting MLLM as a Reranker

In the modern retrieval pipeline, the final performance can be significantly improved by employing a reranker that follows the retriever. However, currently, there is no research to prompt MLLM as a reranker in cross-modal retrieval. Thus, we adapt two existing paradigms to enable MLLMs to ac-

quire reranking capability.

Relevance Generation Paradigm (RGP) Given an image v and a text t , the MLLM is prompted with a template to determine whether they are relevant and outputs “yes” or “no”. Denote the logit of token ‘Yes’ and token ‘No’ with z_Y and z_N , then the ranking score $r(t, v)$ is calculated by applying a softmax function:

$$r(t, v) = \frac{\exp(z_Y)}{\exp(z_Y) + \exp(z_N)}. \quad (7)$$

There are numerous prompt templates for selection. We referenced the templates provided in a systematic study on document reranking (Sun et al., 2025) and adapted them for image-text reranking. Preliminary experiments were conducted to select the optimal prompt template from available options. We report the results for all templates in Appendix D.

Caption Generation Paradigm (CGP) Given an image v and a text t , the log-likelihood of the sentence conditioned on the matching image is calculated as the reranking scores. Normally, we use $\log p(t|v)$ to estimate the image-to-text reranking score and $\log p(v|t)$ to calculate the text-to-image reranking score. Assume that the prior distribution $p(v)$ is uniform, according to Bayes’ rule, $\log p(v|t)$ should be proportional to $\log p(t|v)$. Therefore, we prompt MLLM to calculate the average log-likelihood $\log p(t|v)$ as the reranking score

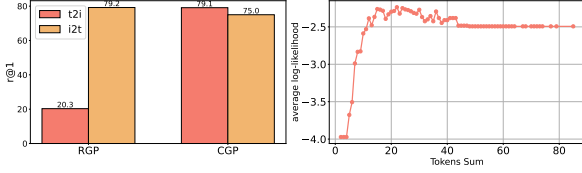


Figure 4: The recall@1 results of RGP and CGP (left) and average log-likelihood generating texts conditioned on no image (right) on Flickr30K.

$r(t, v)$ for both retrieval settings:

$$r(t, v) = \frac{1}{|t|} \sum_i^{|t|} \log p(t_i | t_{<i}, v, \theta) \quad (8)$$

where θ denotes the parameters of MLLM and $|t|$ is the length of text tokens.

Analysis of Performance Disparity We report our results on the Flickr30K validation set in the left of Figure 4, which demonstrates that the optimal methods differ between text-to-image and image-to-text tasks. Specifically, CGP demonstrates significantly better performance in text-to-image reranking, while RGP holds an advantage in image-to-text reranking. Considering MLLM as a language model fundamentally, we believe that the CGP paradigm aligns more closely with MLLM training patterns than RGP, which should yield superior results on both sides. However, as shown in the right of Figure 4, we observed significant discrepancies in the inherent generation probabilities across texts of varying lengths. When presented with one image and multiple candidate texts (i, t_1, \dots, t_N) , the probability $p(t_j | i)$ becomes distorted by the influence of $p(t_j)$'s own probability, introducing bias. In contrast, when presented with one image and multiple candidate texts (t, i_1, \dots, i_N) , the fixed variable t in $p(i_j | t) \propto p(t | i_j)$ eliminates this interference. To address this issue in graph-to-text reordering, two forward passes are required: calculating $p(t_i | i)$ and $p(t_i)$, followed by bias removal. However, this approach significantly increases computational complexity. Therefore, during reranking, **we select RGP for image-to-text and CGP for text-to-image to ensure optimal performance on both sides while maintaining efficiency.**

3.4 Retrieval Pipeline: Re-M

As shown in Figure 3, we design the Re-M pipeline for training-free, zero-shot cross-modal retrieval. In the first stage, we exploit dense representations

for semantic similarity calculation with ANN and sparse representations for token matching with an inverted index separately. The hybrid score calculation method for text and images is consistent with Equation 5, except that the sparse representation weights are replaced by the multi-perspectives prompting strategy (Equation 6). In the second stage, we select the top N candidates with the highest hybrid scores and rerank the candidate texts with RGP, while the candidate images with CGP.

4 Experiment

4.1 Experiment Setup

Dataset and Evaluation We assess baselines and Re-M on widely used image-text retrieval datasets MSCOCO (Lin et al., 2014) and Flickr30K (Plummer et al., 2015) following the Karpathy split (Karpathy and Fei-Fei, 2015). Each image is combined with five sentences. To evaluate the generalizability of our model on domain datasets, we select the CUHK-PEDES (Li et al., 2017), ICFG-PEDES (Ding et al., 2021), and RST-PReid (Zhu et al., 2021) benchmarks in text-based person retrieval. We report recall@k with k=1, 5, 10 for all retrieval datasets.

Baselines We compare Re-M with widely recognized dense models for zero-shot image-text retrieval. We select CLIP with ViT-B and ViT-L (Radford et al., 2021), BLIP with ViT-L (Li et al., 2022), EVA-02-CLIP with 5B parameters (Sun et al., 2023) and E5-V (Jiang et al., 2024b). We also introduce some image lexical models, such as LexLIP (Luo et al., 2023), STAIR (Chen et al., 2023), and D2S (Nguyen et al., 2024), as baselines.

Implementation Details We design the MPP retriever and Re-M pipeline based on two strong baseline MLLM: LLaVA-Next-8B (Li et al., 2024a) and LLaVA-Next-7B (Liu et al., 2024a). We choose the Faiss (Douze et al., 2025) library to construct the ANN index with cosine similarity to rank and use brute force search. We use Pyserini (Lin et al., 2021) to build the inverted index for sparse retrieval. The specific hybrid method follows ranx.fuse (Basani and Romelli, 2022). We set k to 30 for image-text retrieval of MPP, and to 10 for text-based person retrieval. Due to the existence of repeated tokens, this can ensure the length of MPP is shorter than PTT. In the reranking stage, we set the candidate size N to 5 and evaluate the reranker on the recall@1 metric. During the evaluation, we use the

Method		MSCOCO						Flickr30K					
		t2i			i2t			t2i			i2t		
		r@1	r@5	r@10	r@1	r@5	r@10	r@1	r@5	r@10	r@1	r@5	r@10
Dense Baselines	CLIP ViT-B	30.5	56.0	66.8	51.0	74.9	83.5	58.8	83.3	89.8	77.8	95.0	98.2
	CLIP ViT-L	37.0	61.6	71.5	58.1	81.1	87.8	67.3	89.0	93.3	87.2	98.3	99.4
	BLIP ViT-L	48.4	74.4	83.2	63.5	86.5	92.5	70.0	91.2	95.2	75.5	95.1	97.7
	EVA-02-CLIP	51.1	75.0	82.7	68.8	87.8	92.8	78.8	94.2	96.8	93.9	99.4	99.8
	E5-V	<u>52.0</u>	76.5	84.7	<u>62.0</u>	83.6	89.7	<u>79.5</u>	95.0	97.6	<u>88.2</u>	98.7	99.4
Sparse Baselines	LexLIP	53.2	79.1	86.7	70.2	90.7	95.2	78.4	94.6	97.1	91.4	99.2	99.7
	STAIR	41.1	65.4	75.0	57.7	80.5	87.3	66.6	88.7	93.5	81.2	96.1	98.4
	D2S	54.5	80.6	-	-	-	-	79.8	95.9	-	-	-	-
LLaVA-Next-LLaMA-8B	Dense	34.6	60.5	71.2	43.3	68.6	78.2	60.1	83.7	90.3	72.4	90.6	95.6
	MPP Sparse	25.8	51.4	63.3	29.6	57.7	70.4	52.5	79.4	86.4	55.6	83.6	90.4
	MPP Hybrid	37.3	63.7	74.1	46.0	71.2	80.5	64.4	86.7	91.9	74.1	91.9	96.1
	Re-M	51.6	-	-	50.7	-	-	79.1	-	-	79.2	-	-
LLaVA-Next-Mistral-7B	Dense	34.3	59.7	70.3	42.7	67.2	77.1	60.5	83.0	89.0	72.3	90.9	94.6
	MPP Sparse	30.5	56.1	67.5	39.2	66.8	77.7	57.0	82.3	89.1	68.4	89.6	94.5
	MPP Hybrid	38.1	64.5	74.8	46.1	72.4	81.5	65.5	87.5	92.9	75.4	93.9	97.3
	Re-M	52.1	-	-	49.4	-	-	79.8	-	-	80.0	-	-

Table 2: The image-text retrieval results on MSCOCO and Flickr30K. The models painted grey are finetuned on MSCOCO and Flickr30K, and the others are evaluated in the zero-shot setting. Here, the letter 'r' indicates recall.

validation set to determine the optimal weights α (ranging from 0.1 to 0.9) for the hybrid score. All experiments are run on two Nvidia V100 GPUs.

4.2 Main Results

Image-Text Retrieval We present the image-text retrieval results of Re-M and baselines on MSCOCO and Flickr30K in Table 2. The first observation is that MPP hybrid retrieval outperforms both dense and sparse retrieval alone, confirming the effectiveness of the weighted hybrid method in cross-modal retrieval. Recalls are improved by 3 to 4 points, and for $r@1$ with LLaVA-Next-7B, the improvement reaches 5 points. Another finding is that the Re-M reranking performance outperforms strong zero-shot dense baselines. Meanwhile, the reranking pipeline produces competitive results compared to finetuned sparse baselines. Although exceeding EVA-02-CLIP is difficult, our Re-M pipeline maintains its entire generation capability, representing a fundamental advantage.

Text-based Person Retrieval To evaluate the generalization of Re-M, we report results on text-based person retrieval in Table 3, which can be viewed as a domain task. In Table 3, the results demonstrate that (1) MPP hybrid retrievers are still effective on domain tasks. (2) Re-M pipeline can outperform strong dense baselines on TBPR tasks, maintaining similar phenomena with ITR. The $r@1$ results after reranking are even closer to $r@5$ of

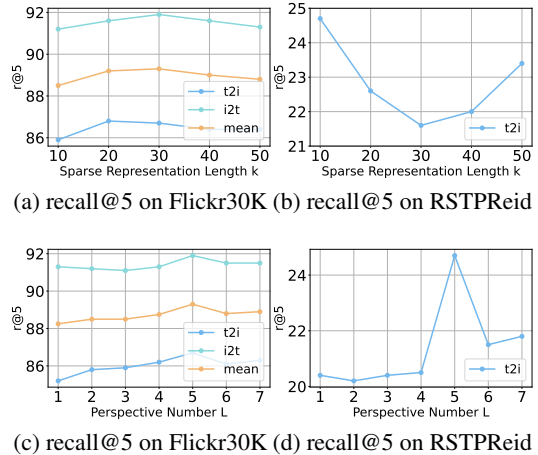


Figure 5: Sparse representation length k (Figure (a) and (b)) and perspective number L (Figure (c) and (d)) of MPP analysis on Flickr30K and RSTPReid, reporting the $r@5$ of MPP hybrid retrieval on LLaVA-Next-8B.

hybrid retrieval. This suggests the generalization of our method without finetuning on domain datasets.

4.3 Ablation Study

Analysis of Sparse Representation Length k Figure 5a and 5b present the trend of $r@5$ for MPP hybrid retrieval on LLaVA-Next-8B across the Flickr30K and RSTPReid datasets, with the increase of length k corresponding to that of the sparse representations w_a from each perspective. In Figure 5a, $r@5$ attains a high value even at a

Method		CUHK-PEDES			ICFG-PEDES			RSTPReid		
		r@1	r@5	r@10	r@1	r@5	r@10	r@1	r@5	r@10
Dense Baselines	CLIP ViT-B	5.9	14.7	21.0	3.2	9.3	13.9	8.9	21.1	31.2
	CLIP ViT-L	10.1	22.8	31.2	5.2	13.3	19.4	11.2	26.4	37.8
	BLIP ViT-L	11.2	26.4	37.8	8.1	19.0	26.1	16.9	38.0	50.9
	EVA-02-CLIP	29.3	47.6	56.4	18.0	33.6	42.2	27.0	50.4	62.0
	E5-V	20.6	37.5	47.0	7.5	17.6	24.3	21.6	42.4	54.8
LLaVA-Next-LLaMA-8B	Dense	5.3	12.7	16.9	1.6	5.0	7.6	7.6	20.7	31.6
	MPP Sparse	5.8	13.3	19.0	1.6	5.5	8.7	6.3	21.3	32.2
	MPP Hybrid	7.0	15.4	21.5	2.2	6.8	10.6	8.6	24.7	36.1
	Re-M	14.9	-	-	5.7	-	-	20.0	-	-
LLaVA-Next-Mistral-7B	Dense	6.8	15.9	22.4	2.3	6.8	10.2	9.6	26.7	36.8
	MPP Sparse	11.7	25.0	32.8	4.8	12.3	17.5	13.4	31.9	43.0
	MPP Hybrid	12.9	26.7	35.7	5.0	12.9	18.4	14.8	34.1	16.9
	Re-M	<u>24.8</u>	-	-	<u>10.7</u>	-	-	<u>26.8</u>	-	-

Table 3: The text-based person retrieval results on CUHK-PEDES, ICFG-PEDES, and RSTPReid. In the table, the letter 'r' indicates recall.

length of 10. As k increases, recall improves gradually and peaks when the length reaches 30. On the text-based person retrieval task (RSTPReid, Figure 5b), the best performance is observed at length 10, with degradation at longer lengths. These results suggest that high sparsity can be maintained while ensuring effective hybrid retrieval.

Analysis of Perspective Number L Figure 5c and 5d report the trend of r@5 for MPP hybrid retrieval on LLaVA-Next-8B across the Flickr30K and RSTPReid datasets under varying numbers L of perspective prompts. The results show that hybrid performance peaks when the number of perspectives reaches five, at which point the perspective set sufficiently covers the dataset content. However, introducing less relevant aspects leads to performance degradation. This indicates that including only the most pertinent perspectives is critical to produce reliable sparse representations, while irrelevant ones should be excluded.

4.4 Discussion of the Three-stage Pipeline

Our evaluation shows that sparse-only MPP retrieval achieves high recall when the retrieval set is sufficiently large. This inspires the conception of a three-stage pipeline that first applies MPP sparse retrieval to select 200 candidates coarsely, followed by hybrid retrieval and reranking. We analyze the time-cost upper bound of this pipeline on the LLaVA-Next-8B, assuming all tokens in the sparse representations are matched, and the maximum length is 150 (30 tokens per perspective). The

pipeline is compared to CLIP ViT-L, and FLOP trends are illustrated in Figure 6 as the data scale increases. When the data scale exceeds about 125 billion, the three-stage pipeline has lower FLOPs than CLIP, indicating its feasibility for super-large-scale search. It is worth noting that Figure 6 reflects the worst-case situation. Through the use of sparser representations, high-performance GPUs, and the KV-cache technique, the time cost can be further reduced for applications on large-scale data.

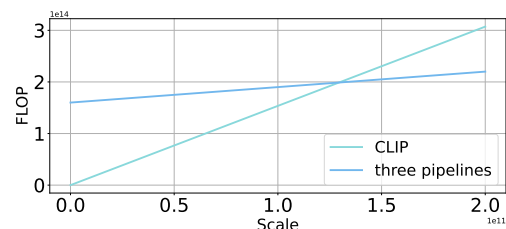


Figure 6: FLOP of three-stage pipeline and CLIP.

5 Conclusion

In this work, we investigate the Re-M pipeline for adapting MLLM to zero-shot cross-modal retrieval and reranking via prompting. We design an efficient hybrid retriever and propose Multi-Perspective Prompting to improve sparse and hybrid retrieval simultaneously. Meanwhile, we introduce two paradigms to prompt MLLM as the reranker to address the asymmetric improvements. Experiments on image-text retrieval datasets and domain text-based person retrieval benchmarks demonstrate the effectiveness of Re-M.

546 Limitations

- 547 • When designing perspectives for datasets, we
548 must investigate the concrete images and sen-
549 tences in the dataset and summarize the de-
550 scription angles of annotators. This process
551 is not automatic, and humans struggle to pro-
552 duce the optimal dataset-specific perspectives
553 manually. However, we did not find a proper
554 method to generate summary perspectives
555 without humankind. For example, we should
556 have considered prompting the LLM to gener-
557 ate and iteratively optimize perspectives, like
558 producing schemas for knowledge graphs au-
559 tomatically. We believe this should be re-
560 searched in the future.
- 561 • Because of the time limitation, we only set
562 the size of the candidate pool for re-ranking
563 to 5 and report $r@1$ in our experiments. This
564 setting constrains the performance of RGP
565 and CGP, so we plan to supplement results
566 to set the size of the candidate pool to 20 and
567 record results on $r@1$, $r@5$, and $r@10$ metrics
568 to further elaborate on the effectiveness of our
569 re-ranking paradigms.
- 570 • In our studies, to achieve improvements on
571 both retrieval settings, we adopt RGP on
572 image-to-text retrieval and CGP on text-to-
573 image retrieval. However, this operation in-
574 dicates that we do not explain why RGP and
575 CGP cannot perform well on both retrieval
576 settings at the same time. In the paper, we
577 simply report the average log-likelihood of
578 various token sequence lengths of sentences
579 conditioned on no images. This may mod-
580 ify the CGP log-likelihood of image-to-text
581 retrieval and improve the results. But there
582 is no experimental evidence to support this.
583 Besides, this paper lacks an analysis of RGP.
- 584 • Due to the absence of an efficiency optimiza-
585 tion mechanism, the process of re-ranking is
586 computationally intensive and leads to huge
587 FLOPs. This influences the application of
588 Re-M on small-scale data. In the future, re-
589 searchers should propose an approach to solve
590 this issue.

References 591

- Md Atabuzzaman, Andrew Zhang, and Chris Thomas. 2025. Zero-shot fine-grained image classification using large vision-language models. *arXiv preprint arXiv:2510.03903*. 592-593-594-595
- Elias Bassani and Luca Romelli. 2022. *ranx.fuse: A python library for metasearch*. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, pages 4808–4812. ACM. 596-597-598-599-600
- Chen Chen, Bowen Zhang, Liangliang Cao, Jiguang Shen, Tom Gunter, Albin Jose, Alexander Toshev, Yantao Zheng, Jonathon Shlens, Ruoming Pang, and Yinfei Yang. 2023. *STAIR: Learning sparse text and image representation in grounded tokens*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15079–15094, Singapore. Association for Computational Linguistics. 601-602-603-604-605-606-607-608-609
- Shijie Chen, Bernal Jimenez Gutierrez, and Yu Su. 2025. Attention in large language models yields efficient zero-shot re-rankers. In *The Thirteenth International Conference on Learning Representations*. 610-611-612-613
- Zefeng Ding, Changxing Ding, Zhiyin Shao, and Dacheng Tao. 2021. Semantically self-aligned network for text-to-image part-aware person re-identification. *arXiv preprint arXiv:2107.12666*. 614-615-616-617
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2025. The faiss library. *IEEE Transactions on Big Data*. 618-619-620-621-622
- Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2025. Colpali: Efficient document retrieval with vision language models. In *ICLR*. 623-624-625-626
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. Splade: Sparse lexical and expansion model for first stage ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2288–2292. 627-628-629-630-631-632
- Yunqi Hong, Sohyun An, Andrew Bai, Neil YC Lin, and Cho-Jui Hsieh. 2025. Unlabeled data improves fine-grained image zero-shot classification with multimodal llms. *arXiv preprint arXiv:2506.03195*. 633-634-635-636
- Ting Jiang, Shaohan Huang, Zhongzhi Luan, Deqing Wang, and Fuzhen Zhuang. 2024a. *Scaling sentence embeddings with large language models*. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3182–3196, Miami, Florida, USA. Association for Computational Linguistics. 637-638-639-640-641-642
- Ting Jiang, Minghui Song, Zihan Zhang, Haizhen Huang, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, and Fuzhen Zhuang. 2024b. E5-v: Universal 643-644-645

646	embeddings with multimodal large language models.	information retrieval research with sparse and dense	701
647	<i>arXiv preprint arXiv:2407.12580</i> .	representations. In <i>Proceedings of the 44th inter-</i>	702
648	Ziyan Jiang, Rui Meng, Xinyi Yang, Semih Yavuz,	<i>national ACM SIGIR conference on research and</i>	703
649	Yingbo Zhou, and Wenhua Chen. 2024c. Vlm2vec:	<i>development in information retrieval</i> , pages 2356–	704
650	Training vision-language models for massive	2362.	705
651	multimodal embedding tasks. <i>arXiv preprint</i>	Tsung-Yi Lin, Michael Maire, Serge Belongie, James	706
652	<i>arXiv:2410.05160</i> .	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,	707
653	Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-	and C Lawrence Zitnick. 2014. Microsoft coco:	708
654	semantic alignments for generating image descrip-	Common objects in context. In <i>European confer-</i>	709
655	tions. In <i>Proceedings of the IEEE Conference on</i>	<i>ence on computer vision</i> , pages 740–755. Springer.	710
656	<i>Computer Vision and Pattern Recognition (CVPR)</i> .	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan	711
657	Dahun Kim and Anelia Angelova. 2025. Context-	Zhang, Sheng Shen, and Yong Jae Lee. 2024a. <i>Llava-</i>	712
658	adaptive multi-prompt embedding with large lan-	<i>next: Improved reasoning, ocr, and world knowledge</i> .	713
659	guage models for vision-language alignment. <i>arXiv</i>	Wenhan Liu, Yutao Zhu, and Zhicheng Dou. 2024b.	714
660	<i>preprint arXiv:2508.02762</i> .	Demorank: Selecting effective demonstrations for	715
661	Yibin Lei, Di Wu, Tianyi Zhou, Tao Shen, Yu Cao,	large language models in ranking task. <i>arXiv preprint</i>	716
662	Chongyang Tao, and Andrew Yates. 2024. <i>Meta-task</i>	<i>arXiv:2406.16332</i> .	717
663	<i>prompting elicits embeddings from large language</i>	Kehan Long, Shasha Li, Chen Xu, Jintao Tang, and	718
664	<i>models</i> . In <i>Proceedings of the 62nd Annual Meeting</i>	Ting Wang. 2025. Precise zero-shot pointwise rank-	719
665	<i>of the Association for Computational Linguistics (Vol-</i>	ing with llms through post-aggregated global context	720
666	<i>ume 1: Long Papers)</i> , pages 10141–10157, Bangkok,	information. In <i>Proceedings of the 48th International</i>	721
667	Thailand. Association for Computational Linguistics.	<i>ACM SIGIR Conference on Research and Develop-</i>	722
668	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Ren-	<i>ment in Information Retrieval</i> , pages 2384–2394.	723
669	rui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and	Ziyang Luo, Pu Zhao, Can Xu, Xiubo Geng, Tao Shen,	724
670	Chunyu Li. 2024a. <i>Llava-next: Stronger llms su-</i>	Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin	725
671	<i>percharge multimodal capabilities in the wild</i> .	Jiang. 2023. Lexlip: Lexicon-bottlenecked language-	726
672	Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia	image pre-training for large-scale image-text sparse	727
673	Shao. 2023. Making large language models a bet-	retrieval. In <i>Proceedings of the IEEE/CVF Interna-</i>	728
674	ter foundation for dense retrieval. <i>arXiv preprint</i>	<i>tional Conference on Computer Vision (ICCV)</i> , pages	729
675	<i>arXiv:2312.15503</i> .	11206–11217.	730
676	Chaofan Li, Zheng Liu, Shitao Xiao, Yingxia Shao, and	Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhua	731
677	Defu Lian. 2024b. Llama2vec: Unsupervised adap-	Chen, and Jimmy Lin. 2024. Unifying multimodal	732
678	tation of large language models for dense retrieval.	retrieval via document screenshot embedding. In <i>Pro-</i>	733
679	In <i>Proceedings of the 62nd Annual Meeting of the</i>	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	734
680	<i>Association for Computational Linguistics (Volume</i>	<i>ods in Natural Language Processing</i> , pages 6492–	735
681	<i>1: Long Papers)</i> , pages 3490–3500.	6505.	736
682	Jieran Li, Xiuyuan Hu, Yang Zhao, Shengyao Zhuang,	Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and	737
683	and Hao Zhang. 2025. Leveraging reference docu-	Jimmy Lin. 2023. Zero-shot listwise document	738
684	ments for zero-shot ranking via large language mod-	reranking with a large language model. <i>arXiv</i>	739
685	els. <i>arXiv preprint arXiv:2506.11452</i> .	<i>preprint arXiv:2305.02156</i> .	740
686	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	Thong Nguyen, Mariya Hendriksen, Andrew Yates, and	741
687	Hoi. 2022. <i>BLIP: Bootstrapping language-image pre-</i>	Maarten de Rijke. 2024. Multimodal learned sparse	742
688	<i>training for unified vision-language understanding</i>	retrieval with probabilistic expansion control. In <i>Eu-</i>	743
689	<i>and generation</i> . In <i>Proceedings of the 39th Interna-</i>	<i>ropean Conference on Information Retrieval</i> , pages	744
690	<i>tional Conference on Machine Learning</i> , volume 162	448–464. Springer.	745
691	of <i>Proceedings of Machine Learning Research</i> , pages	Zhijie Nie, Richong Zhang, and Zhanyu Wu. 2025. <i>A</i>	746
692	12888–12900. PMLR.	<i>text is worth several tokens: Text embedding from</i>	747
693	Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu	<i>LLMs secretly aligns well with the key tokens</i> . In	748
694	Yue, and Xiaogang Wang. 2017. Person search with	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	749
695	natural language description. In <i>Proceedings of the</i>	<i>sociation for Computational Linguistics (Volume 1:</i>	750
696	<i>IEEE Conference on Computer Vision and Pattern</i>	<i>Long Papers)</i> , pages 7683–7694, Vienna, Austria.	751
697	<i>Recognition (CVPR)</i> .	Association for Computational Linguistics.	752
698	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-	Bryan A. Plummer, Liwei Wang, Chris M. Cervantes,	753
699	Hong Yang, Ronak Pradeep, and Rodrigo Nogueira.	Juan C. Caicedo, Julia Hockenmaier, and Svetlana	754
700	2021. Pyserini: A python toolkit for reproducible	Lazebnik. 2015. Flickr30k entities: Collecting	755

756	region-to-phrase correspondences for richer image-to-sentence models. In <i>Proceedings of the IEEE International Conference on Computer Vision (ICCV)</i> .	Wenjie Li, and Min Zhang. 2024. Gme: Improving universal multimodal retrieval by multimodal llms. <i>arXiv preprint arXiv:2412.16855</i> .	812
757			813
758			814
759	Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, and 1 others. 2024. Large language models are effective text rankers with pairwise ranking prompting. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1504–1518.	Aichun Zhu, Zijie Wang, Yifeng Li, Xili Wan, Jing Jin, Tian Wang, Fangqiang Hu, and Gang Hua. 2021. Dssl: Deep surroundings-person separation learning for text-based person retrieval. In <i>Proceedings of the 29th ACM international conference on multimedia</i> , pages 209–217.	815
760			816
761			817
762			818
763			819
764			820
765			
766	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual outputs from natural language supervision . In <i>Proceedings of the 38th International Conference on Machine Learning</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 8748–8763. PMLR.	Honglei Zhuang, Zhen Qin, Kai Hui, Junru Wu, Le Yan, Xuanhui Wang, and Michael Bendersky. 2024a. Beyond yes and no: Improving zero-shot LLM rankers via scoring fine-grained relevance labels . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 358–370, Mexico City, Mexico. Association for Computational Linguistics.	821
767			822
768			823
769			824
770			825
771			826
772			827
773			828
774			829
775	Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Shengyao Zhuang, Bing Liu, Bevan Koopman, and Guido Zuccon. 2023. Open-source large language models are strong zero-shot query likelihood models for document ranking . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8807–8817, Singapore. Association for Computational Linguistics.	830
776			831
777			832
778			833
779			834
780			835
781			836
782			
783	Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023. Eva-clip: Improved training techniques for clip at scale . <i>arXiv preprint arXiv:2303.15389</i> .	Shengyao Zhuang, Xueguang Ma, Bevan Koopman, Jimmy Lin, and Guido Zuccon. 2024b. PromptReps: Prompting large language models to generate dense and sparse representations for zero-shot document retrieval . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 4375–4391, Miami, Florida, USA. Association for Computational Linguistics.	837
784			838
785			839
786			840
787	Shuoqi Sun, Shengyao Zhuang, Shuai Wang, and Guido Zuccon. 2025. An investigation of prompt variations for zero-shot llm-based rankers. In <i>European Conference on Information Retrieval</i> , pages 185–201. Springer.		841
788			842
789			843
790			844
791			
792	Kaiye Wang, Qiyue Yin, Wei Wang, Shu Wu, and Liang Wang. 2016. A comprehensive survey on cross-modal retrieval. <i>arXiv preprint arXiv:1607.06215</i> .	A Dataset Statistics	845
793			
794			
795	Han Xiao. 2024. Scaling test-time compute for embedding models . Jina AI Blog.	We report the data split statistics in Table 4. For ICFG-PEDES, we use the train dataset as the validation dataset to select the best hybrid weight.	846
796			847
797			848
798			
799			
800	Ran Xu, Wenqi Shi, Yue Yu, Yuchen Zhuang, Yanqiao Zhu, May Dongmei Wang, Joyce C. Ho, Chao Zhang, and Carl Yang. 2024. BMRetriever: Tuning large language models as better biomedical text retrievers . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 22234–22254, Miami, Florida, USA. Association for Computational Linguistics.	B Prompt Templates	849
801			
802			
803			
804			
805	Xudong Yan, Songhe Feng, Yang Zhang, Jian Yang, Yueguan Lin, and Haojun Fei. 2024. Leveraging mllm embeddings and attribute smoothing for compositional zero-shot learning. <i>arXiv preprint arXiv:2411.12584</i> .	For MSCOCO and Flickr30K datasets in image-text retrieval, we use the following perspective prompts to summarize information:	850
806			851
807			852
808			
809			
810	Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi Dai, Dingkun Long, Pengjun Xie, Meishan Zhang,	<ul style="list-style-type: none"> • <code><image>/<sent></code>\n Summary the people or objects in above image/sentence in one word: • <code><image>/<sent></code>\n Summary the relations, such as belongings or spatial position, between main people or objects in above image/sentence in one word: • <code><image>/<sent></code>\n Summary the environment, weather or places in above image/sentence in one word: 	853
811			854
			855
			856
			857
			858
			859
			860
			861

- 862 • `<image>/<sent>`\n Summary the actions or
863 movements of main people or objects in above
864 image/sentence in one word:
- 865 • `<image>/<sent>`\n Summary the appearance,
866 such as color, material, decoration and so
867 on, of main people or objects in above im-
868 age/sentence in one word:

869 For CUHK-PEDEs, ICFG-PEDES, and RST-
870 PReid datasets in text-based person retrieval, we
871 use the following prompts to summarize informa-
872 tion:

- 873 • `<image>/<sent>`\n Summary the gender of per-
874 son in above image/sentence in one word:
- 875 • `<image>/<sent>`\n Summary the actions or
876 movements of person in above image/sentence
877 in one word:
- 878 • `<image>/<sent>`\n Summary the objects in
879 above image/sentence in one word:
- 880 • `<image>/<sent>`\n Summary the wearing of
881 person in above image/sentence in one word:
- 882 • `<image>/<sent>`\n Summary the appearance
883 and decoration details of person, such as color,
884 pattern and so on, in above image/sentence in
885 one word:

886 For RGP, we exploit the following prompts to
887 demand MLLM to generate relevance score.

888 When using RGP:

- 889 • text-to-image: For the following query sen-
890 tence and candidate image, judge whether they
891 are relevant. Output 'Yes' or 'No'. \n Query
892 Sentence: `<sent>` Candidate Image: `<image>`
893 Output:
- 894 • image-to-text: For the following query image
895 and candidate sentence, judge whether they
896 are relevant. Output 'Yes' or 'No'. \n Query
897 Image: `<image>` Candidate Sentence: `<sent>`
898 Output:

899 For CGP, we exploit the following prompts to
900 demand MLLM to generate a relevance score:

- 901 • Image: `<image>`\n What is the caption of the
902 above image? `<sent>`"

Dataset	train	val	test
MSCOCO	113200	5000	5000
Flickr30K	29800	1000	1000
CUHK-PEDES	34054	3078	3074
ICFG-PEDES	34674	0	19848
RSTPReid	18505	1000	1000

Table 4: Dataset statistics.

C Results of Basic SSR, PTT, and MPP

We report the whole retrieval results of SSR and PTT on dense, sparse, and hybrid retrieval in Table C. From this table, we can clearly see that simply utilizing SSR or PTT cannot achieve simultaneous improvements of sparse and hybrid retrieval. In addition, although MPP with SSR(T) and PTT(I) obtains similar hybrid results compared to MPP with PTT(T) and PTT(I), the sparse result is very low. Therefore, we adopt the MPP(PTT(T) PTT(I)) used in this paper. Furthermore, well-performing sparse results provide the possibility of a three-stage pipeline.

D Results of RGP Prompt Variations

Here we provide some other RGP prompt variations:

RGP₁:

- text-to-image: Given a candidate and a query, predict whether the candidate includes an answer to the query by producing either 'Yes' or 'No'. \n Candidate: `<image>`\n Query: `<sent>`\n Does the candidate answer the query? Answer:
- image-to-text: Given a candidate and a query, predict whether the candidate includes an answer to the query by producing either 'Yes' or 'No'. \n Candidate: `<sent>`\n Query: `<image>`\n Does the candidate answer the query? Answer:

RGP₂:

- text-to-image: Query: `<sent>`\n Candidate: `<image>`\n Does the candidate answer the query? Answer 'Yes' or 'No'. Answer:
- image-to-text: Query: `<image>`\n Candidate: `<sent>`\n Does the candidate answer the query? Answer 'Yes' or 'No'. Answer:

Method		MSCOCO						Flickr30K					
		t2i			i2t			t2i			i2t		
		r@1	r@5	r@10	r@1	r@5	r@10	r@1	r@5	r@10	r@1	r@5	r@10
LLaVA-Next-LLaMA-8B	Dense	34.6	60.5	71.2	43.3	68.6	78.2	60.1	83.7	90.3	72.4	90.6	95.6
	SSR(T) PTT(I) Sparse	17.2	37.5	48.3	5.5	18.8	29.8	30.4	56.7	68.4	12.3	32.4	44.6
	SSR(T) PTT(I) Hybrid	36.5	63.1	73.2	44.1	70.0	78.6	62.1	85.8	91.5	73.0	90.6	95.8
	PTT(T) PTT(I) Sparse	26.0	50.9	62.3	32.4	58.3	69.7	47.6	74.9	83.2	57.5	82.6	89.7
	PTT(T) PTT(I) Hybrid	34.4	60.3	71.0	43.0	68.2	78.2	59.8	83.3	90.3	72.0	90.4	95.2
	MPP(SSR(T) PTT(I) Sparse	18.3	38.7	49.2	7.4	21.7	31.1	32.2	57.6	67.1	16.6	36.9	49.3
	MPP(SSR(T) PTT(I) Hybrid	38.6	65.1	75.3	<u>44.8</u>	<u>70.1</u>	<u>79.0</u>	64.6	87.4	92.3	<u>73.2</u>	<u>91.3</u>	<u>95.7</u>
	MPP(PTT(T) PTT(I) Sparse	25.8	51.4	63.3	29.6	57.7	70.4	52.5	79.4	86.4	55.6	83.6	90.4
	MPP(PTT(T) PTT(I) Hybrid	<u>37.3</u>	<u>63.7</u>	<u>74.1</u>	46.0	71.2	80.5	64.4	86.7	91.9	74.1	91.9	96.1
LLaVA-Next-Mistral-7B	Dense	34.3	59.7	70.3	42.7	67.2	77.1	60.5	83.0	89.0	72.3	90.9	94.6
	SSR(T) PTT(I) Sparse	17.1	36.3	46.7	5.5	18.3	27.3	34.8	59.1	68.5	15.3	36.4	48.4
	SSR(T) PTT(I) Hybrid	35.9	61.6	72.0	42.9	67.7	77.7	63.1	85.4	90.6	73.7	91.0	95.1
	PTT(T) PTT(I) Sparse	25.4	49.2	60.7	31.3	56.0	66.9	48.6	74.8	81.9	59.4	83.4	90.7
	PTT(T) PTT(I) Hybrid	34.1	59.6	70.2	42.2	67.1	76.7	60.2	82.8	88.7	71.9	90.6	94.9
	MPP(SSR(T) PTT(I) Sparse	18.5	37.6	47.5	12.4	30.6	41.8	34.2	58.9	68.0	28.2	51.9	63.3
	MPP(SSR(T) PTT(I) Hybrid	38.7	65.0	75.1	<u>43.8</u>	<u>68.8</u>	<u>79.2</u>	<u>65.2</u>	<u>87.2</u>	<u>92.1</u>	<u>73.6</u>	<u>90.6</u>	<u>95.6</u>
	MPP(PTT(T) PTT(I) Sparse	30.5	56.1	67.5	39.2	66.8	77.7	57.0	82.3	89.1	68.4	89.6	94.5
	MPP(PTT(T) PTT(I) Hybrid	<u>38.1</u>	<u>64.5</u>	<u>74.8</u>	46.1	72.4	81.5	65.5	87.5	92.9	75.4	93.9	97.3

Table 5: The dense, sparse, and hybrid results of SSR, PTT, and MPP on MSCOCO and Flickr30K in a zero-shot setting.

We present the results of RGP variations in the Table D. In this table, we find that the RGP prompt template seems to be sensitive to LLaVA. With various prompt variations, RGP reranking method performs quite differently. An observation is that an RGP prompt can only achieve improvement on one of the text-to-image and image-to-text retrieval settings. Of course, this situation is suboptimal, and it is better to achieve increases under both retrieval settings at the same time. We believe this is an interesting research point to further optimize the zero-shot cross-modal retrieval results.

E Ethical Statement & Risks

This study complied with all applicable ethical guidelines. Formal ethical approval was not required as the research did not involve human participants or animal experiments.

This study introduces some open-source multi-modal large language models on HuggingFace and cross-modal retrieval datasets for evaluation. Considering this, we may violate the protocol to cite these works. It must be declared that we have no intention whatsoever to violate the citation protocol. If anyone finds contents that violate the protocol, the corresponding research will be modified.

As for private information, this study uses images containing person or company objects from datasets. However, these datasets are widely used

in cross-modal retrieval and text-based person retrieval tasks. Considering this, it is conceived that this study does not violate privacy laws or contain other offensive information. In addition, because of the use of open-source benchmarks, this study does not attend to any recruitment and payment.

When writing this paper, this study uses AI assistance for polishing, translation, and fixing language bugs. However, we ensure that our research is entirely the product of our own ideas and not generated by AI.

Method		MSCOCO		Flickr30K	
		t2i	i2t	t2i	i2t
LLaVA-Next-LLaMA-8B	Dense	34.6	43.3	60.1	72.4
	MPP Sparse	25.8	29.6	52.5	55.6
	MPP Hybrid	37.3	46.0	64.4	74.1
	MPP+RGP ₁	47.6	32.3	75.8	55.4
	MPP+RGP ₂	14.2	49.0	21.9	77.3
	MPP+RGP	14.1	50.7	20.3	79.2
LLaVA-Next-Mistral-7B	Dense	34.3	42.7	60.5	72.3
	MPP Sparse	30.5	39.2	57.0	68.4
	MPP Hybrid	38.1	46.1	65.5	75.4
	MPP+RGP ₁	48.4	29.3	76.9	53.1
	MPP+RGP ₂	12.8	49.8	18.1	81.4
	MPP+RGP	12.7	49.4	19.1	80.0

Table 6: The r@1 reranking results of different RGP variations.