

# MCA: Modality Composition Awareness for Robust Composed Multimodal Retrieval

Anonymous ACL submission

## Abstract

Multimodal retrieval, which seeks to retrieve relevant content across modalities such as text or image, supports applications from AI search to contents production. Despite the success of separate-encoder approaches like CLIP aligning modality-specific embeddings with contrastive learning, recent multimodal large language models (MLLMs) enable a unified encoder that directly processes composed inputs. While flexible and advanced, we identify that unified encoders trained with conventional contrastive learning are prone to learn modality shortcut, leading to poor robustness under distribution shifts. We propose a modality composition awareness framework to mitigate this issue. Concretely, it consists of a preference loss enforces multimodal embeddings to outperform their unimodal counterparts, and a composition regularization objective aligns multimodal embeddings with prototypes composed from its unimodal parts. These objectives explicitly model structural relationships between the composed representation and its unimodal counterparts. Experiments on various benchmarks show gains in out-of-distribution retrieval, highlighting modality composition awareness as a effective principle for robust composed multimodal retrieval when utilizing MLLMs as the unified encoder.

## 1 Introduction

Multimodal retrieval, which aims to retrieve semantically relevant contents across multiple modalities such as text, image and audio, is a fundamental task in various information fields. Its applications span a wide range of domains, such as text-vision retrieval (Huynh et al., 2025; Wang et al., 2021), music retrieval (Doh et al., 2023), product search (Goenka et al., 2022; Zhu et al., 2024b), and multimodal retrieval-augmented generation (Yasunaga et al., 2023; Ghosh et al., 2024; Yang et al., 2024; Jeong et al., 2025). The core ability of multimodal retrieval is to represent multimodal inputs

in a shared and comparable embedding space. A prevailing approach to this problem is to adopt unimodal encoders and align the encoded embeddings through contrastive learning (CL). Models following this separate-encoder paradigm, such as CLIP (Radford et al., 2021) and CLAP (Elizalde et al., 2023), have demonstrated the effectiveness of CL, achieving strong performance across various multimodal retrieval tasks. On the other hand, with the rapid development of multimodal large language models (MLLMs) (Alayrac et al., 2022; Li et al., 2023; Bai et al., 2023; Chu et al., 2023; Liu et al., 2023a; Chen et al., 2024), there has been growing interest in employing MLLMs as encoders for multimodal retrieval (Jiang et al., 2024; Zhang et al., 2024b; Huang et al., 2025; Jiang et al., 2025). Unlike separate-encoder frameworks, MLLMs are capable of processing inputs from different modalities, as well as their compositions, within a unified architecture, providing several advantages such as a shared semantic space across modalities, flexible handling of composed queries and documents, e.g., text+image, and the ability to leverage powerful pretrained representations including rich language knowledge and multimodal understanding.

However, this flexibility also makes the model more prone to **modality shortcut learning**. Since all modalities are jointly processed within a shared encoder, the training loss can be minimized by over-relying on the stronger modality signal, while ignoring the complementary one. An example in Figure 1 illustrates that the model suffers from modality shortcuts when processing highly similar images, while ignoring the textual instruction: “put up convertible roof, remove snow, place SUV standing on flat asphalt”. Hence, the architecture shift from separate-encoder to unified-encoder makes the direct application of conventional CL objective to unified MLLM encoders limiting to robustness. Furthermore, the modality shortcut issue naturally extends to scenarios involving multiple modalities,



Figure 1: Illustration of modality shortcut problem in multimodal retrieval. Given the inputs consists of an image and textual instruction, **left**: baseline model trained with conventional CL loss retrieve a visually similar image, focusing too much on vision modality; **right**: model trained with modality composition awareness retrieves the relevant image following the instruction.

particularly as recent work aims to unify various modalities into a single framework (Girdhar et al., 2023; Zhu et al., 2024a; Xu et al., 2025). As the number of modalities increases, so does the risk of the model collapsing to rely on a single dominant modality, making explicit modality composition-aware constraints crucial for robust generalization.

We argue that **modality composition** is the key to mitigate this issue. In this paper, we propose a *modality composition awareness* (MCA) framework that modeling the structural relationship between multimodal and unimodal representations, as shown in Figure 2. MCA consists of two complementary objectives from the preference and consistency perspective, respectively. First, a *preference loss* enforces that the embeddings of a multimodal composition should be more discriminative than any of its unimodal counterparts, thereby discouraging modality shortcut learning. Second, a *composition regularization* objective encourages the consistency between the composed embedding produced by the unified encoder and a compositional prototype constructed from its unimodal embeddings, ensuring that composed representations remain grounded in their constituent modalities. Together, these objectives explicitly model the structural relationship between multimodal and unimodal inputs, leading to more robust representation.

Figure 5b visualizes the representation learned with MCA. It demonstrates that composed queries, unimodal queries and targets have a clearer boundary, showing that MCA reduces shortcut reliance.

To comprehensive evaluate MCA, we conduct extensive experiments on both in-domain (IND) and OOD benchmarks, covering retrieval and grounding tasks. The results demonstrate the MCA improves robustness with OOD improvements under distribution shifts while maintaining IND performance. Through ablations, we show that both the preference and regularization components contribute complementary benefits. In addition, we also highlight the interaction between input richness and the strength of our proposed losses. Those results suggest that explicitly modeling modality composition could be a general principle with broader implications for the unified multimodal retrieval using MLLMs.

## 2 Related Work

**Multimodal Retrieval** Classical multimodal retrieval typically adopts separate-encoder architectures, learning a text encoder and an image/audio/video encoder whose outputs are aligned in a share space via CL loss. This line of work has achieved strong performance at scale and efficient retrieval (Radford et al., 2021; Li et al., 2022; Zhai et al., 2023). Beyond unimodal inputs, composed retrieval focuses inputs with multiple modalities, such as text+image or text+audio, which better capture use intent in real-world scenarios. This setting has been studied in contexts like text-guided image retrieval (Yu et al., 2016), fashion production search (Wu et al., 2021) and multimodal modification tasks (Liu et al., 2021b). Separate-encoder approaches are principally designed for unimodal queries. When the query or candidate is a composition of modalities, an extra fusion module (Wei et al., 2024; Huynh et al., 2025) has to be trained to fuse the modalities. While these methods demonstrate the feasibility of handling composed queries, they are typically build upon conventional CL loss and do not model the cross-modal interaction between the modalities, leading to a performance sacrifice (Huang et al., 2025).

**MLLMs for Retrieval** Recent MLLMs extend pretrained LLMs with multimodal adapters or cross-modal modules, enabling an integrated architecture to process and understand over text images, audios and videos (Alayrac et al., 2022; Li

et al., 2023; Bai et al., 2023; Chu et al., 2023; Liu et al., 2023a; Chen et al., 2024). Given the unified processing of multiple modalities and rich pre-learned knowledge, MLLMs are increasingly applied to retrieval tasks, especially composed retrieval. Most approaches adapt the unified encoder directly with contrastive learning for embedding-based retrieval (Jiang et al., 2024; Zhang et al., 2024b; Huang et al., 2025; Jiang et al., 2025), while others employ MLLMs as re-rankers for the later re-ranking (Lin et al., 2024). In this work, we study the embedding-based approaches. An advantage of MLLMs for retrieval is the capability of jointly understanding contents in multiple modalities with a unified encoder. However, adopting the unified encoder in contrastive learning could also lead to shortcut learning (Geirhos et al., 2020; Wu et al., 2022). As a result, these recent methods generally inherit conventional contrastive loss and therefore remain vulnerable to modality shortcut when handling composed inputs. Orthogonal to these previous works, we firstly study the modality shortcut problem in this setting and our proposed MCA framework mitigates issue by introducing modality composition-aware objectives that can be easily integrated with MLLM-based retriever to enhance the robustness.

### 3 Methodology

#### 3.1 Preliminary: MLLM Embedding

Multimodal retrieval aims to search for target documents by a given query. The embedding model, which is central to multimodal retrieval and the focus of our approach, is introduced in this section with the encoding process and training objectives.

**Unified Multimodal Encoder** Unlike conventional separate-encoder methods such as CLIP (Radford et al., 2021; Zhai et al., 2023), we study the approaches (Lin et al., 2024; Jiang et al., 2024; Zhang et al., 2024b; Huang et al., 2025; Jiang et al., 2025) using a *unified encoder*  $f_\theta(\cdot)$ , such as a MLLM. Formally, given an input  $x$  that could comprises multiple modalities,  $f_\theta(\cdot)$  encodes the input into a shared space as:  $\mathbf{h} = f_\theta(x)$ , where  $\mathbf{h} \in \mathbb{R}^d$  is the embedding in  $d$  dimensions. In addition,  $f_\theta(\cdot)$  includes the necessary processes such as tokenization and pooling, and we omit those details because it does not affect the modeling.

**Contrastive Learning** Then we can obtain the embedding of a query or document by the uni-

fied encoder. Prior works usually adopt the conventional CL loss to align the representations of queries and documents. To achieve this, we first define a function  $\text{sim}_\theta(\cdot, \cdot)$  to measure the similarity between two inputs  $x$  and  $y$  as follows,

$$\text{sim}_\theta(x, y) := \frac{\text{sim}[f_\theta(x), f_\theta(y)]}{\tau}, \quad (1)$$

where  $\tau$  is the contrastive temperature parameter. As  $f_\theta(\cdot)$  is a unified encoder,  $x$  and  $y$  can be any modalities or the composition of them. Next, the CL objective (Radford et al., 2021) is optimized for aligning the query and positive document. Given the dataset  $\mathcal{D}$  that is constructed by queries and documents, in which each query is paired with its positive document, the CL loss can be formalized as follows,

$$\mathcal{L}^{\text{CL}} = \mathbb{E}_{x, y^+, Y \sim \mathcal{D}} \left[ -\log \frac{e^{\text{sim}_\theta(x, y^+)}}{e^{\text{sim}_\theta(x, y^+)} + \sum_{y \in Y} e^{\text{sim}_\theta(x, y)}} \right] \quad (2)$$

where  $x$  is the query,  $y^+$  is the paired positive document and  $Y$  is the batch of negative documents. By doing so, all examples are encoded into a shared space, within which the query is expected to be pulled closer to the positive document but pushed away from negative ones.

#### 3.2 MCA: Multimodal Embedding with Modality Composition Awareness

Although the conventional CL has demonstrated remarkable effectiveness in retrieval (Radford et al., 2021; Baldrati et al., 2022; Liu et al., 2023c), it is inherently designed under the separate-encoder paradigm, where each modality is encoded independently and alignment is enforced at the representation level. With the rapid development of MLLMs, the utilization of unified encoder for retrieval is becoming promising (Jiang et al., 2024; Zhang et al., 2024b; Huang et al., 2025; Jiang et al., 2025). However, those MLLM-based approaches still use conventional CL as the training objective, which is less adequate under composed scenario, where queries or documents may consist of multiple modalities. Such approaches enable models to naturally handle composed inputs within a unified architecture, which could lead to modality shortcut learning as we introduced in §1. Hence, a robust multimodal retrieval paradigm, particularly one that leverages MLLMs and adapts to a wider range of composed scenarios, requires not only aligning across modalities, but also explicitly modeling the *modality composition*.

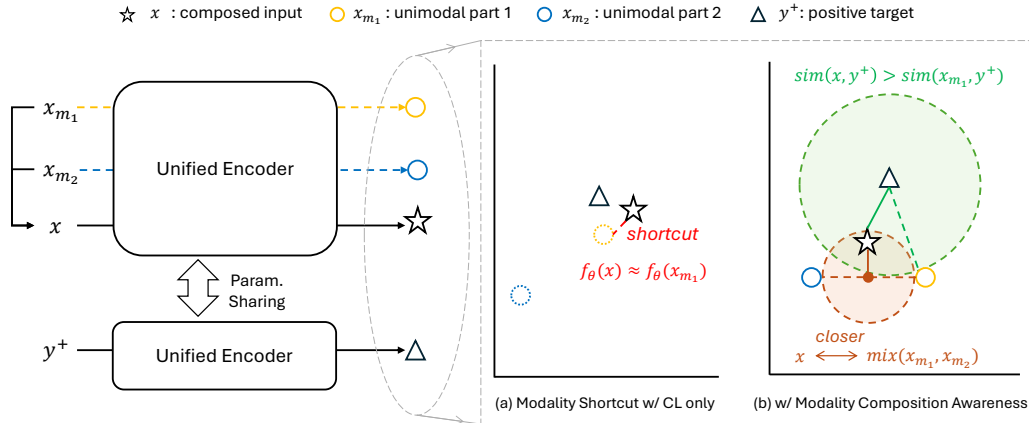


Figure 2: Modality Composition Awareness (MCA, §3.2). Unimodal parts are shown as dotted circles when not explicitly modeled. (a) Training with vanilla CL: the composed embedding can be close to the target but still align disproportionately with one unimodal part, leading to modality shortcuts. (b) Training with CL and MCA: the composed embedding is explicitly constrained to be closer to the target than any of its unimodal parts by MCP (§3.2.1), and anchored to a compositional prototype mixed (Equation 5) from unimodal embeddings by MCR (§3.2.2) to mitigate modality shortcut.

Consequently, we propose *Modality Composition Awareness* to mitigate this issue. It consists of two components: (1) *Modality Composition Preference* (MCP, §3.2.1) and (2) *Modality Composition Regularization* (MCR, §3.2.2). These two training objectives are illustrated in Figure 2 and will be introduced in the following sections. We first formalize the notation of modalities as  $\mathcal{M}^1$ . Given a composed input  $x$ , we define the set of its unimodal counterparts as follows,

$$\mathcal{U}(x) := \{x_m \mid m \in \mathcal{M}_x\}, \quad (3)$$

where  $\mathcal{M}_x \subset \mathcal{M}$  is the modality set of  $x$ , and  $x_m$  denotes its unimodal part for modality  $m$ . For example, if  $x$  comprises text and image modalities, then  $\mathcal{U}(x) = \{x_{\text{text}}, x_{\text{image}}\}$ .

The modality shortcut problem can be formalized as  $f_\theta(x) \approx f_\theta(x_m)$  regardless of other modalities after the optimization. In such a case the model only learns unimodal representation even composed inputs are provided. Formally, this key characteristic of modality shortcut is that the representation satisfies  $f_\theta((x_m, x_{\text{other}})) \approx f_\theta((x_m, x'_{\text{other}}))$ , meaning the model ignores variations in the other modality. As a result, when critical information resides in other modality under OOD conditions, the generalization performance could degrade. Note that modality shortcuts only arise when at least one side of the input pair is a composed input. If both the query and document are unimodal, the

<sup>1</sup>In this paper, we use text and image as available modalities, i.e.,  $\mathcal{M} = \{\text{text}, \text{image}\}$ , but the definition can be extended to more modalities.

task reduces to standard cross-modal retrieval with minimal shortcut risk. Therefore, the following proposed losses are only applied when either the query or document is composed.

### 3.2.1 MCP: Modality Composition Preference.

We propose a preference loss that explicitly enforces composed inputs to be more discriminative than their unimodal counterparts. Concretely, given a pair where either side is a composed input, we compute similarities between the composed embedding and candidate targets, as well as between its corresponding unimodal parts and the same target. The MCP loss is formalized as,

$$\mathcal{L}^{\text{MCP}} = \mathbb{E}_{x, y^+ \sim \mathcal{D}} \left[ - \sum_{x_m \in \mathcal{U}(x)} \log \frac{e^{\text{sim}_\theta(x, y^+)}}{e^{\text{sim}_\theta(x_m, y^+)}} \right], \quad (4)$$

where  $x$  and  $y^+$  are the paired query-document sampled from dataset  $\mathcal{D}$ . The MCP loss encourages the composed similarity higher than the unimodal similarities, ensuring that the model leverages complementary signals from multiple modalities rather than relying on one dominant modality. In other words, the loss formulates a preference-style constraint that whenever a composed input and its unimodal counterparts compete on the same retrieval task, the composed representation should be preferred. For symmetry, we compute the loss in two directions, i.e.,  $x$  is the document and  $y^+$  is the paired query, ensuring the preference is enforced for both query and document representation.

### 3.2.2 MCR: Modality Composition Regularization

In addition to the preference constraint, we further propose MCR loss to encourage the consistency between the composed embedding and a simple prototype composition of its unimodal embeddings. The intuition is that the representation produced by the encoder for a multimodal input should not deviate arbitrarily from the semantic space spanned by its unimodal parts. To this end, we redefine the similarity function as follows,

$$\text{sim}_{\theta, \phi}^{\text{mix}}(x, \mathcal{U}(x')) := \frac{\text{sim}[f_{\theta}(x), \text{mix}_{\phi}(\{f_{\theta}(x'_m) | x'_m \in \mathcal{U}(x')\})]}{\tau}, \quad (5)$$

where  $x$  and  $x'$  are two inputs to compare.  $x$  is the one we want to regularize, and  $\mathcal{U}(x')$  denotes the set of unimodal components that we defined in Equation 3.  $\text{mix}(\cdot)$  is a mixer module we defined that aggregates multiple unimodal embeddings into a composed prototype  $\mathbf{h}' \in \mathbb{R}^d$  with the same output dimension as  $f_{\theta}(\cdot)$ . Thus, the redefined similarity function measures the similarities between the composed prototype derived from  $x'$  and the composed embedding of  $x$ . In addition, to ensure that the regularization does not dominate learning, the mixer is deliberately kept simple, such as mean pooling or gated fusion, so that the encoder must learn to align composed inputs with the composed prototype rather than overfitting to the mixer.

The MCR loss is then defined as a contrastive objective, where the composed embedding is pulled closer to its mixed prototype than to other in-batch negatives,

$$\mathcal{L}^{\text{MCR}} = \mathbb{E}_{x, X \sim \mathcal{D}} \left[ -\log \frac{e^{\text{sim}_{\theta, \phi}^{\text{mix}}(x, \mathcal{U}(x))}}{e^{\text{sim}_{\theta, \phi}^{\text{mix}}(x, \mathcal{U}(x))} + \sum_{x' \in X} e^{\text{sim}_{\theta, \phi}^{\text{mix}}(x, \mathcal{U}(x'))}} \right], \quad (6)$$

where  $x$  is a sampled query or document, and  $X$  is a batch of composed inputs that could be either queries or documents for symmetry. Intuitively, the MCR loss enforces that the composed embedding remains anchored to the space formed by its constituent unimodal embeddings, thereby reducing the risk of spurious modality shortcut learning. Finally, we combine the CL loss with the two proposed losses. The overall training objective for optimizing  $\theta$  and  $\phi^2$  is formulated as,

$$\mathcal{L} = \mathcal{L}^{\text{CL}} + \alpha \times \mathcal{L}^{\text{MCP}} + \beta \times \mathcal{L}^{\text{MCR}}, \quad (7)$$

where  $\alpha$  and  $\beta$  are weighting coefficients controlling the relative strength of auxiliary terms. The impact of weighting is discussed in §4.2.

<sup>2</sup> $\phi$  is optimized only when the mixer contains learnable parameters.

## 4 Experiments

To evaluate the proposed MCA, we conduct extensive experiments on multimodal retrieval. The details of the datasets and benchmarks are introduced in §A.2. The training settings are introduced in §A.3. For evaluation we consider both IND benchmarks, which share the same distribution as the training data, and OOD benchmarks that test the generalization to unseen compositions, domains and tasks. OOD performance directly evaluates the shortcut learning problem (Geirhos et al., 2020), because the ideal solution generalizes to OOD test sets, but a shortcut-driven model may only performs well on training and IND test sets, which is illustrated in Figure 6. The OOD benchmarks cover two types of tasks: 1) *OOD retrieval* that measures performance under distribution shifts such as domain and modality composition changes, and 2) *OOD grounding* measures cross-task generalization. Both types of tasks require robust modality composition. Detailed evaluation settings are introduced in §A.4.

### 4.1 Main Results

#### 4.1.1 Overall Results

Our results in Table 1 demonstrate the superior robustness of MCA in OOD retrieval and zero-shot grounding scenarios. Note that some baselines differ in training data and backbone models. For instance, VLM2Vec was explicitly trained on Grounding tasks so it is not zero-shot performance. The completely fair comparison is between MCA and our reproduced Qwen2-VL, which share identical experimental settings.

In general, MLLM-based methods outperform dual-encoder models because their unified encoder enables deeper cross-modal understanding, leading to better performance on all scenarios. However, MLLMs can suffer from modality shortcut issues, where the model relies too heavily on one modality as we introduced in §1. Our MCA is designed to address this limitation and further enhance robustness. Specifically, MCA consistently achieves the highest average in OOD retrieval at 57.4, which is 3.2 points higher than the best baseline, showing a strong ability to generalize to unseen domains. In zero-shot grounding, MCA also leads with an average score of 61.0, outperforming the best baseline by 3.1 points, highlighting its advanced compositional understanding and effective transfer to tasks without explicit training. These results il-

Method	In-domain Retrieval								OOD Retrieval				Zero-shot Grounding					
	VisDial	CIRR	VisualNews_12i	VisualNews_12t	MSCOCO_12i	MSCOCO_12t	NIGHTS	WebQA	Avg.	OVEN	FashionIQ	EDIS	Avg.	MSCOCO	Visual7W-P	RefCOCO	RefCOCO-M	Avg.
<i>Dual-Encoder</i>																		
CLIP (Radford et al., 2021)	30.7	12.6	78.9	79.6	59.5	57.7	60.4	67.5	55.8	41.1	11.4	81.0	44.5	33.8	55.1	56.9	61.3	51.7
OpenCLIP (Cherti et al., 2023)	25.4	15.4	74.0	78.0	63.6	62.1	66.1	62.1	55.8	45.0	13.8	77.5	45.4	34.5	56.3	54.2	68.3	53.3
SigLIP (Zhai et al., 2023)	21.5	15.1	51.0	52.4	58.3	55.0	62.9	58.1	46.7	56.0	20.1	23.6	33.2	46.4	70.1	70.8	50.8	59.5
BLIP2 (Li et al., 2023)	18.0	9.8	48.1	13.5	53.7	20.3	56.5	55.4	39.5	39.3	9.3	54.4	34.4	28.9	52.0	47.4	59.5	46.9
MagicLens (Zhang et al., 2024a)	24.8	39.1	50.7	21.1	54.1	40.0	58.1	43.0	41.3	11.2	1.6	62.6	25.1	22.1	23.4	22.8	35.6	25.9
<i>Unified MLLM Encoder</i>																		
E5-V (Jiang et al., 2024)	9.2	6.1	13.5	8.1	20.7	14.0	4.2	17.7	11.6	5.9	2.8	26.8	11.8	10.8	14.3	11.9	38.9	18.9
*VLM2Vec (Jiang et al., 2025)	80.9	49.9	75.4	80.0	75.7	73.1	65.5	87.6	73.5	56.5	16.2	87.8	53.5	-	-	-	-	-
†Qwen2-VL (Wang et al., 2024)	81.3	51.6	75.7	78.6	73.9	71.9	65.6	86.8	73.2	63.0	15.6	84.0	54.2	41.0	60.7	60.3	69.6	57.9
MCA (Ours)	80.3	50.5	76.3	78.7	74.6	72.1	66.1	86.8	73.2 ( $\Delta 0.0$ )	66.7	19.3	86.1	57.4 ( $\Delta + 3.2$ )	42.7	65.4	65.0	70.9	61.0 ( $\Delta + 3.1$ )

Table 1: Overall results. \* indicates the baseline that explicitly trained on zero-shot grounding datasets so the performance is not comparable. † indicates the baselines reproduced with identical settings, backbone LLMs, and training data, and  $\Delta$  indicates the difference compared to the completely comparable baselines.

413 illustrate that modality composition awareness not  
 414 only boosts robustness to distribution shifts but  
 415 also enables the model to excel in zero-shot scenarios.  
 416 Importantly, MCA maintains competitive  
 417 in-domain retrieval performance, matching the best  
 418 baseline at 73.2. This balance confirms that the  
 419 gains in robustness do not come at the expense of  
 420 accuracy on familiar data. More qualitative analysis  
 421 is introduced in §A.1.

#### 4.1.2 Loss Breakdown Analysis

422 As MCA consists of two optimization targets, we  
 423 further analyze the contribution of the two auxiliary  
 424 losses separately. As shown in Figure 3, both MCP  
 425 and MCR individually lead to consistent positive  
 426 gains on OOD benchmarks. Notably, their standalone  
 427 improvements are relatively smaller. When  
 428 combined, the two losses yield substantially larger  
 429 gains of +5.9% on OOD retrieval and 5.4% on zero-  
 430 shot grounding, over the model trained with only  
 431 contrastive learning, demonstrating the necessity  
 432 of applying them together. This observation aligns  
 433 with our design intuition in §3.2 that MCP encourages  
 434 preference against unimodal shortcuts while MCR  
 435 enforces structural consistency, and only their joint  
 436 application forms the ideal constraint for robust  
 437 modality composition, as illustrated in Figure 2.  
 438 On IND benchmarks, all model variants converge to  
 439 almost identical accuracy levels. Note that the  
 440 proposed MCA introduces auxiliary losses to enhance  
 441 model robustness. Therefore, similar in-domain  
 442 performance compared to the baseline is expected  
 443 and indicates that the conventional CL is sufficient  
 444 to make the model well-converged on IND data.  
 445 The main focus should be on improve-

ments in OOD and zero-shot scenarios.

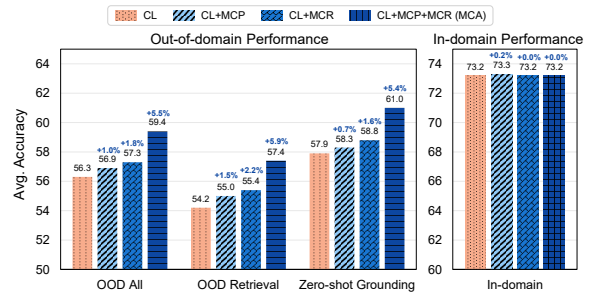


Figure 3: Breakdown of MCA through loss ablation.

#### 4.1.3 Convergence Analysis

447 To assess the impact of the proposed MCA loss on  
 448 training dynamics, we perform a convergence analysis  
 449 by comparing the loss curves and validation  
 450 performance of both the baseline and the proposed  
 451 method. This ensures that any differences in OOD  
 452 and zero-shot results are not due to incomplete op-  
 453 timization or training instability.  
 454  
 455

456 **Loss Curve** We first examine whether the pro-  
 457 posed objectives affect the optimization dynamics  
 458 of CL loss. Figure 4a compares the curves of the  
 459 CL loss values with and without MCA during the  
 460 training. Both exhibit smooth and monotonic de-  
 461 crease, indicating that introducing MCA does not  
 462 hinder the convergence of the main contrastive ob-  
 463 jective. We then analyze the internal loss dynamics  
 464 of MCA in Figure 4b. The CL, proposed MCP  
 465 and MCR losses all converge stably. We note that  
 466 both MCR and MCP losses exhibit noticeable fluc-  
 467 tuations at the beginning of training. The possi-  
 468 ble reason could be unstable alignment between

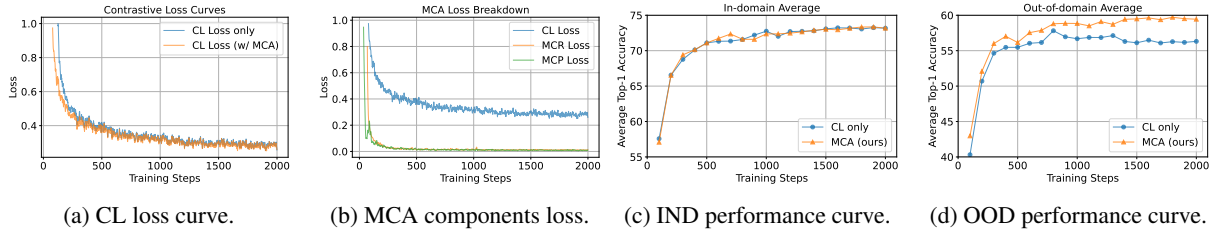


Figure 4: Convergence analysis of MCA. For a fair comparison, the y-axis range for the loss is set to 0–1, and the range for benchmark performance is consistently set to 20.

unimodal and multimodal embeddings in the early stage. Once the representation space becomes more coherent, both losses stabilize at small values.

**Performance Curve** To further track model performance during training, we measure average retrieval accuracy every 100 training steps. Figure 4c shows the IND results, where performance improves steadily and eventually converges to a level comparable to the baseline. In contrast, Figure 4d presents the OOD results, where MCA consistently outperforms the baseline. Moreover, the performance curves reveal that the benefits of MCA appear after a few hundred of steps. While IND accuracy converges similarly for models trained with and without MCA, the OOD gap emerges within 500 steps and continues to widen throughout training. This suggests that MCA enhances generalization by mitigating modality shortcuts. Above results show that MCA converges smoothly, with auxiliary losses behaving as a regularizer, and leads to consistent performance improvements on OOD benchmark where the robustness is most critical.

## 4.2 Input Resolution and MCA Weighting

To better understand the behavior of MCA, we investigate its sensitivity to two critical factors that are relevant: the resolution of image inputs and the weighting of proposed losses. These two factors directly affect the balance between modalities. Specifically, when input information is degraded, the regularization strength might become crucial prevent shortcuts, while richer input information may require lighter regularization to avoid overconstraining the model. Table 2 summarizes the results across three input image resolution and different weighting values. For simplicity, we keep the coefficients of MCP and MCR loss equal, and very their shared ratio to the CL loss.

Overall, we observe that higher resolution consistently yield stronger performance as expected, since richer visual detail provides more reliable vi-

MCA Ratio	IND Avg.	OOD Avg.	OOD Retrieval Avg.	Zero-shot Grounding Avg.
<i>Low resolution (128 × 128)</i>				
0	69.7	46.7	42.5	49.9
0.01	69.5 (-0.3%)	44.0 (-5.8%)	37.6 (-11.6%)	48.8 (-2.2%)
0.10	69.3 (-0.6%)	49.5 (+5.9%)	46.8 (+10.1%)	51.4 (+3.0%)
1.00	69.5 (-0.3%)	50.9 (+8.9%)	48.1 (+13.1%)	52.9 (+6.0%)
<i>Mid resolution (672 × 672)</i>				
0	72.0	53.9	49.9	56.8
0.01	71.3 (-1.0%)	54.7 (+1.4%)	53.5 (+7.2%)	55.5 (-2.3%)
0.10	71.5 (-0.7%)	54.9 (+1.8%)	51.3 (+2.8%)	57.6 (+1.4%)
1.00	70.9 (-1.6%)	55.8 (+3.5%)	53.9 (+8.0%)	57.3 (+0.8%)
<i>High resolution (1344 × 1344)</i>				
0	73.2	56.3	54.2	57.9
0.01	73.2 (+0.0%)	59.4 (+5.5%)	57.4 (+5.9%)	61.0 (+5.4%)
0.10	73.0 (-0.3%)	56.6 (+0.5%)	53.4 (-1.5%)	59.1 (+2.0%)
1.00	73.2 (+0.0%)	58.3 (+3.5%)	57.0 (+5.1%)	59.3 (+2.4%)

Table 2: Impact of weighting MCA on varying input image resolution. Darker colors denote larger deviations from the vanilla CL.

sual grounding. Interestingly, the relative gain from MCA grows larger as the resolution decreases. In lower resolution setting, where visual information is degraded and the risk of modality shortcuts is higher, MCA provides significant improvements by enforcing the use of complementary textual cues. This pattern suggests that MCA is potentially more effective in scenarios with *imbalanced modality* quality, where it helps the model resist collapsing onto a single modality and instead maintain faithful composition. Nevertheless, we also observe notable exceptions. In particular, under low-resolution inputs with very small weights 0.01, the model performs significantly worse than the vanilla CL baseline on OOD benchmarks. The possible reason is, the additional losses introduce small but inconsistent gradients, which are amplified by the noisy low-resolution image representations. As a result, the model is shifted away from the vanilla CL optimum without receiving sufficient corrective signal, causing a collapse even more severe than the baseline. In the mid-resolution setting, we also observe a modest degradation on IND benchmarks. This can

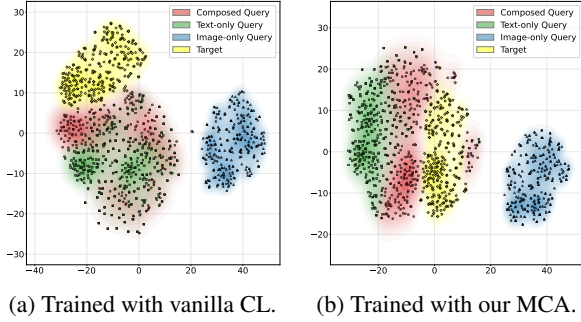


Figure 5: Visualizations of queries and targets.

be interpreted as the classic bias-variance trade-off. This degradation remains small, and the overall benefit of MCA becomes more evident when robustness is prioritized.

In addition, the results generally highlight an interaction between input richness and the strength of MCA regularization. When visual inputs are degraded, the risk of modality shortcuts is amplified, and stronger weighting is required for MCA to be effective. In contrast, when inputs provide richer information, the model relies less on shortcuts and lighter weighting suffices to stabilize the training. This trend suggests a practical guideline: the richer the input information, the smaller the weighting can lead to a better generalization, whereas weaker inputs demand stronger regularization to preserve robust modality composition.

### 4.3 Embedding Visualization

In Figure 5a, we visualize unimodal queries and targets randomly sampled from the CIRRD dataset trained with vanilla CL and MCA, respectively. Details of the t-SNE implementation are provided in §A.5. The visualization shows that, under vanilla CL, composed queries (red) are not well separated and frequently cluster close to text-only queries. This suggests that the model tends to exploit modality shortcuts rather than learning robust representations for composed queries. In contrast, with MCA training, composed queries and text-only queries are clearly separated, indicating a more meaningful and robust representation distribution.

### 4.4 Mixer Implementation

To examine the effect of different implementations of the mixer that defined in Equation 5, we compare several simple modules, including the gated fusion module by default, the non-parametric mean pooling and the multimodal factorized bilinear pooling (MFB) (Yu et al., 2017). All mixer choices introduce only a negligible number of additional

parameters, ensuring that the observed gains are not due to increased model capacity. As shown in Table 3, the choice of mixer has little effect on IND benchmarks, suggesting that MCA does not alter the standard retrieval behavior. However, the improvements differ in OOD settings. Both gated fusion and MFB variants yield improvements over the CL-only baseline, demonstrating the explicitly anchoring composed embeddings to the prototype from their unimodal parts is beneficial. In contrast, the mean pooling provides only marginal improvements and even slightly degrades performance on cross-task grounding. This indicates that overly simple designs may be insufficient in practice to construct the prototype from unimodal embeddings. Among them, the default gated fusion achieves the best performance. Above experiments suggest that the effectiveness of MCR does not heavily rely on the specific choice of the mixer, but rather on the general principle of enforcing compositional consistency, although the gain could vary depending on the mixer implementation.

Mixer	In-domain Avg.	OOD Avg.	OOD Retrieval Avg.	Zero-shot Grounding Avg.
Vanilla CL	73.2	56.3	54.2	57.9
Mean Pooling	73.3 (+0.1%)	56.5 (+0.3%)	54.8 (+1.1%)	57.8 (-0.2%)
Gated Fusion	73.2 (+0.0%)	59.4 (+5.5%)	57.4 (+5.9%)	61.0 (+5.4%)
MFB	73.0 (-0.2%)	57.7 (+2.4%)	54.6 (+0.7%)	60.0 (+3.6%)

Table 3: Impact of mixer. Darker colors denote larger deviations from the vanilla CL.

## 5 Conclusion

In this work, we introduced modality composition awareness (MCA) for robust multimodal retrieval. By explicitly modeling the relationship between unimodal and multimodal inputs, MCA incorporates two objectives: (1) modality composition preference, which discourages modality shortcuts by ensuring that composed representations are more discriminative than unimodal ones, and (2) modality composition regularization, which stabilizes embedding by aligning multimodal presentations with compositional prototypes. Extensive empirical experiments across IND and OOD benchmarks demonstrate that MCA achieves consistent gains under distribution shifts and unseen zero-shot tasks, while maintaining comparable IND accuracy. These results highlight MCA as an effective principle for mitigating modality shortcut problem when using MLLMs as unified encoder and improving generalization in multimodal retrieval.

## 613 Limitations

614 While our proposed modality composition aware-  
615 ness show consistent improvements across vari-  
616 ous OOD benchmarks, several limitations remain.  
617 First, the design of our losses assume that modality  
618 shortcuts primarily arise when composed inputs  
619 are present, and thus the framework is not directly  
620 applied to purely unimodal queries. This assump-  
621 tion may overlook shortcut phenomena that could  
622 still exist in separate-encoder approaches, although  
623 the risk is notably lower. Second, the experimental  
624 results still show that our proposed MCA could  
625 lead to performance degradation under particular  
626 weighting and input resolution setting. Although  
627 rare, this phenomenon suggests that it may require  
628 setting adjust to make the proposed loss work in  
629 practice. Third, in this work we only explore sim-  
630 ple mixer implementation for fair comparison, and  
631 extending the design to richer mixers is an inter-  
632 esting open direction. Lastly, although we have  
633 identified and studied the modality shortcut prob-  
634 lem for vision-text, due to the lack of composed  
635 training data and benchmarks, the generalization  
636 to broader modalities, e.g., audio, has not yet been  
637 validated. We expect the problem of modality short-  
638 cut to become even more prominent in the future  
639 with more available composed training data and  
640 benchmarks in MLLM-based retrieval.

## 641 Ethical Statement

642 This work focuses on improving multimodal repre-  
643 sentation learning by mitigating modality shortcut  
644 problem. Our experiments are conducted entirely  
645 on publicly available datasets and model check-  
646 points. Our method does not directly generate new  
647 contents but instead focuses on retrieval, as a re-  
648 sult, it poses minimal risk of misuse in generating  
649 harmful or deceptive contents. Nevertheless, we  
650 emphasize the importance of responsible use of the  
651 retrieval model because the MLLMs could have  
652 unexpected biases.

## 653 References

654 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,  
655 Antoine Miech, Iain Barr, Yana Hasson, Karel  
656 Lenc, Arthur Mensch, Katherine Millican, Malcolm  
657 Reynolds, and 1 others. 2022. Flamingo: a visual  
658 language model for few-shot learning. *Advances in  
659 neural information processing systems*, 35:23716–  
660 23736.

661 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,

Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
Huang, and 1 others. 2023. Qwen technical report.  
*arXiv preprint arXiv:2309.16609*.

Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and  
Alberto Del Bimbo. 2022. Effective conditioned and  
composed image retrieval combining clip-based fea-  
tures. In *Proceedings of the IEEE/CVF conference  
on computer vision and pattern recognition*, pages  
21466–21474.

Yingshan Chang, Mridu Narang, Hisami Suzuki, Gui-  
hong Cao, Jianfeng Gao, and Yonatan Bisk. 2022.  
Webqa: Multihop and multimodal qa. In *Proceed-  
ings of the IEEE/CVF conference on computer vision  
and pattern recognition*, pages 16495–16504.

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo  
Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,  
Xizhou Zhu, Lewei Lu, and 1 others. 2024. Internvl:  
Scaling up vision foundation models and aligning  
for generic visual-linguistic tasks. In *Proceedings of  
the IEEE/CVF conference on computer vision and  
pattern recognition*, pages 24185–24198.

Mehdi Cherti, Romain Beaumont, Ross Wightman,  
Mitchell Wortsman, Gabriel Ilharco, Cade Gordon,  
Christoph Schuhmann, Ludwig Schmidt, and Jenia  
Jitsev. 2023. Reproducible scaling laws for con-  
trastive language-image learning. In *Proceedings  
of the IEEE/CVF conference on computer vision and  
pattern recognition*, pages 2818–2829.

Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shil-  
iang Zhang, Zhijie Yan, Chang Zhou, and Jingren  
Zhou. 2023. Qwen-audio: Advancing universal  
audio understanding via unified large-scale audio-  
language models. *arXiv preprint arXiv:2311.07919*.

Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh,  
Deshraj Yadav, José MF Moura, Devi Parikh, and  
Dhruv Batra. 2017. Visual dialog. In *Proceedings of  
the IEEE conference on computer vision and pattern  
recognition*, pages 326–335.

SeungHeon Doh, Minz Won, Keunwoo Choi, and Juhan  
Nam. 2023. Toward universal text-to-music retrieval.  
In *ICASSP 2023-2023 IEEE International Confer-  
ence on Acoustics, Speech and Signal Processing  
(ICASSP)*, pages 1–5. IEEE.

Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Is-  
mail, and Huaming Wang. 2023. Clap learning  
audio concepts from natural language supervision.  
In *ICASSP 2023-2023 IEEE International Confer-  
ence on Acoustics, Speech and Signal Processing  
(ICASSP)*, pages 1–5. IEEE.

Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy  
Chai, Richard Zhang, Tali Dekel, and Phillip Isola.  
2023. Dreamsim: Learning new dimensions of hu-  
man visual similarity using synthetic data. In *Ad-  
vances in Neural Information Processing Systems*,  
volume 36, pages 50742–50768.

662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716

717	Robert Geirhos, Jörn-Henrik Jacobsen, Claudio	Training vision-language models for massive multi-	773
718	Michaelis, Richard Zemel, Wieland Brendel,	modal embedding tasks. In <i>The Thirteenth Interna-</i>	774
719	Matthias Bethge, and Felix A Wichmann. 2020.	<i>tional Conference on Learning Representations</i> .	775
720	Shortcut learning in deep neural networks. <i>Nature</i>		
721	<i>Machine Intelligence</i> , 2(11):665–673.		
722	Sreyan Ghosh, Sonal Kumar, Chandra Kiran Reddy	Sahar Kazemzadeh, Vicente Ordonez, Mark Matten,	776
723	Evuru, Ramani Duraiswami, and Dinesh Manocha.	and Tamara Berg. 2014. Referitgame: Referring to	777
724	2024. Recap: Retrieval-augmented audio captioning.	objects in photographs of natural scenes. In <i>Proceed-</i>	778
725	In <i>ICASSP 2024-2024 IEEE International Confer-</i>	<i>ings of the 2014 conference on empirical methods in</i>	779
726	<i>ence on Acoustics, Speech and Signal Processing</i>	<i>natural language processing (EMNLP)</i> , pages 787–	780
727	<i>(ICASSP)</i> , pages 1161–1165. IEEE.	798.	781
728	Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Man-	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	782
729	nat Singh, Kalyan Vasudev Alwala, Armand Joulin,	2023. Blip-2: Bootstrapping language-image pre-	783
730	and Ishan Misra. 2023. Imagebind: One embed-	training with frozen image encoders and large lan-	784
731	ding space to bind them all. In <i>Proceedings of the</i>	guage models. In <i>International conference on ma-</i>	785
732	<i>IEEE/CVF conference on computer vision and pat-</i>	<i>chine learning</i> , pages 19730–19742. PMLR.	786
733	<i>tern recognition</i> , pages 15180–15190.		
734	Sonam Goenka, Zhaoheng Zheng, Ayush Jaiswal,	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	787
735	Rakesh Chada, Yue Wu, Varsha Hedau, and Pradeep	Hoi. 2022. Blip: Bootstrapping language-image pre-	788
736	Natarajan. 2022. Fashionvlp: Vision language trans-	training for unified vision-language understanding	789
737	former for fashion retrieval with feedback. In <i>Pro-</i>	and generation. In <i>International conference on ma-</i>	790
738	<i>ceedings of the IEEE/CVF Conference on Computer</i>	<i>chine learning</i> , pages 12888–12900. PMLR.	791
739	<i>Vision and Pattern Recognition</i> , pages 14105–14115.		
740	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-	Sheng-Chieh Lin, Chankyu Lee, Mohammad Shoeybi,	792
741	Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu	Jimmy Lin, Bryan Catanzaro, and Wei Ping. 2024.	793
742	Chen. 2022. LoRA: Low-rank adaptation of large	Mm-embed: Universal multimodal retrieval with	794
743	language models. In <i>International Conference on</i>	multimodal llms. <i>arXiv preprint arXiv:2411.02571</i> .	795
744	<i>Learning Representations</i> .		
745	Hexiang Hu, Yi Luan, Yang Chen, Urvashi Khandel-	Tsung-Yi Lin, Michael Maire, Serge Belongie, James	796
746	wal, Mandar Joshi, Kenton Lee, Kristina Toutanova,	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,	797
747	and Ming-Wei Chang. 2023. Open-domain visual	and C Lawrence Zitnick. 2014. Microsoft coco:	798
748	entity recognition: Towards recognizing millions of	Common objects in context. In <i>European confer-</i>	799
749	wikipedia entities. In <i>Proceedings of the IEEE/CVF</i>	<i>ence on computer vision</i> , pages 740–755. Springer.	800
750	<i>International Conference on Computer Vision</i> , pages		
751	12065–12075.	Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente	801
752	Lang Huang, Qiyu Wu, Zhongtao Miao, and Toshihiko	Ordonez. 2021a. Visual news: Benchmark and chal-	802
753	Yamasaki. 2025. Joint fusion and encoding: Advanc-	enges in news image captioning. In <i>Proceedings</i>	803
754	ing multimodal retrieval from the ground up. <i>arXiv</i>	<i>of the 2021 Conference on Empirical Methods in</i>	804
755	<i>preprint arXiv:2502.20008</i> .	<i>Natural Language Processing</i> , pages 6761–6771.	805
756	Chuong Huynh, Jinyu Yang, Ashish Tawari, Mubarak	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	806
757	Shah, Son Tran, Raffay Hamid, Trishul Chilimbi, and	Lee. 2023a. Visual instruction tuning. <i>Advances</i>	807
758	Abhinav Shrivastava. 2025. Collm: A large language	<i>in neural information processing systems</i> , 36:34892–	808
759	model for composed image retrieval. In <i>Proceed-</i>	34916.	809
760	<i>ings of the Computer Vision and Pattern Recognition</i>		
761	<i>Conference</i> , pages 3994–4004.	Siqi Liu, Weixi Feng, Tsu-jui Fu, Wenhu Chen, and	810
762	Soyeong Jeong, Kangsan Kim, Jinheon Baek, and	William Yang Wang. 2023b. Edis: Entity-driven	811
763	Sung Ju Hwang. 2025. Videorag: Retrieval-	image search over multimodal web content. <i>arXiv</i>	812
764	augmented generation over video corpus. <i>arXiv</i>	<i>preprint arXiv:2305.13631</i> .	813
765	<i>preprint arXiv:2501.05874</i> .		
766	Ting Jiang, Minghui Song, Zihan Zhang, Haizhen	Zhenghao Liu, Chenyan Xiong, Yuanhuiyi Lv, Zhiyuan	814
767	Huang, Weiwei Deng, Feng Sun, Qi Zhang, Deqing	Liu, and Ge Yu. 2023c. Universal vision-language	815
768	Wang, and Fuzhen Zhuang. 2024. E5-v: Universal	dense retrieval: Learning a unified representation	816
769	embeddings with multimodal large language models.	space for multi-modal retrieval. In <i>The Eleventh In-</i>	817
770	<i>arXiv preprint arXiv:2407.12580</i> .	<i>ternational Conference on Learning Representations</i> .	818
771	Ziyan Jiang, Rui Meng, Xinyi Yang, Semih Yavuz,	Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney,	819
772	Yingbo Zhou, and Wenhu Chen. 2025. VLM2vec:	and Stephen Gould. 2021b. Image retrieval on	820
		real-life images with pre-trained vision-and-language	821
		models. In <i>Proceedings of the IEEE/CVF interna-</i>	822
		<i>tional conference on computer vision</i> , pages 2125–	823
		2134.	824
		Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	825
		Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	826
		try, Amanda Askell, Pamela Mishkin, Jack Clark, and	827

828	1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PmlR.		
829			
830			
831			
832	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .		
833			
834			
835			
836			
837			
838	Xiaohan Wang, Linchao Zhu, and Yi Yang. 2021. T2vlad: global-local sequence alignment for text-video retrieval. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 5079–5088.		
839			
840			
841			
842			
843	Cong Wei, Yang Chen, Haonan Chen, Hexiang Hu, Ge Zhang, Jie Fu, Alan Ritter, and Wenhu Chen. 2024. Uniir: Training and benchmarking universal multimodal information retrievers. In <i>European Conference on Computer Vision</i> , pages 387–404. Springer.		
844			
845			
846			
847			
848			
849	Hui Wu, Yupeng Gao, Xiaoxiao Guo, Ziad Al-Halah, Steven Rennie, Kristen Grauman, and Rogerio Feris. 2021. Fashion iq: A new dataset towards retrieving images by natural language feedback. In <i>Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition</i> , pages 11307–11317.		
850			
851			
852			
853			
854			
855	Qiyu Wu, Chongyang Tao, Tao Shen, Can Xu, Xiubo Geng, and Daxin Jiang. 2022. PCL: Peer-contrastive learning with diverse augmentations for unsupervised sentence embeddings. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 12052–12066, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		
856			
857			
858			
859			
860			
861			
862			
863	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, and 1 others. 2025. Qwen2. 5-omni technical report. <i>arXiv preprint arXiv:2503.20215</i> .		
864			
865			
866			
867	Mu Yang, Bowen Shi, Matthew Le, Wei-Ning Hsu, and Andros Tjandra. 2024. Audiobox tta-rag: Improving zero-shot and few-shot text-to-audio with retrieval-augmented generation. <i>arXiv preprint arXiv:2411.05141</i> .		
868			
869			
870			
871			
872	Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Richard James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-Tau Yih. 2023. Retrieval-augmented multimodal language modeling. In <i>International Conference on Machine Learning</i> , pages 39755–39769. PMLR.		
873			
874			
875			
876			
877			
878	Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions. In <i>European conference on computer vision</i> , pages 69–85. Springer.		
879			
880			
881			
882	Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. 2017. Multi-modal factorized bilinear pooling with		
883			
		co-attention learning for visual question answering. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 1821–1830.	884 885 886
		Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 11975–11986.	887 888 889 890 891
		Kai Zhang, Yi Luan, Hexiang Hu, Kenton Lee, Siyuan Qiao, Wenhu Chen, Yu Su, and Ming-Wei Chang. 2024a. Magiclens: Self-supervised image retrieval with open-ended instructions. <i>arXiv preprint arXiv:2403.19651</i> .	892 893 894 895 896
		Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi Dai, Dingkun Long, Pengjun Xie, Meishan Zhang, Wenjie Li, and Min Zhang. 2024b. Gme: Improving universal multimodal retrieval by multimodal llms. <i>arXiv preprint arXiv:2412.16855</i> .	897 898 899 900 901
		Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, WANG HongFa, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Cai Wan Zhang, Zhifeng Li, Wei Liu, and Li Yuan. 2024a. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. In <i>The Twelfth International Conference on Learning Representations</i> .	902 903 904 905 906 907 908 909
		Xinliang Zhu, Sheng-Wei Huang, Han Ding, Jinyu Yang, Kelvin Chen, Tao Zhou, Tal Neiman, Ouyue Xie, Son Tran, Benjamin Yao, and 1 others. 2024b. Bringing multimodality to amazon visual search system. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 6390–6399.	910 911 912 913 914 915 916
		Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 4995–5004.	917 918 919 920 921

## A Appendix

### A.1 Qualitative analysis

Figure 7 presents qualitative comparisons between the baseline, i.e., training with only contrastive learning, and MCA model, to highlight the behavior of modality shortcut. Across diverse cases from multiple benchmarks, the baseline fails by relying on partial cues in dominant modalities. For example, in #1 it select an irrelevant baseball scene because of the word “*catcher*”, ignoring the visual grounding on the original image. In contrast, MCA retrieves the image with textual description “*white uniform in front*”. In #2 and #3, the baseline retrieves examples relying on image similarity only, but ignoring the textual instructions. “*wearing a blue collar*” in #2 and “*standing on the ground on all fours*” in #3 are not reflected, while MCA grasps those information. In #4, the baseline mistakenly retrieves a zebra facing right, over-relying on the text cue and failing to ground the original image. MCA, by contrast, localizes the “*right zebra*” in the scene. In the news retrieval case #5, the query mentions “*David Cameron speaking in front of 10 Downing St*”, yet the baseline retrieves the pair containing “*David Cameron*” in both text and image but ignores the “*speaking in front of 10 Downing St*”. MCA composes the texts and image thus finds the best match. Similarly, in the fashion example #6, the baseline retrieves an image satisfying “*polka dot and white*” but ignoring the original clothing style in the image. In contrast, MCA integrates both modalities and selects the desired target. Those examples demonstrate how MCA mitigate modality shortcuts with the potential for more robust multimodal retrieval under various composed scenarios. In addition, we also present visualization of learned embeddings demonstrating that MCA effectively reduce the modality shortcut in §A.5.

### A.2 Datasets

**Training data.** We utilize six multimodal retrieval datasets including both *cross-modal* and *composed* tasks for training: MSCOCO (Lin et al., 2014), VisualNews (Liu et al., 2021a), VisDial (Das et al., 2017), CIRR (Liu et al., 2021b), NIGHTS (Fu et al., 2023), and WebQA (Chang et al., 2022). Table 4 shows the number for examples and input modalities of training data. We use

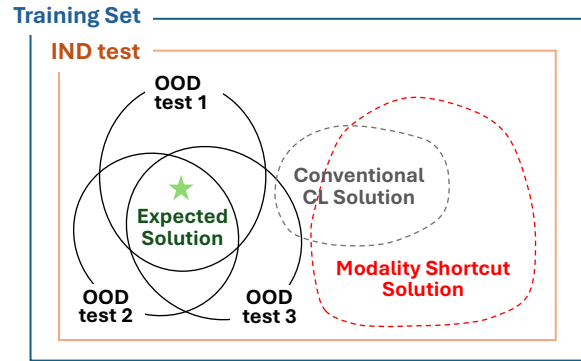


Figure 6: A conceptual diagram showing how OOD benchmarks evaluate modality shortcut.

the training splits from MMEB<sup>3</sup> without additional curation. All models are trained jointly on the union of these datasets. Although MCA is designed for composed inputs, unimodal training pairs are also included. Each training batch contains a mix of unimodal and composed inputs. This ensures that the model learns basic representation ability for each modality, which is essential modeling of structural relationships among modalities that is required by our objectives.

**Benchmarks.** We evaluate the models on a range of multimodal retrieval benchmarks, which we group into *IND* and *out-of-domain* (OOD) settings. The *IND* benchmarks correspond directly to the test splits of the six datasets used for training as mentioned in the previous paragraph and Table 4. Since models that rely on modality shortcuts often exhibit poor generalization, we also assess this issue through robustness under OOD settings. The OOD benchmarks are deliberately selected: OVEN (Hu et al., 2023) combines visual and textual modalities, where each instance consists of an image paired with a visual recognition question, alongside a reference Wikipedia image and its textual description including the title and first 100 tokens that serve as the target candidate for answering the question. FashionIQ (Hu et al., 2023) features a multimodal composition of fashion product images paired with crowd-sourced textual descriptions that specify differences between products, enabling composed image retrieval where a reference image and modification text are combined to retrieve target images. EDIS (Liu et al., 2023b) combines entity-rich text queries with multimodal candidates consisting of news images

<sup>3</sup><https://huggingface.co/datasets/TIGER-Lab/MMEB-train>



















No.	#1	#2	#3	#4	#5	#6
Query	 Select the portion of the image that follows the language expressions: "the catcher"	 Reduce the number of dogs to one dog wearing a blue collar.	 Reduce the number of monkeys to one monkey standing on the ground on all fours.	 Select the portion of the image that follows the language expressions: "right zebra"	Find a news image that matches the provided caption: Britain's Prime Minister David Cameron speaks to members of the media in front of 10 Downing St in London.	 Find an image to match the fashion image and style note: Has a polka dot pattern and is white.
Retrieved Examples	Baseline  catcher				 Assuming no sudden resignations, the Tories will be led by David Cameron during the general election campaign and Labour will be led by Gordon Brown.	
	MCA  White uniform in front				 Europe Scotland votes to remain part of United Kingdom.	 

Figure 7: Qualitative examples. The central part of the retrieved image by MCA for #6 is zoomed in for a better visibility.

Table 4: Statistics of training datasets. I and T are the abbreviation of image and text, respectively.

Dataset	MSCOCO		VisualNews		VisDial	CIRR	NIGHTS	WebQA
Input modalities	I → T	T → I	I → T	T → I	T → I	(I+T) → I	I → I	T → (I+T)
# of training pairs	113K	100K	100K	100K	123K	26K	16K	17K

paired with their headlines, requiring models to understand both textual entities/events and visual content for retrieval. MSCOCO-Grounding (Lin et al., 2014; Jiang et al., 2025) transforms object detection into a multimodal ranking task where queries combine an image with a textual object name to retrieve cropped images of the specified object, with distractors sourced from other objects within the same image and from different images. Visual7W-Pointing (Zhu et al., 2016) combines textual questions with images to establish semantic links between descriptions and image regions, enabling both visual question answering through text-image composition and visual grounding through multimodal object localization. RefCOCO (Kazemzadeh et al., 2014) employs multimodal queries combining images with referring language expressions to identify specific objects, pairing image-text queries to retrieve cropped object images. RefCOCO-Matching (Kazemzadeh et al., 2014; Jiang et al., 2025) uses the same source datasets, while repurposed to match identical image-text compositions where both query and

target contain the same object with its referring expression. We directly use the test splits from MMEB<sup>4</sup>. Note that since modality shortcut in MLLMs, the problem we study in this paper, only happens in composed queries or documents, all of our OOD benchmarks focus on *composed retrieval*, in which either query or document side comprise multiple modalities.

### A.3 Default training settings

By default, we adopt Qwen2-VL-2B-Instruct<sup>5</sup> as the unified encoder backbone. The contrastive temperature  $\tau$  is set to 0.02. All models are fine-tuned with LoRA (Hu et al., 2022) adapter. Unless otherwise specified, the LoRA rank is set to 8. The overall loss combines conventional contrastive learning and our proposed MCR and MCP terms as introduced in Equation 7, where  $\alpha$  and  $\beta$  are weighting hyper-parameters. These auxiliary terms act as regularization signals rather than main training ob-

<sup>4</sup><https://huggingface.co/datasets/TIGER-Lab/MMEB-eval>

<sup>5</sup><https://huggingface.co/Qwen/Qwen2-VL-2B-Instruct>

jectives. We empirically adopt a small coefficient between 0.01 and 1.0 to stabilize training. For simplicity, unless the otherwise specified, we tie  $\alpha$  and  $\beta$  by setting  $\alpha = \beta = 0.01$  by default so that both the preference and regularization losses contribute equally. Following previous work (Jiang et al., 2025), we use AdamW optimizer with a learning rate of  $2e - 5$  and a linear schedule with warm-up. We also fixed the global batch size as 1024 with gradient accumulation for varying numbers of GPU and memory. The total training steps are 2000 and warm-up steps are 200. The training takes around 200 hours on  $8 \times$  Nvidia A100 40G. For MCR loss, the mixer function  $\text{Mix}_\phi$  in Equation 5 is instantiated as a simple gated fusion module by default, where each unimodal embedding is reweighted by a learnable gating coefficient before aggregation. For ablation, we also experiment with alternatives such as non-parametric mean pooling and multi-modal factorized bilinear pooling (Yu et al., 2017) in §4.4.

#### A.4 Default Evaluation Settings

Unless otherwise specified, we report results using the final checkpoint of training for a fair comparison across all methods, as we have verified that the trend is stable across checkpoints in §4.1.3. We report standard retrieval metrics  $\text{accuracy}@1$  for all tasks.

#### A.5 tSNE Implementation

In Figure 5a and 5b, we visualize the learned embedding spaces using t-distributed Stochastic Neighbor Embedding (t-SNE) to analyze the distribution characteristics of different query types. Embeddings were extracted from the CIRR dataset, which contains image-text compositions for retrieval tasks. For dimensionality reduction, we employed t-SNE with perplexity=100, 1,000 iterations, and a random seed of 42 to ensure reproducibility. We randomly sampled 300 instances from each embedding type to maintain visual clarity while preserving distribution characteristics. To ensure fair comparison, we used identical sampling indices across all embedding types for baseline and our model. The visualization combines kernel density estimation (KDE) with scatter plots to represent both the overall distribution and individual data points. We implemented a custom transparency gradient for the KDE plots, making low-density regions increasingly transparent to highlight areas of embedding concentration. Differ-

ent marker shapes distinguish between embedding types: circles for composed queries, squares for text-only queries, triangles for image-only queries, and diamonds for target images. The resulting visualization effectively captures the relationships between different modalities in the embedding space, revealing patterns of overlap and separation between composed, unimodal, and target representations. Comparing Figure 5a and 5b, it clearly indicates that vanilla CL could lead to a over-mixed distribution for composed and text-only embeddings, while MCA learns a more robust representation with a clear boundary among different types of representations with a competitive IND performance and a better OOD performance.

#### A.6 LLM Usage

During the preparation of this paper, we utilized LLMs to assist with English proofreading. We also use LLM assistant for improving the efficiency of experiments such as organizing job execution scripts. Additionally, we used LLM assistant to support processing experimental results and generating figures.