

MultiModal Code-Switching: Interleaving Visual Objects into Language for Explicit Object-Level Alignment

Anonymous ACL submission

Abstract

Existing Multimodal Large Language Models (MLLMs) predominantly rely on image-text pairs for modality alignment pretraining, mapping global image representations to long textual descriptions. However, this image-level alignment suffers from *referential ambiguity*: the model must struggle to infer the correspondences between multiple visual objects and textual entities from the global representation, leading to data inefficiency and suboptimal semantic grounding. To address this, we propose MultiModal Code-Switching (MMCS), a novel pretraining paradigm that enables explicit object-level supervision. Inspired by linguistic code-switching, MMCS interleaves vision and language by replacing textual entity embeddings with embeddings of their corresponding visual objects, enforcing local visual-textual grounding during pretraining. We further develop a scalable data synthesis pipeline to generate 773k samples with accurate object-entity correspondences. Experiments across model scales show that MMCS is highly data-efficient: with only 50k samples, it matches or surpasses models trained on 600k image-caption pairs, while consistently improving visual grounding and perception capabilities¹.

1 Introduction

Multimodal Large Language Models (MLLMs) have established a new state-of-the-art in vision-language understanding, demonstrating exceptional performance across visual question answering (Goyal et al., 2017; Hudson and Manning, 2019), document understanding (Mathew et al., 2021), and visual grounding (Kazemzadeh et al., 2014; Mao et al., 2016). The dominant architecture for these models (Alayrac et al., 2022; Liu et al., 2023; Dai et al., 2023; Bai et al., 2025b; Zhu et al., 2025) generally comprises a vision encoder,

a Large Language Model (LLM) backbone, and a projector (e.g., MLP or Q-Former) that bridges the modality gap. To unify these components, the standard training paradigm follows a two-stage process: 1) modality alignment pretraining, mapping visual features into the LLM’s semantic space via image-text pairs; and 2) multimodal supervised instruction tuning (SFT), optimizing the model for downstream task execution. Within this framework, modality alignment is foundational, as the fidelity of this alignment dictates the upper bound of the model’s multimodal capabilities (Liu et al., 2023; McKinzie et al., 2024).

The prevailing consensus in recent research is to utilize dense image captions for modality alignment, providing detail-rich linguistic signals (Chen et al., 2024a,b; Li et al., 2024, 2025; Deitke et al., 2025; Xing et al., 2025). In this paradigm, the model encodes the image into a global image representation, and then is forced to predict a lengthy text sequence based on this global representation, thereby achieving alignment at the *image level*. However, natural images are inherently complex, often containing multiple objects and background elements. As highlighted in Figure 1 (right), standard dense caption datasets contain an average of 6.8 to 11.0 distinct entities per sample (Chen et al., 2024b; Onoe et al., 2024; Li et al., 2024; Garg et al., 2024). While the captions meticulously enumerate specific objects and attributes, the vision encoder and projector compress the entire scene into a generic, global representation.

This discrepancy leads to a *referential ambiguity* issue: the model must implicitly infer the correspondence between specific visual regions and the corresponding textual phrases from global representations. From a computational perspective, this “many-to-many” mapping forces the model to rely on statistical co-occurrences rather than genuine semantic grounding. Consequently, this implicit alignment paradigm is highly data-inefficient, ne-

¹We provide an anonymous code repository at https://anonymous.4open.science/r/MMCS_anonymous-4C7A/

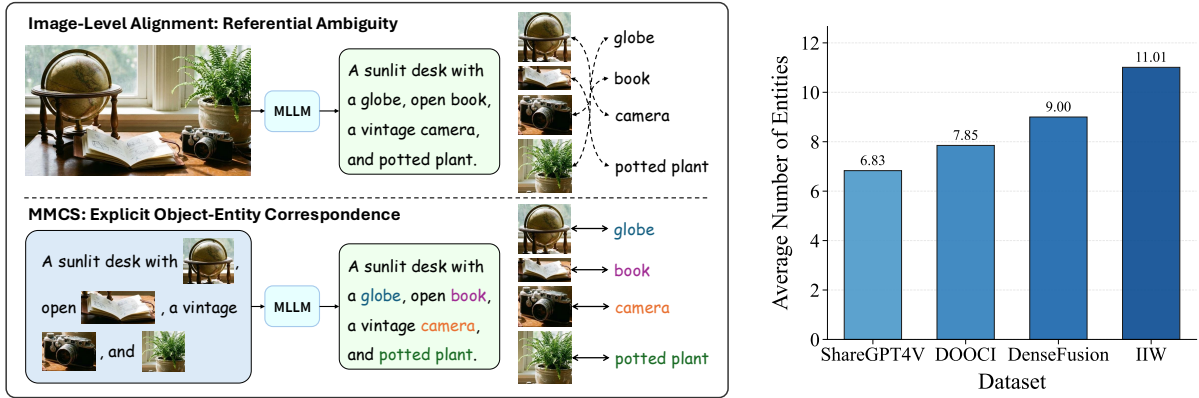


Figure 1: **Left:** Illustration of the referential ambiguity in standard image-level alignment (top), contrasted with our MultiModal Code-Switching (MMCS) paradigm (bottom). MMCS resolves ambiguity by replacing textual entities with their corresponding visual objects to provide explicit correspondence signals. **Right:** The large average number of entities in dense caption datasets highlights the complexity of natural images.

082 cessitating massive-scale datasets to learn robust
 083 object-entity associations (McKinzie et al., 2024;
 084 Dong et al., 2025). Our analysis in Section 5 sub-
 085 stantiates this by revealing that standard models
 086 exhibit diffuse attention patterns and suboptimal
 087 representation consistency when multiple localized
 088 objects are involved.

089 To address these limitations, we propose Multi-
 090 modal Code-Switching (MMCS), a novel pre-
 091 training paradigm that introduces *explicit object-*
 092 *level supervision*. Inspired by the linguistic phe-
 093 nomenon of code-switching (Poplack, 1981; Thara
 094 and Poornachandran, 2018), MMCS treats vision
 095 and language as distinct “codes”. Instead of re-
 096 lying solely on global image context, we create
 097 interleaved representations by substituting the em-
 098 beddings of textual entities with the embeddings of
 099 their corresponding visual objects (Figure 1, left).
 100 By conditioning the generation of the immediate
 101 textual context directly on these local visual fea-
 102 tures, MMCS imposes a structural constraint that
 103 enforces explicit grounding of textual entities to
 104 corresponding visual regions. This eliminates the
 105 need for the model to infer correspondences from
 106 global representations, thereby facilitating object-
 107 level alignment efficiently.

108 To implement this, we develop a data synthe-
 109 sis pipeline to generate 773k high-quality samples
 110 with accurate object-entity correspondences. Given
 111 an image, our pipeline generates a detailed caption,
 112 extracts textual entities, and employs a grounding
 113 model to localize the corresponding visual objects.
 114 Empirically, MMCS demonstrates extraordinary
 115 data efficiency compared to standard image-level
 116 pretraining. Notably, with only 50k MMCS sam-

117 ples, our model achieves performance exceeding
 118 models pretrained on 600k standard image-caption
 119 pairs—a 12-fold improvement in efficiency. More-
 120 over, MMCS maintains consistent gains when scal-
 121 ing up both the dataset size and model capacity,
 122 yielding average improvements of 7.9% on visual
 123 grounding and 2.1% on perception-centric bench-
 124 marks in data-sufficient scenarios. Further in-depth
 125 analysis confirms that these gains stem from higher-
 126 fidelity representation alignment and sharper atten-
 127 tion distributions.

128 Our contributions are summarized as follows:

- 129 • We introduce MultiModal Code-Switching
 130 (MMCS), a pretraining paradigm that shifts
 131 alignment from implicit image-level associa-
 132 tions to **explicit object-level supervision**.
- 133 • We develop a **scalable data synthesis**
 134 **pipeline** that generates 773k samples with pre-
 135 cise object-entity correspondences, bypassing
 136 the need for manual annotation.
- 137 • We conduct extensive experiments across vari-
 138 ous model scales and vision encoders, demon-
 139 strating the effectiveness of MMCS. We fur-
 140 ther provide **mechanistic insights** into how
 141 MMCS enhances the internal feature space
 142 topology of MLLMs.

143 2 Related Works

144 **Multimodal Large Language Models** Current
 145 mainstream MLLMs adopt the ViT-MLP-LLM
 146 paradigm (An et al., 2025; Bai et al., 2025a; Zhu
 147 et al., 2025; Guo et al., 2025). Specifically, this

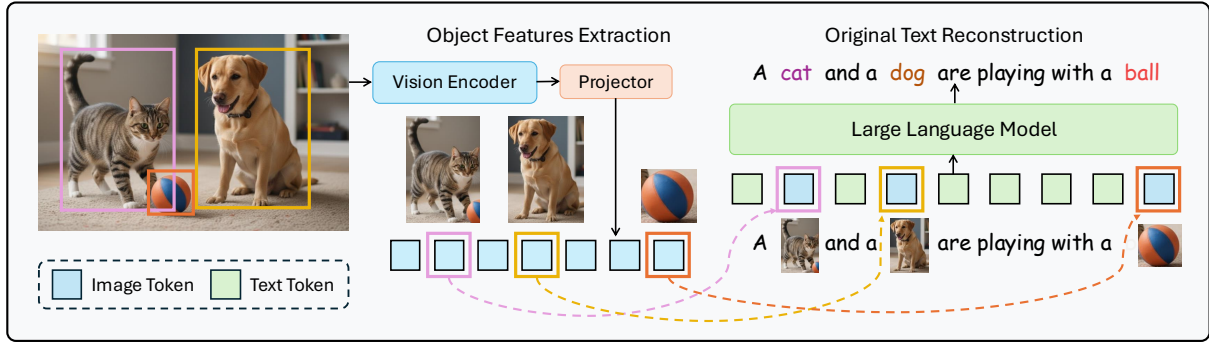


Figure 2: **Overview of MMCS pretraining paradigm.** We construct an interleaved image-text sequence by substituting textual entities with corresponding extracted visual objects. The generation of these entities and the subsequent context is conditioned on its visual counterparts, thereby facilitating object-level alignment.

architecture employs an MLP-based projector to align image features from a pretrained vision encoder with the input embedding space of an LLM backbone. These mapped visual tokens serve as soft prompts for conditional text generation. The training typically follows a two-stage strategy: where the projector is first pretrained on image-caption pairs before unfreezing the LLM for SFT (Liu et al., 2023, 2024a; Ye et al., 2023; Dong et al., 2025). The prevailing consensus in recent research is to utilize dense image captions during the modality alignment stage to provide detail-rich linguistic signals (Chen et al., 2024a,b; Li et al., 2024, 2025; Deitke et al., 2025).

Despite these advancements, models trained on holistic detailed captions suffer from a referential ambiguity issue: they must implicitly infer the correspondence between textual entities and specific visual regions from global representations. Consequently, these models struggle to disentangle specific visual concepts from complex scenes, relying instead on statistical co-occurrences. This results in data inefficiency and suboptimal semantic grounding. In contrast, our approach seeks to provide explicit object-entity mappings, enforcing the model to ground textual descriptions to their corresponding visual regions.

Fine-Grained Multimodal Alignment The limitations of image-level alignment have prompted increasing interest in fine-grained methodologies. For instance, SEA (Yin et al., 2025) utilizes CLIP (Radford et al., 2021) to assign semantic labels to visual patches and introduces a contrastive loss for patch-level alignment. Similarly, Patch Aligned Training (Jiang et al., 2025) employs off-the-shelf vision expert models (Zhang et al., 2024; Liu et al., 2024c; Kirillov et al., 2023) to generate text labels

for visual patches, directly maximizing the cosine similarity between features of these visual patches and corresponding text embeddings.

While patch-level methods demonstrate improvements over image-level alignment, they overlook a semantic mismatch: a single visual patch often lacks independent semantic completeness, whereas a word carries distinct meaning. This fragmentation results in noisy alignment signals. In contrast, our method aligns text with complete visual objects rather than fragmented patches. By focusing on objects as the fundamental unit of visual semantics, our method ensures a cleaner, more semantically coherent alignment signal.

3 Methods

3.1 Insight: Code-Switching

Our method draws inspiration from the linguistic phenomenon of Code-Switching (CS), defined as the interleaving of two or more languages within a single utterance (Thara and Poornachandran, 2018). For example, in the sentence “The *Klavier* is a versatile keyboard instrument,” the German term *Klavier* (piano) functions as a code-switched entity embedded within English context.

Conceptually, this phenomenon implies that the switched term must be semantically compatible with its surrounding context. We leverage this principle by treating visual objects as distinct “codes” carrying specific semantic information. Just as a bilingual speaker selects the most appropriate word from either language, we substitute textual entities with their visual counterparts. This imposes a constraint where the model must resolve the visual representation to satisfy the semantic expectations of the sentence, thereby shifting the learning signal from implicit global correlation to explicit object-

level grounding.

3.2 Multimodal Code-Switching Pretraining

Building on this insight, we introduce Multimodal Code-Switching (MMCS), a novel alignment paradigm that establishes correspondence between visual objects and textual entities. As illustrated in Figure 2, we construct interleaved code-switching image-text sequences by replacing the entity tokens with their corresponding visual objects. This substitution strategy creates a direct dependency chain: the generation of the entity and the subsequent context is strictly conditioned on its visual counterparts. Therefore, the model is compelled to ground entities to corresponding visual regions, eliminating the need to infer correspondences from global representations.

Formally, let $\mathbf{V} = \{v_n\}_{n=1}^N$ denote the visual tokens of the full image, where each v_n corresponds to a visual patch with spatial coordinate b_n . Given the annotated bounding box b_{object} for a specific visual object, we extract the subset of patches that spatially intersect with the object region. The object tokens $\mathbf{v}_{\text{object}}$ is defined as:

$$\mathbf{v}_{\text{object}} = \{v_n \in \mathbf{V} \mid \text{Area}(b_n \cap b_{\text{object}}) > 0\}. \quad (1)$$

Let \mathbf{X} denote the original text tokens. Consider a textual entity segment $\mathbf{e} = [x^i, \dots, x^{i+m}]$ starting at index i . The interleaved image-text sequence \mathbf{X}_{MMCS} is constructed by replacing the textual entity with the extracted object tokens:

$$\mathbf{X}_{\text{MMCS}} = \text{Concat}(\mathbf{X}^{<i}, \mathbf{v}_{\text{object}}, \mathbf{X}^{>i+m}). \quad (2)$$

We employ a language modeling objective where the loss is computed exclusively on text tokens, treating the visual tokens as conditioning context for next-token prediction:

$$\mathcal{L}_{\text{LM}} = - \sum_{x_{\text{MMCS}}^t \in \text{Text}} \log p_{\theta}(x_{\text{MMCS}}^t \mid \mathbf{X}_{\text{MMCS}}^{<t}). \quad (3)$$

To further enforce explicit supervision on object-entity correspondence, we additionally minimize the negative log-likelihood of the original textual entity segment given the preceding context:

$$\mathcal{L}_{\text{entity}} = - \log p_{\theta}(\mathbf{e} \mid \mathbf{X}_{\text{MMCS}}^{<i}, \mathbf{v}_{\text{object}}), \quad (4)$$

where \mathbf{e} represents the substituted textual entity tokens in the original text sequence. Our MMCS pretraining combines both the language modeling loss and the entity reconstruction loss:

$$\mathcal{L}_{\text{MMCS}} = \mathcal{L}_{\text{LM}} + \mathcal{L}_{\text{entity}}. \quad (5)$$

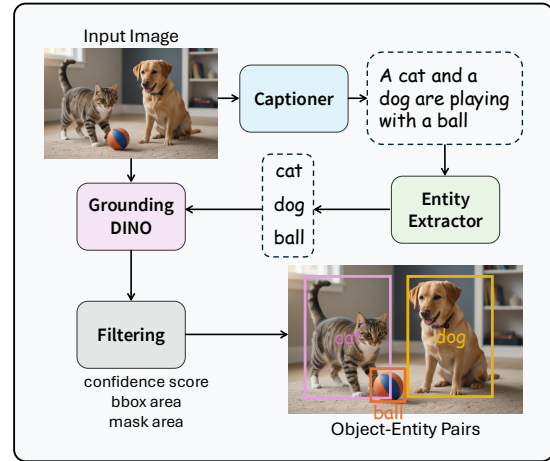


Figure 3: Data synthesis pipeline.

3.3 Data Synthesis Pipeline

As shown in Figure 3, we develop a data synthesis pipeline to generate interleaved code-switching image-text sequences with precise object-entity correspondence. For each input image, the pipeline proceeds through the following steps:

Detailed Image Captioning We begin by utilizing Qwen3-VL-32B-Instruct (Bai et al., 2025a) to generate detailed image captions. To capture granular visual details, we prompt the caption model to exhaustively describe the attributes of all visible elements within the scene.

Textual Entity Extraction Given the generated detailed captions, we employ Qwen2.5-72B-Instruct (Yang et al., 2024) to extract textual entities. The output comprises a list of noun phrases, often accompanied by brief attribute descriptions (e.g., a white mug labeled “O.CO”). These attributes are essential for disambiguating instances where multiple objects of the same category are present in a single image.

Visual Object Grounding We leverage Grounding DINO (Liu et al., 2024c) to anchor these textual entities to their corresponding visual regions. For each successfully localized entity, the output is a tuple containing the bounding box coordinates, the text label, and an associated confidence score.

Filtering To ensure high-quality object-entity correspondence, we implement a multi-stage filtering process. First, bounding boxes with low confidence scores are discarded. Subsequently, we calculate the area of each remaining bounding box and prune those violating spatial constraints (e.g.,

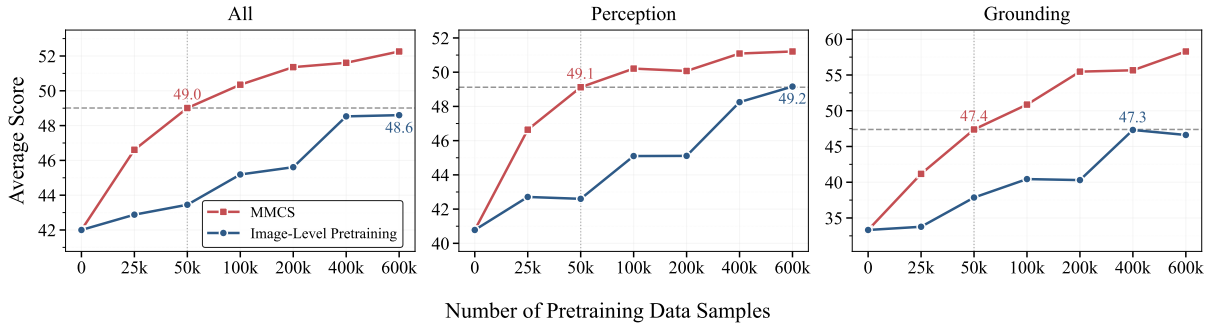


Figure 4: **Analysis of data efficiency.** We utilize Qwen2.5-3B as the LLM backbone and evaluate data efficiency by varying the pretraining dataset size from 0 to 600k, while keeping the SFT dataset fixed at 200k. The curves depict average scores across our benchmark suite, demonstrating that MMCS achieves superior data efficiency over image-level pretraining.

smaller than a single visual patch or exceeding 50% of the total image area). Finally, SAM-2.1 (Ravi et al., 2025) is employed to generate segmentation masks. Objects yielding mask areas below a predefined threshold are discarded to avoid the inclusion of heavily occluded instances.

We curate a diverse collection of images and apply the aforementioned synthesis pipeline. The generated pretraining dataset comprises 773,779 samples for pretraining, each containing a source image, its comprehensive caption, and a set of visual objects associated with bounding boxes and textual entity labels. A statistical summary of the synthesized dataset is presented in Table 6.

4 Experiments

4.1 Experiment Setup

Model Architecture 1) **Vision Encoders:** Unless otherwise specified, we utilize SigLIP2-SO-400M (Tschannen et al., 2025) equipped with a tile-wise dynamic high-resolution strategy (Liu et al., 2024b). 2) **LLMs:** MMCS is integrated with Qwen2.5-3B-Instruct (Yang et al., 2024), Qwen3-8B (Yang et al., 2025), and Llama3-8B-Instruct (Team, 2024). 3) **Projector:** We utilize a two-layer MLP with a 2×2 pixel unshuffle operation (Chen et al., 2024d) to reduce the number of image tokens.

Training Dataset For pretraining, we utilize the 773k dataset synthesized via the pipeline described in Section 3.3 for both MMCS and image-level pretraining. For SFT, we utilize the LLaVA-NeXT (Liu et al., 2024b) dataset, which contains 779k high-quality instruction-following samples.

Evaluation Our evaluation covers a comprehensive suite of benchmarks categorized into three do-

main: 1) **Visual Grounding.** We evaluate region understanding and grounding capabilities on referring expression comprehension (REC) tasks. 2) **Visual Perception,** including diagram/chart understanding, OCR-related tasks and real-world fine-grained perception. 3) **General VQA,** covering broad visual understanding and reasoning tasks.

For more details on training and evaluation, please refer to Appendix A.

4.2 Analysis of Data Efficiency

We first investigate the data efficiency of MMCS by varying the scale of the pretraining dataset. Specifically, we compare MMCS against standard image-level supervision across dataset sizes ranging from 0 to 600k samples. To isolate the impact of pretraining, we utilize a constant subset of 200k examples from LLaVA-NeXT dataset for SFT. As illustrated in Figure 4, MMCS exhibits superior data efficiency compared to pretraining with image-caption pairs. Remarkably, with only 50k samples, MMCS achieves downstream performance that surpasses image-level pretraining on 600k samples. We attribute this advantage to the explicit object-entity supervision introduced in our method, which enables the model to efficiently align visual objects with their corresponding textual entities instead of relying on massive-scale statistical co-occurrences.

4.3 Scaling Up Training Data and Model Sizes

To further investigate the scalability and universality of MMCS in data-abundant scenarios, we scale up the training to utilize the full 773k pretraining and 779k instruction-following datasets. To verify robustness across different model sizes and architectures, we extend our evaluation to larger LLMs, including Qwen3-8B and Llama3-8B.

LLM	Method	RefCOCO			RefCOCO+			RefCOCOg		AVG
		testA	testB	val	testA	testB	val	test	val	
Qwen2.5 3B	Caption	65.00	54.64	61.68	63.88	47.49	56.62	61.27	62.39	59.12
	MMCS	80.66	68.75	75.18	74.35	56.23	65.49	70.32	71.31	70.29
Qwen3 8B	Caption	85.40	74.15	81.51	80.89	64.98	72.82	76.61	78.76	76.89
	MMCS	85.02	77.80	84.21	83.93	68.71	76.98	80.11	81.58	79.79
Llama3 8B	Caption	73.29	56.23	65.94	70.10	50.62	60.59	63.86	64.65	63.16
	MMCS	83.08	70.99	78.79	77.35	57.84	68.51	72.38	74.30	72.91

Table 1: **Performance on referring expression comprehension.** “Caption” denotes image-level pretraining with image-caption pairs. We report Acc@0.5 across all splits.

LLM	ViT	Method	AI2D	ChartQA	CVBench	OCR	RWQA	VQA ^T	V*	Avg
Qwen2.5 3B	SigLIP2	Caption	66.13	56.52	61.26	44.70	55.29	58.29	48.17	55.20
		MMCS	68.72	57.36	66.87	47.60	56.99	59.35	50.79	57.46
Qwen2.5 3B	Qwen2.5ViT	Caption	69.40	70.36	56.33	58.80	55.16	62.87	57.59	62.75
		MMCS	71.15	71.92	59.78	62.40	56.86	64.97	57.59	64.67
Qwen3 8B	SigLIP2	Caption	72.51	62.48	69.71	50.60	58.17	64.37	50.79	60.94
		MMCS	73.93	64.96	72.82	53.30	60.13	65.23	53.40	62.90
Llama3 8B	SigLIP2	Caption	72.31	58.28	65.66	47.00	59.87	62.59	51.31	58.98
		MMCS	71.96	60.20	69.48	48.90	59.08	64.50	53.93	60.69

Table 2: **Performance on perception-centric benchmarks.** We report results across multiple benchmarks covering a diverse range of perception-centric tasks. “OCR”, “RWQA”, “VQA^T”, and “V*” denote OCRBench, RealWorldQA, TextVQA, and V-Star, respectively.

Visual Grounding We evaluate visual grounding capabilities on referring expression comprehension tasks. As shown in Table 1, MMCS exhibits considerable improvements over the image captioning baseline across RefCOCO/+g datasets, achieving an average gain of 7.9%. Notably, since MMCS does not introduce bounding box coordinates of visual objects during pretraining, these improvements stem directly from enhanced object recognition and region understanding capability, highlighting the effectiveness of object-level alignment.

Visual Perception As shown in Table 2, MMCS consistently outperforms standard image captioning paradigm in perception tasks, maintaining robust gains as data volumes and model sizes increase. Our method demonstrates considerable performance gains in fine-grained visual perception tasks, achieving average improvements of 4.0% on CVBench, 2.8% on OCRBench and 2.0% on V-Star. These results suggest that by grounding textual descriptions to specific visual regions, MMCS effectively mitigates the referential ambiguity inherent in global representations, thereby facilitating precise visual parsing that generalizes to diverse scenarios.

General VQA Table 3 shows that MMCS yields improvements on general VQA tasks as well. We attribute these gains to a dual mechanism: first, enhanced visual perception and region understanding provide a more accurate visual context; second, the high-quality modality alignment establishes a stronger foundation for obtaining multimodal reasoning capability in the SFT stage.

4.4 Compatibility with Different Vision Encoders

Beyond tile-wise approaches, native dynamic resolution (Wang et al., 2024; Bai et al., 2025b) constitutes another prominent strategy for processing high-resolution images. This strategy encodes images into a variable number of visual tokens while preserving original aspect ratios (Dehghani et al., 2023). We demonstrate our method’s general effectiveness using Qwen2.5ViT, the vision encoder utilized in Qwen2.5-VL (Bai et al., 2025b) which employs this native dynamic resolution mechanism. As evidenced in Table 2 and 3, MMCS maintains its superior performance in this setting. These results validate our method’s compatibility across diverse vision encoders.

LLM	ViT	Method	MMB	MME	MMStar	MMMU	MMVet	GQA	Avg
Qwen2.5 3B	SigLIP2	Caption	67.53	60.36	40.60	42.77	35.64	60.22	51.19
		MMCS	68.56	62.75	40.90	43.98	33.30	61.57	51.84
Qwen2.5 3B	Qwen2.5ViT	Caption	67.70	62.82	39.20	45.21	39.27	60.12	52.39
		MMCS	68.47	65.68	42.90	45.31	39.31	60.66	53.72
Qwen3 8B	SigLIP2	Caption	74.48	68.14	47.10	48.67	44.72	62.54	57.61
		MMCS	76.63	68.54	48.80	51.49	41.56	63.02	58.34
Llama3 8B	SigLIP2	Caption	70.79	63.43	39.20	44.96	39.36	63.17	53.48
		MMCS	68.99	65.54	40.90	43.15	42.06	63.40	54.01

Table 3: **Performance on general VQA benchmarks.** We report the normalized score for MME.

Method	General	Perception	Grounding	All
Caption	51.19	55.77	59.12	55.36
Patch Aligned	51.78	56.43	66.55	58.26
MMCS	51.84	58.24	70.29	60.12
w/o $\mathcal{L}_{\text{entity}}$	51.45	54.23	67.55	57.74
w/o \mathcal{L}_{LM}	51.37	56.26	65.44	57.69

Table 4: **Performance comparison with patch-level alignment and ablation study.** All experiments utilize Qwen2.5-3B-Instruct as the LLM backbone.

4.5 Comparison with Patch-Level Alignment

We provide an empirical comparison between MMCS and a method for patch-level alignment, Patch Aligned (Jiang et al., 2025). To ensure a fair comparison, we reproduce Patch Aligned using the identical pretraining and SFT datasets as our method. As presented in Table 4, MMCS consistently outperforms the patch-level baseline across all downstream tasks, achieving an average improvement of 1.9%. These results confirm that mapping textual entities to distinct visual objects rather than arbitrary patches establishes a more robust and semantically coherent alignment.

4.6 Ablation Study

We perform an ablation study on the components of Eq. 5. The results in Table 4 demonstrate that both objectives are crucial for optimal performance.

Role of $\mathcal{L}_{\text{entity}}$: Enforcing Semantic Precision.

Ablating the entity reconstruction loss ($\mathcal{L}_{\text{entity}}$) results in the most severe drop in visual perception performance (-4.0%). This indicates that $\mathcal{L}_{\text{entity}}$ acts as a *semantic anchor*. By forcing the model to explicitly translate visual features back into textual entities, this objective compels the projector to encode precise and discriminative visual details, thereby enhancing visual perception.

Role of \mathcal{L}_{LM} : Facilitating Contextual Integration.

Conversely, removing the language modeling loss (\mathcal{L}_{LM}) primarily degrades visual grounding capabilities (-4.9%). We attribute this to \mathcal{L}_{LM} 's role in *contextual integration*. It encourages the model to reason about the relationship between the visual object and its descriptive attributes, which is critical for both visual perception and referring expression comprehension.

5 Further Discussion

In this section, we investigate the mechanisms underlying the observed improvements by analyzing the MLLM after modality alignment pretraining but prior to the SFT stage.

5.1 Representation Alignment Measurement

To investigate whether MMCS achieves superior multimodal alignment at the feature level compared to standard image-level pretraining, we directly quantify vision-language representational alignment within the MLLM. We employ three metrics: CKA (Kornblith et al., 2019) to evaluate *global* geometric correspondence, and Mutual k-NN and CKNNA (Huh et al., 2024) to assess *local* neighborhood consistency. Detailed definitions of these metrics are provided in Appendix B.

Specifically, we input 1,000 image-caption pairs sampled from the COCO2014 validation set into the MLLM and extract hidden states from all layers. These states are subsequently partitioned into visual and textual components to compute layer-wise alignment metrics. As illustrated in Figure 5, MMCS achieves superior representational alignment than image-level supervision. These findings align with the hypothesis that representational alignment correlates with model capability (Huh et al., 2024), offering a rationale for the performance improvements observed in downstream tasks.

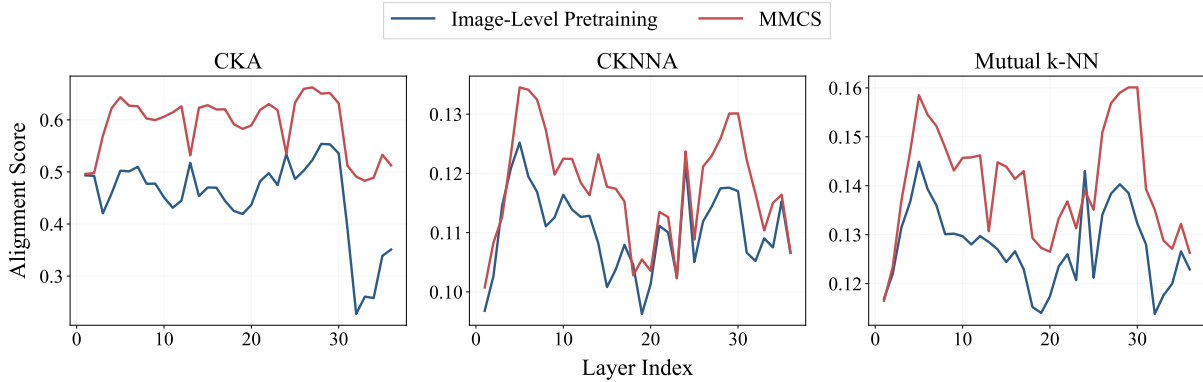


Figure 5: **Layer-wise representation alignment.** To evaluate cross-modal integration, we extract hidden states from each LLM decoder layer for input image-text pairs and partition them by modality. We then quantify the alignment using three distinct metrics: CKA, CKNNA and Mutual k-NN. MMCS exhibits superior representation alignment across the majority of layers over the standard image-level pretraining.

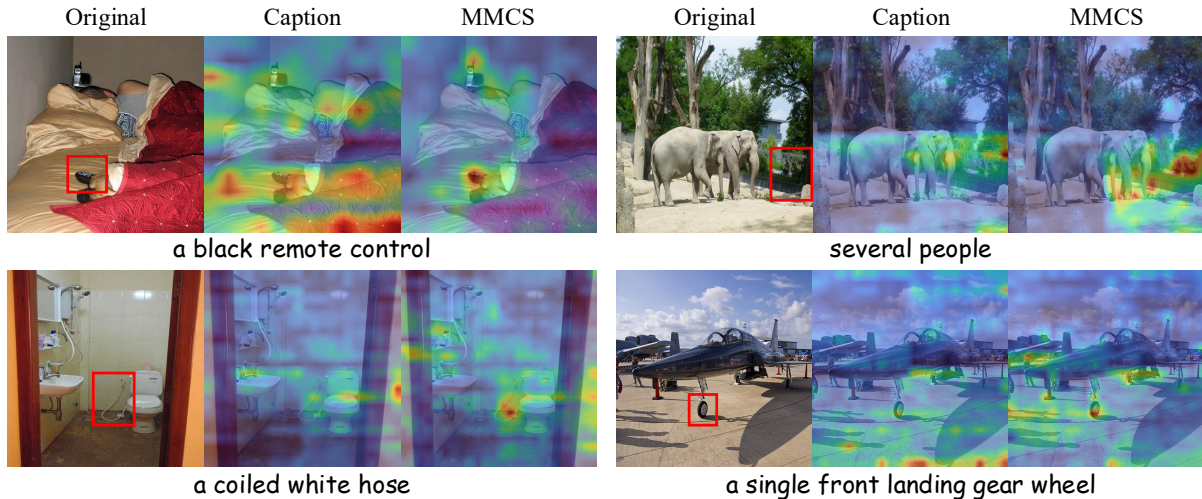


Figure 6: **Visualization of attention maps.** Each example displays a triplet containing: (Left) the original image with a red bounding box highlighting the target object; (Middle) the attention map from an MLLM pretrained with standard captioning; and (Right) the attention map from our MMCS method. Warmer colors indicate higher attention weights (min-max normalized). The specific text entity driving the attention is noted below each example.

5.2 Attention Map Analysis

Figure 6 presents a qualitative comparison of attention maps. We analyze the cross-modal attention distribution from textual entity tokens to visual features by feeding image-caption pairs into the MLLM. As observed, the model pretrained with MMCS accurately attends to visual regions corresponding to specific textual entity descriptions. In contrast, the model trained with image-caption pairs often produces diffuse or misaligned attention patterns. Given that the LLM backbone is kept frozen during pretraining, these results indicate that training the projector with explicit object-entity correspondence induces the learning of more semantically precise and interpretable features. These features align better with the LLM’s pre-existing

semantic space.

6 Conclusion

In this paper, we propose MMCS, a modality alignment pretraining paradigm that establishes explicit object-entity correspondence. Our extensive experiments reveal that MMCS demonstrates considerable improvements in both data efficiency and downstream performance compared to image-level pretraining. Furthermore, we investigate the advantages of MMCS in facilitating multimodal representation alignment from the perspective of the internal feature space topology. Our findings highlight the importance of object-level alignment in developing data-efficient MLLMs with advanced performance.

512 Limitations

513 While MMCS demonstrates considerable improve-
514 ments on data efficiency and downstream perfor-
515 mance, our current implementation primarily fo-
516 cuses on visual objects within natural images. The
517 core methodology of establishing explicit corre-
518 spondence between local visual features and tex-
519 tual entities is generalizable to broader scenarios.
520 Extending MMCS to domains such as chart under-
521 standing or scene text recognition will be further
522 explored in our future works.

523 We also acknowledge that the quality of the pre-
524 training data is currently bounded by the perfor-
525 mance of the vision expert tools employed. Devel-
526 oping a more scalable, high-fidelity data synthesis
527 pipeline to mitigate noise from these upstream mod-
528 els remains a critical direction for future research.

529 References

530 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
531 Antoine Miech, Iain Barr, Yana Hasson, Karel
532 Lenc, Arthur Mensch, Katherine Millican, Malcolm
533 Reynolds, Roman Ring, Eliza Rutherford, Serkan
534 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei,
535 Marianne Monteiro, Jacob L. Menick, Sebastian
536 Borgeaud, and 8 others. 2022. [Flamingo: a visual
537 language model for few-shot learning](#). In *Advances
538 in Neural Information Processing Systems 35: Annual
539 Conference on Neural Information Processing
540 Systems 2022, NeurIPS 2022, New Orleans, LA, USA,
541 November 28 - December 9, 2022*.

542 Xiang An, Yin Xie, Kaicheng Yang, Wenkang Zhang,
543 Xiuwei Zhao, Zheng Cheng, Yirui Wang, Songcen
544 Xu, Changrui Chen, Chunsheng Wu, Huajie Tan,
545 Chunyuan Li, Jing Yang, Jie Yu, Xiyao Wang, Bin
546 Qin, Yumeng Wang, Zizhen Yan, Ziyong Feng, and 3
547 others. 2025. [Llava-onevision-1.5: Fully open frame-
548 work for democratized multimodal training](#). *CoRR*,
549 abs/2509.23661.

550 Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen,
551 Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei
552 Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi-
553 fang Guo, Qidong Huang, Jie Huang, Fei Huang,
554 Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng
555 Li, and 45 others. 2025a. [Qwen3-vl technical report](#).
556 *Preprint*, arXiv:2511.21631.

557 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
558 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
559 Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Ming-
560 Hsuan Yang, Zhaohai Li, Jianqiang Wan, Pengfei
561 Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 oth-
562 ers. 2025b. [Qwen2.5-vl technical report](#). *CoRR*,
563 abs/2502.13923.

564 Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Jun-
565 ying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen,

Jianquan Li, Xiang Wan, and Benyou Wang. 2024a. [Allava: Harnessing gpt4v-synthesized data for A lite vision-language model](#). *CoRR*, abs/2402.11684. 566
567
568

569 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Cong-
570 hui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2024b. [Sharegpt4v: Improving large multi-modal
571 models with better captions](#). In *Computer Vision
572 - ECCV 2024 - 18th European Conference, Milan,
573 Italy, September 29-October 4, 2024, Proceedings,
574 Part XVII*, volume 15075 of *Lecture Notes in Com-
575 puter Science*, pages 370–387. Springer. 576

577 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang
578 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang,
579 Yu Qiao, Dahua Lin, and Feng Zhao. 2024c. [Are we
580 on the right way for evaluating large vision-language
581 models?](#) In *Advances in Neural Information Pro-
582 cessing Systems 38: Annual Conference on Neural
583 Information Processing Systems 2024, NeurIPS 2024,
584 Vancouver, BC, Canada, December 10 - 15, 2024*.

585 Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu,
586 Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye,
587 Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang,
588 Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo,
589 Jiahao Wang, Tan Jiang, Bo Wang, and 21 others. 2024d. [Expanding performance boundaries of open-
590 source multimodal models with model, data, and
591 test-time scaling](#). *CoRR*, abs/2412.05271. 592

593 Wenliang Dai, Junnan Li, Dongxu Li, Anthony
594 Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
595 Boyang Li, Pascale Fung, and Steven C. H. Hoi.
2023. [Instructblip: Towards general-purpose vision-
596 language models with instruction tuning](#). In *Ad-
597 vances in Neural Information Processing Systems
598 36: Annual Conference on Neural Information Pro-
599 cessing Systems 2023, NeurIPS 2023, New Orleans,
600 LA, USA, December 10 - 16, 2023*.

602 Mostafa Dehghani, Basil Mustafa, Josip Djolonga,
603 Jonathan Heek, Matthias Minderer, Mathilde Caron,
604 Andreas Steiner, Joan Puigcerver, Robert Geirhos,
605 Ibrahim M. Alabdulmohsin, Avital Oliver, Piotr
606 Padlewski, Alexey A. Gritsenko, Mario Lucic, and
607 Neil Houlsby. 2023. [Patch n’ pack: Navit, a vision
608 transformer for any aspect ratio and resolution](#). In
609 *Advances in Neural Information Processing Systems
610 36: Annual Conference on Neural Information Pro-
611 cessing Systems 2023, NeurIPS 2023, New Orleans,
612 LA, USA, December 10 - 16, 2023*.

613 Matt Deitke, Christopher Clark, Sangho Lee, Rohun
614 Tripathi, Yue Yang, Jae Sung Park, Mohammadreza
615 Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini,
616 Jiasen Lu, Taira Anderson, Erin Bransom, Kiana
617 Ehsani, Huong Ngo, Yen-Sung Chen, Ajay Patel,
618 Mark Yatskar, Chris Callison-Burch, and 31 others.
2025. [Molmo and pixmo: Open weights and open
619 data for state-of-the-art vision-language models](#). In
620 *IEEE/CVF Conference on Computer Vision and Pat-
621 tern Recognition, CVPR 2025, Nashville, TN, USA,
622 June 11-15, 2025*, pages 91–104. Computer Vision
623 Foundation / IEEE. 624

625	Hongyuan Dong, Zijian Kang, Weijie Yin, LiangXiao	and enhancing vision understanding in multimodal	682
626	LiangXiao, ChaoFeng ChaoFeng, and Ran Jiao. 2025.	language models. <i>CoRR</i> , abs/2505.17316.	683
627	Scalable vision language model training via high		
628	quality data curation. In <i>Proceedings of the 63rd An-</i>		
629	<i>annual Meeting of the Association for Computational</i>	Kushal Kafle, Brian L. Price, Scott Cohen, and Christo-	684
630	<i>Linguistics (Volume 1: Long Papers), ACL 2025, Vi-</i>	pher Kanan. 2018. <i>DVQA: understanding data vi-</i>	685
631	<i>enna, Austria, July 27 - August 1, 2025</i> , pages 33272–	<i>ualizations via question answering</i> . In <i>2018 IEEE</i>	686
632	33293. Association for Computational Linguistics.	<i>Conference on Computer Vision and Pattern Recog-</i>	687
		<i>nition, CVPR 2018, Salt Lake City, UT, USA, June</i>	688
		<i>18-22, 2018</i> , pages 5648–5656. Computer Vision	689
633	Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin,	Foundation / IEEE Computer Society.	690
634	Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jin-		
635	rui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Ron-	Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and	691
636	grong Ji. 2023. <i>MME: A comprehensive evaluation</i>	Tamara L. Berg. 2014. <i>Referitgame: Referring to</i>	692
637	<i>benchmark for multimodal large language models.</i>	<i>objects in photographs of natural scenes</i> . In <i>Proceed-</i>	693
638	<i>CoRR</i> , abs/2306.13394.	<i>ings of the 2014 Conference on Empirical Methods in</i>	694
		<i>Natural Language Processing, EMNLP 2014, Octo-</i>	695
639	Roopal Garg, Andrea Burns, Burcu Karagol Ayan,	<i>ber 25-29, 2014, Doha, Qatar, A meeting of SIGDAT,</i>	696
640	Yonatan Bitton, Ceslee Montgomery, Yasumasa	<i>a Special Interest Group of the ACL</i> , pages 787–798.	697
641	Onoe, Andrew Bunner, Ranjay Krishna, Jason	ACL.	698
642	Baldrige, and Radu Soricut. 2024. <i>Imageinwords:</i>	Aniruddha Kembhavi, Mike Salvato, Eric Kolve,	699
643	<i>Unlocking hyper-detailed image descriptions</i> . In <i>Pro-</i>	Min Joon Seo, Hannaneh Hajishirzi, and Ali Farhadi.	700
644	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	2016. <i>A diagram is worth a dozen images</i> . In <i>Com-</i>	701
645	<i>ods in Natural Language Processing, EMNLP 2024,</i>	<i>puter Vision - ECCV 2016 - 14th European Confer-</i>	702
646	<i>Miami, FL, USA, November 12-16, 2024</i> , pages 93–	<i>ence, Amsterdam, The Netherlands, October 11-14,</i>	703
647	127. Association for Computational Linguistics.	<i>2016, Proceedings, Part IV</i> , volume 9908 of <i>Lecture</i>	704
		<i>Notes in Computer Science</i> , pages 235–251. Springer.	705
648	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv	Geewook Kim, Teakgyu Hong, Moonbin Yim,	706
649	Batra, and Devi Parikh. 2017. <i>Making the V in VQA</i>	JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Won-	707
650	<i>matter: Elevating the role of image understanding in</i>	seok Hwang, Sangdoon Yun, Dongyoon Han, and Se-	708
651	<i>visual question answering</i> . In <i>2017 IEEE Conference</i>	unghyun Park. 2022. <i>Ocr-free document understand-</i>	709
652	<i>on Computer Vision and Pattern Recognition, CVPR</i>	<i>ing transformer</i> . In <i>Computer Vision - ECCV 2022</i>	710
653	<i>2017, Honolulu, HI, USA, July 21-26, 2017</i> , pages	<i>- 17th European Conference, Tel Aviv, Israel, Octo-</i>	711
654	6325–6334. IEEE Computer Society.	<i>ber 23-27, 2022, Proceedings, Part XXVIII</i> , volume	712
		13688 of <i>Lecture Notes in Computer Science</i> , pages	713
655	Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng,	498–517. Springer.	714
656	Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang,	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi	715
657	Jianguo Jiang, Jiawei Wang, Jingji Chen, Jingjia	Mao, Chloé Rolland, Laura Gustafson, Tete Xiao,	716
658	Huang, Kang Lei, Liping Yuan, Lishu Luo, Pengfei	Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo,	717
659	Liu, Qinghao Ye, Rui Qian, Shen Yan, and 81 oth-	Piotr Dollár, and Ross B. Girshick. 2023. <i>Segment</i>	718
660	ers. 2025. <i>Seed1.5-vl technical report</i> . <i>CoRR</i> ,	<i>anything</i> . In <i>IEEE/CVF International Conference on</i>	719
661	abs/2505.07062.	<i>Computer Vision, ICCV 2023, Paris, France, October</i>	720
		<i>1-6, 2023</i> , pages 3992–4003. IEEE.	721
662	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Simon Kornblith, Mohammad Norouzi, Honglak Lee,	722
663	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	and Geoffrey E. Hinton. 2019. <i>Similarity of neural</i>	723
664	Weizhu Chen. 2022. <i>Lora: Low-rank adaptation of</i>	<i>network representations revisited</i> . In <i>Proceedings of</i>	724
665	<i>large language models</i> . In <i>The Tenth International</i>	<i>the 36th International Conference on Machine Learn-</i>	725
666	<i>Conference on Learning Representations, ICLR 2022,</i>	<i>ing, ICML 2019, 9-15 June 2019, Long Beach, Cali-</i>	726
667	<i>Virtual Event, April 25-29, 2022</i> . OpenReview.net.	<i>fornia, USA</i> , volume 97 of <i>Proceedings of Machine</i>	727
		<i>Learning Research</i> , pages 3519–3529. PMLR.	728
668	Drew A. Hudson and Christopher D. Manning. 2019.	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin John-	729
669	<i>GQA: A new dataset for real-world visual reason-</i>	son, Kenji Hata, Joshua Kravitz, Stephanie Chen,	730
670	<i>ing and compositional question answering</i> . In <i>IEEE</i>	Yannis Kalantidis, Li-Jia Li, David A. Shamma,	731
671	<i>Conference on Computer Vision and Pattern Recogni-</i>	Michael S. Bernstein, and Li Fei-Fei. 2017. <i>Vi-</i>	732
672	<i>tion, CVPR 2019, Long Beach, CA, USA, June 16-20,</i>	<i>sual genome: Connecting language and vision us-</i>	733
673	<i>2019</i> , pages 6700–6709. Computer Vision Founda-	<i>ing crowdsourced dense image annotations</i> . <i>Int. J.</i>	734
674	<i>tion / IEEE</i> .	<i>Comput. Vis.</i> , 123(1):32–73.	735
675	Minyoung Huh, Brian Cheung, Tongzhou Wang, and	Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper	736
676	Phillip Isola. 2024. <i>Position: The platonic represen-</i>	R. R. Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab	737
677	<i>tation hypothesis</i> . In <i>Forty-first International Con-</i>	Kamali, Stefan Popov, Matteo Mallocci, Alexander	738
678	<i>ference on Machine Learning, ICML 2024, Vienna,</i>		
679	<i>Austria, July 21-27, 2024</i> . OpenReview.net.		
680	Jiachen Jiang, Jinxin Zhou, Bo Peng, Xia Ning, and		
681	Zhihui Zhu. 2025. <i>Analyzing fine-grained alignment</i>		

739	Kolesnikov, Tom Duerig, and Vittorio Ferrari. 2020.	Part VI, volume 15064 of <i>Lecture Notes in Computer Science</i> , pages 216–233. Springer.	796
740	The open images dataset V4 . <i>Int. J. Comput. Vis.</i> ,		797
741	128(7):1956–1981.		
742	Xiangtai Li, Tao Zhang, Yanwei Li, Haobo Yuan,	Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang,	798
743	Shihao Chen, Yikang Zhou, Jiahao Meng, Yueyi	Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-	799
744	Sun, Shilin Xu, Lu Qi, Tianheng Cheng, Yi Lin,	Lin Liu, Lianwen Jin, and Xiang Bai. 2024e. Ocr-	800
745	Zilong Huang, Wenhao Huang, Jiashi Feng, and	bench: on the hidden mystery of OCR in large multi-	801
746	Guang Shi. 2025. Denseworld-1m: Towards detailed	modal models . <i>Sci. China Inf. Sci.</i> , 67(12).	802
747	dense grounded caption in the real world . <i>CoRR</i> ,		
748	abs/2506.24102.	Junhua Mao, Jonathan Huang, Alexander Toshev, Oana	803
749	Xiaotong Li, Fan Zhang, Haiwen Diao, Yuezhe	Camburu, Alan L. Yuille, and Kevin Murphy. 2016.	804
750	Wang, Xinlong Wang, and Lingyu Duan. 2024.	Generation and comprehension of unambiguous ob-	805
751	Densefusion-1m: Merging vision experts for com-	ject descriptions . In <i>2016 IEEE Conference on Com-</i>	806
752	prehensive multimodal perception . In <i>Advances in</i>	<i>puter Vision and Pattern Recognition, CVPR 2016,</i>	807
753	<i>Neural Information Processing Systems 38: Annual</i>	<i>Las Vegas, NV, USA, June 27-30, 2016</i> , pages 11–20.	808
754	<i>Conference on Neural Information Processing Sys-</i>	IEEE Computer Society.	809
755	<i>tems 2024, NeurIPS 2024, Vancouver, BC, Canada,</i>		
756	<i>December 10 - 15, 2024</i> .	Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R.	810
757	Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James	Joty, and Enamul Hoque. 2022. Chartqa: A bench-	811
758	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,	mark for question answering about charts with visual	812
759	and C. Lawrence Zitnick. 2014. Microsoft COCO:	and logical reasoning . In <i>Findings of the Association</i>	813
760	common objects in context . In <i>Computer Vision -</i>	<i>for Computational Linguistics: ACL 2022, Dublin,</i>	814
761	<i>ECCV 2014 - 13th European Conference, Zurich,</i>	<i>Ireland, May 22-27, 2022</i> , pages 2263–2279. Associ-	815
762	<i>Switzerland, September 6-12, 2014, Proceedings,</i>	ation for Computational Linguistics.	816
763	<i>Part V, volume 8693 of Lecture Notes in Computer</i>		
764	<i>Science</i> , pages 740–755. Springer.	Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawa-	817
765	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	har. 2021. Docvqa: A dataset for VQA on document	818
766	Lee. 2024a. Improved baselines with visual instruc-	images . In <i>IEEE Winter Conference on Applications</i>	819
767	tion tuning . In <i>IEEE/CVF Conference on Computer</i>	<i>of Computer Vision, WACV 2021, Waikoloa, HI, USA,</i>	820
768	<i>Vision and Pattern Recognition, CVPR 2024, Seat-</i>	<i>January 3-8, 2021</i> , pages 2199–2208. IEEE.	821
769	<i>tle, WA, USA, June 16-22, 2024</i> , pages 26286–26296.		
770	IEEE.	Brandon McKinzie, Zhe Gan, Jean-Philippe Faucon-	822
771	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan	nier, Sam Dodge, Bowen Zhang, Philipp Dufter,	823
772	Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llava-	Dhruti Shah, Xianzhi Du, Futang Peng, Anton Be-	824
773	next: Improved reasoning, ocr, and world knowledge .	lyi, Haotian Zhang, Karanjeet Singh, Doug Kang,	825
774	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	Hongyu Hè, Max Schwarzer, Tom Gunter, Xiang	826
775	Lee. 2023. Visual instruction tuning . In <i>Advances in</i>	Kong, Aonan Zhang, Jianyu Wang, and 10 others.	827
776	<i>Neural Information Processing Systems 36: Annual</i>	2024. MM1: methods, analysis and insights from	828
777	<i>Conference on Neural Information Processing Sys-</i>	multimodal LLM pre-training . In <i>Computer Vision</i>	829
778	<i>tems 2023, NeurIPS 2023, New Orleans, LA, USA,</i>	<i>- ECCV 2024 - 18th European Conference, Milan,</i>	830
779	<i>December 10 - 16, 2023</i> .	<i>Italy, September 29-October 4, 2024, Proceedings,</i>	831
780	Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao	<i>Part XXIX, volume 15087 of Lecture Notes in Com-</i>	832
781	Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jian-	<i>puter Science</i> , pages 304–323. Springer.	833
782	wei Yang, Hang Su, Jun Zhu, and Lei Zhang. 2024c.	Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh,	834
783	Grounding DINO: marrying DINO with grounded	and Anirban Chakraborty. 2019. OCR-VQA: visual	835
784	pre-training for open-set object detection . In <i>Com-</i>	question answering by reading text in images . In	836
785	<i>puter Vision - ECCV 2024 - 18th European Confer-</i>	<i>2019 International Conference on Document Analy-</i>	837
786	<i>ence, Milan, Italy, September 29-October 4, 2024,</i>	<i>sis and Recognition, ICDAR 2019, Sydney, Australia,</i>	838
787	<i>Proceedings, Part XLVII, volume 15105 of Lecture</i>	<i>September 20-25, 2019</i> , pages 947–952. IEEE.	839
788	<i>Notes in Computer Science</i> , pages 38–55. Springer.	Yasumasa Onoe, Sunayana Rane, Zachary Berger,	840
789	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li,	Yonatan Bitton, Jaemin Cho, Roopal Garg, Alexander	841
790	Songyang Zhang, Wangbo Zhao, Yike Yuan, Ji-	Ku, Zarana Parekh, Jordi Pont-Tuset, Garrett Tanzer,	842
791	aqi Wang, Conghui He, Ziwei Liu, Kai Chen, and	Su Wang, and Jason Baldridge. 2024. DOCCI: de-	843
792	Dahua Lin. 2024d. Mmbench: Is your multi-modal	scriptions of connected and contrasting images . In	844
793	model an all-around player? In <i>Computer Vision</i>	<i>Computer Vision - ECCV 2024 - 18th European Con-</i>	845
794	<i>- ECCV 2024 - 18th European Conference, Milan,</i>	<i>ference, Milan, Italy, September 29-October 4, 2024,</i>	846
795	<i>Italy, September 29-October 4, 2024, Proceedings,</i>	<i>Proceedings, Part LX, volume 15118 of Lecture</i>	847
		<i>Notes in Computer Science</i> , pages 291–309. Springer.	848
		Shana Poplack. 1981. Syntactic structure and social	849
		function of code-switching. In Richard P. Durán, edi-	850
		tor, <i>Latino Language and Communicative Behavior</i> ,	851
		pages 169–184. Ablex, Norwood, NJ.	852

853	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 models from natural language supervision . In <i>Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 8748–8763. PMLR.	<i>Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	912 913 914 915
854			
855			
856			
857			
858			
859			
860			
861			
862			
863	Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryal, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloé Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross B. Girshick, Piotr Dollár, and Christoph Feichtenhofer. 2025. SAM 2: Segment anything in images and videos . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	Michael Tschannen, Alexey A. Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier J. Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. 2025. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features . <i>CoRR</i> , abs/2502.14786.	916 917 918 919 920 921 922 923
864			
865			
866			
867			
868			
869			
870			
871			
872			
873	Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: an open large-scale dataset for training next generation image-text models . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution . <i>CoRR</i> , abs/2409.12191.	924 925 926 927 928 929 930
874			
875			
876			
877			
878			
879			
880			
881			
882			
883			
884			
885	Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. 2019. Objects365: A large-scale, high-quality dataset for object detection . In <i>2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019</i> , pages 8429–8438. IEEE.	Penghao Wu and Saining Xie. 2024. V*: Guided visual search as a core mechanism in multimodal llms . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 13084–13094. IEEE.	931 932 933 934 935
886			
887			
888			
889			
890			
891			
892	Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019</i> , pages 8317–8326. Computer Vision Foundation / IEEE.	xAI Team. 2024. Grok-1.5 Vision Preview .	936
893			
894			
895			
896			
897			
898			
899	Llama Team. 2024. The llama 3 herd of models . <i>CoRR</i> , abs/2407.21783.	Long Xing, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Jianze Liang, Qidong Huang, Jiaqi Wang, Feng Wu, and Dahua Lin. 2025. Caprl: Stimulating dense image caption capabilities via reinforcement learning . <i>CoRR</i> , abs/2509.22647.	937 938 939 940 941
900			
901	S. Thara and Prabaharan Poornachandran. 2018. Code-mixing: A brief survey . In <i>2018 International Conference on Advances in Computing, Communications and Informatics, ICACCI 2018, Bangalore, India, September 19-22, 2018</i> , pages 2382–2388. IEEE.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40 others. 2025. Qwen3 technical report . <i>CoRR</i> , abs/2505.09388.	942 943 944 945 946 947 948
902			
903			
904			
905			
906	Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Iyer, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, Xichen Pan, Rob Fergus, Yann LeCun, and Saining Xie. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms . In <i>Advances in Neural</i>	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jixi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. Qwen2.5 technical report . <i>CoRR</i> , abs/2412.15115.	949 950 951 952 953 954 955
907			
908			
909			
910			
911			
912			
913			
914			
915			
916			
917			
918			
919			
920			
921			
922			
923			
924			
925			
926			
927			
928			
929			
930			
931			
932			
933			
934			
935			
936			
937			
938			
939			
940			
941			
942			
943			
944			
945			
946			
947			
948			
949			
950			
951			
952			
953			
954			
955			
956			
957			
958			
959			
960			
961			
962			
963			
964			
965			
966			
967			
968			

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. [From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions](#). *Trans. Assoc. Comput. Linguistics*, 2:67–78.

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2024. [Mm-vet: Evaluating large multimodal models for integrated capabilities](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024. [MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 9556–9567. IEEE.

Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2025. [Lmms-eval: Reality check on the evaluation of large multimodal models](#). In *Findings of the Association for Computational Linguistics: NAACL 2025, Albuquerque, New Mexico, USA, April 29 - May 4, 2025*, pages 881–916. Association for Computational Linguistics.

Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, Yandong Guo, and Lei Zhang. 2024. [Recognize anything: A strong image tagging model](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024 - Workshops, Seattle, WA, USA, June 17-18, 2024*, pages 1724–1732. IEEE.

Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, and 32 others. 2025. [InternV3: Exploring advanced training and test-time recipes for open-source multimodal models](#). *CoRR*, abs/2504.10479.

A Implementation Details

A.1 Training Hyperparameters

Table 5 details the hyperparameters employed for both modality alignment pretraining and SFT stages. High-resolution images are divided into smaller image tiles of the resolution that the ViT is originally trained for, and encoded independently. All experiments are conducted on 8 NVIDIA

Hyperparameter	Pretrain	Finetune
Trainable	MLP	MLP+LLM
Global batch size	128	64
Learning rate	1e-3	2e-4
Epochs	1	
Max image tiles	4 + 1	
Optimizer	AdamW	
LR schedule	Cosine decay	
Warm up ratio	0.03	
Weight decay	0	
LoRA r	-	32
LoRA alpha	-	64
LoRA dropout	-	0.05

Table 5: **Hyperparameters for model training.** “4+1” indicates that the high-resolution image is divided into at most 4 tiles with an additional thumbnail tile.

# Samples	# Objects	Avg. Chars	Avg. Objects
773,779	5,145,630	961.65	6.65

Table 6: **Statistics of the synthesized pretraining dataset.** “Avg. Chars” denotes the average character count per detailed caption, while “Avg. Objects” represents the average number of visual objects identified per sample after filtering.

A6000 GPUs. During the SFT stage, we apply LoRA (Hu et al., 2022) to the LLM backbone. For 8B-scale models, pretraining completes within 10 hours, while SFT completes within 20 hours.

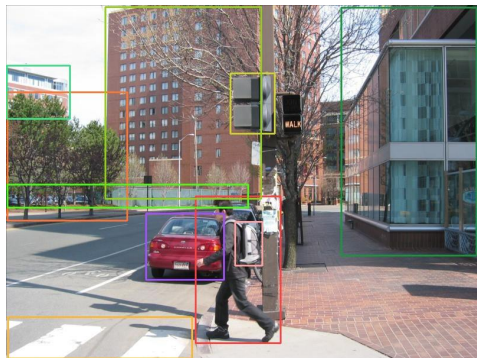
A.2 Training Datasets

The image sources utilized in Section 3.3 include MS COCO (Lin et al., 2014), Flickr30k (Young et al., 2014), GQA (Hudson and Manning, 2019), Objects365 (Shao et al., 2019), OpenImages (Kuznetsova et al., 2020) and SA-1B (Kirillov et al., 2023). We select 773k images from these sources and apply our data synthesis pipeline to construct the pretraining dataset. Table 6 presents a statistical summary of the dataset, and Figure 7 provides two illustrative examples.

For SFT, we utilize the LLaVA-NeXT dataset (Liu et al., 2024b), which comprises 779k instruction-following samples covering general VQA, OCR-related tasks and document/chart understanding. Specifically, this dataset incorporates data from AI2D (Kembhavi et al., 2016), ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), DVQA (Kafle et al., 2018), GQA, LAION-GPT4V (Schuhmann et al., 2022), OCR-VQA (Mishra et al., 2019), ShareGPT4V (Chen et al., 2024b), SynthDoG-EN (Kim et al., 2022) and Visual Genome (Krishna et al., 2017).



A black CRT computer monitor is turned on, displaying a windowed interface with lists of text. In front of the monitor sits a beige keyboard with raised keys and a blue pen resting on it. To the right of the keyboard is a white computer mouse on a blue mousepad. A black landline telephone with a coiled cord is placed near the mouse. The desk surface is wooden and cluttered with various items including a green and black notebook with a circuit board design, yellow sticky notes, a stack of white papers, and a small photo frame showing a person. Another smaller photo frame with a picture of two people is positioned near the monitor. A white desktop computer tower is under the desk. On the left side of the desk, there is a container holding pens and markers, a paper clip holder, and a tray with scissors and other office supplies. A calendar for October is hanging on the wall to the right. An orange sweater or jacket is draped over the back of a black office chair. The wall behind the desk is plain and light-colored, with visible vertical paneling



A man wearing a dark jacket, black pants, and black shoes is walking across a crosswalk. He has short dark hair and is carrying a gray backpack. The crosswalk consists of white painted stripes on the asphalt road. A pedestrian signal mounted on a metal pole displays the word "WALK" in orange letters. Behind the man, a red compact car with a visible license plate is stopped at the curb. The sidewalk is made of red brick pavers and runs alongside a modern building with large glass windows featuring a blue-tinted pattern. Several bare trees line the sidewalk and street. In the background, a tall red-brick apartment or office building with many windows stands behind a green chain-link fence. To the left, there are green-leaved trees and another multi-story building with a blue facade. The sky is bright and mostly clear. A few other vehicles are parked along the street further back.

Figure 7: Illustrative examples from the pretraining dataset. **Left:** The original image with annotated bounding boxes of visual objects. **Right:** The dense caption of the image with successfully localized textual entities. The textual entities are color-coded to match the corresponding grounded visual objects.

A.3 Benchmarks

We conduct a comprehensive evaluation across a suite of benchmarks organized into three distinct domains:

- **Visual Grounding:** We evaluate referring expression comprehension capability using RefCOCO, RefCOCO+ (Kazemzadeh et al., 2014), and RefCOCOG (Mao et al., 2016).
- **Perception-centric tasks:** We employ AI2D and ChartQA for diagram and chart understanding; OCRBench (Liu et al., 2024e) and TextVQA (Singh et al., 2019) to assess OCR capabilities; and CVBench (Tong et al., 2024), RealWorldQA (xAI Team, 2024), and V-Star (Wu and Xie, 2024) for real-world visual perception.
- **General VQA:** We encompass MMBench (Liu et al., 2024d), MME (Fu et al., 2023), MMStar (Chen et al., 2024c), and MMVet (Yu et al., 2024) for general visual understanding; MMMU (Yue et al., 2024) for multi-

disciplinary reasoning requiring college-level knowledge; and GQA for real-world compositional visual reasoning.

To ensure reproducibility, all evaluations are conducted using the Imms-eval framework (Zhang et al., 2025).

B Representation Alignment Metrics

In this section, we provide detailed definitions of the representation alignment metrics used in Section 5.1. We adopt the notation from Huh et al. (2024).

CKA CKA (Centered Kernel Alignment) reflects the *global* similarity between two representation spaces by comparing their kernel matrices. Let $\phi_i \in \mathbb{R}^n$ and $\psi_i \in \mathbb{R}^m$ be vectorized features of two modalities (e.g. language and vision). Let $\mathbf{K}_{ij} = \kappa(\phi_i, \phi_j)$ and $\mathbf{L}_{ij} = \kappa(\psi_i, \psi_j)$ be the kernel matrices computed from a dataset using some kernel-function κ . For an inner-product kernel, the ij -th entry of the centered counterpart of these ker-

nel matrices is given by

$$\begin{aligned}\bar{\mathbf{K}}_{ij} &= \langle \phi_i, \phi_j \rangle - \mathbb{E}_l[\langle \phi_i, \phi_l \rangle], \\ \bar{\mathbf{L}}_{ij} &= \langle \psi_i, \psi_j \rangle - \mathbb{E}_l[\langle \psi_i, \psi_l \rangle].\end{aligned}\quad (6)$$

Then the cross-covariance of \mathbf{K} and \mathbf{L} is:

$$\text{HSIC}(\mathbf{K}, \mathbf{L}) = \frac{1}{(n-1)^2} \text{Trace}(\bar{\mathbf{K}}\bar{\mathbf{L}}). \quad (7)$$

Finally, CKA is obtained by normalizing this quantity:

$$\text{CKA}(\mathbf{K}, \mathbf{L}) = \frac{\text{HSIC}(\mathbf{K}, \mathbf{L})}{\sqrt{\text{HSIC}(\mathbf{K}, \mathbf{K})\text{HSIC}(\mathbf{L}, \mathbf{L})}}. \quad (8)$$

CKNNA CKNNA (Centered Kernel Nearest-Neighbor Alignment) is a relaxed variation of CKA that emphasizes *local* structural alignment. It modifies the measure by replacing $\text{HSIC}(\mathbf{K}, \mathbf{L})$ with $\text{Align}(\mathbf{K}, \mathbf{L})$, which computes Eq. 7 considering only the k -nearest neighbors in the dataset:

$$\text{Align}(\mathbf{K}, \mathbf{L}) = \sum_i \sum_j \alpha(i, j) \bar{\mathbf{K}}_{ij} \bar{\mathbf{L}}_{ij}, \quad (9)$$

where $\alpha(i, j)$ is an indicator function that selects common nearest neighbors:

$$\alpha(i, j) = \mathbb{1}[\phi_j \in \text{knn}(\phi_i) \wedge \psi_i \in \text{knn}(\psi_j) \wedge i \neq j]. \quad (10)$$

This term acts as a mask, preserving only the interactions between sample i and j if j is a neighbor of i in both representation spaces.

Mutual k-NN Mutual k-NN (Mutual k-Nearest Neighbor) measures the average overlap of nearest neighbor sets of representations. Let $\{\phi_i, \psi_i\}_{i=1}^b$ denote a mini-batch of paired features from two modalities, where the collections of these features are denoted as $\Phi = \{\phi_1, \dots, \phi_b\}$ and $\Psi = \{\psi_1, \dots, \psi_b\}$. For each feature pair (ϕ_i, ψ_i) , we compute the respective nearest neighbor sets $\mathcal{S}(\phi_i)$ and $\mathcal{S}(\psi_i)$:

$$\begin{aligned}\mathcal{S}(\phi_i) &= d_{\text{knn}}(\phi_i, \Phi \setminus \phi_i), \\ \mathcal{S}(\psi_i) &= d_{\text{knn}}(\psi_i, \Psi \setminus \psi_i),\end{aligned}\quad (11)$$

where d_{knn} returns the set of indices of the k -nearest neighbors. We then measure the alignment via the average intersection:

$$m_{\text{NN}}(\phi_i, \psi_i) = \frac{1}{k} |\mathcal{S}(\phi_i) \cap \mathcal{S}(\psi_i)|, \quad (12)$$

where $|\cdot|$ denotes the cardinality of the set.

In Section 5.1, CKA is calculated using all 1000 image-caption pairs in the dataset. For CKNNA and Mutual k-NN, we report results with $k = 10$.