

# Predictive Regularization Against Visual Representation Degradation in Multimodal Large Language Models

## Supplementary Material

### 7. Implementation Details

**Training and evaluation configuration.** Our training configuration simply follows llava-1.5 [43] and llava-next [44]. Specifically, in Stage 1, we fine-tune the projector using a learning rate of 1e-3. The total batch size is 256, achieved by distributing a per-device batch size of 16 across 8 devices with a gradient accumulation step of 2. For Stage 2, two configurations are explored: if the vision encoder is frozen, the projector and LLM are fine-tuned with a learning rate of 2e-5; alternatively, if the entire model is fine-tuned, the vision encoder’s learning rate is set to 2e-6, while the remaining model components (projector and LLM) use 1e-5. Stage 2 utilizes a total batch size of 128 (per-device batch size of 4 across 8 devices, with a gradient accumulation step of 4). Common to both stages are the use of DeepSpeed-ZeRO-3 for memory optimization, a cosine learning rate scheduler, AdamW optimizer, and weight decay and warmup ratio set to 0 and 0.03, respectively. Each stage is trained for a single epoch. For PRe parameters, we set  $\lambda = 0.5$  as the default and select the middle layer of LLM as the target regularization layer, the ablation experiments are discussed in the Sec. 5 of the main paper. We evaluate most datasets using Imms-eval, except for evaluating MMVP using the standard scripts provided by Cambrian-1. A simple training algorithm of our PRe is shown in Algorithm 1

**Implementation details of analysis experiments.** In Sec. 4.1, we conduct extensive diagnostic experiments to point out the problem of visual representation degradation. For the *Global Functional Degradation* analysis (Figure 2), we use the publicly available, pre-trained weights of several popular MLLMs from huggingface: Qwen-2.5-VL-7B-Instruct, LLaVA-1.5-7B, LLaVA-1.6-7B, and InternVL-3.5-4B. This allowed us to demonstrate the generality of the degradation phenomenon across different models. For the diagnostic analyses in *Patch Structure Degradation* (Fig. 3 and Fig. 4) and *Degradation as a Trade-off for Language Capability* (Fig. 5 and Fig. 6), we trained a baseline model employs the Qwen-2.5-7B-Instruct as LLM paired with a CLIP-ViT-L/14@336 vision encoder. The vision encoder was kept frozen during the initial pre-training stage but was made trainable during the visual instruction tuning stage. All analyses in these paragraphs were conducted on this consistently trained baseline to ensure a controlled experimental environment.

**Details of ‘patch structure degradation’ paragraph.**

---

#### Algorithm 1 PRe Training Step

---

**Input:** Image batch  $\mathbf{I}$ , Text batch  $\mathbf{T}$ , MLLM  $M$ , PRe head  $f_{pred}$ .  
**Hyperparameters:** Target layer  $l_{target}$ , weight  $\lambda$ .  
**1. Get initial & intermediate visual representations**  
 $\mathbf{H}_v^0 \leftarrow M.encode\_image(\mathbf{I})$   $\triangleright$  Features after projection  
 $\mathbf{H}_t^0 \leftarrow M.tokenizer(\mathbf{T})$   $\triangleright$  Get text tokens  
 $\mathbf{H}^l \leftarrow M.forward\_llm(\mathbf{H}_v^0, \mathbf{H}_t^0)$   $\triangleright$  Get all layer hidden states  
 $\mathbf{H}_v^{l_{target}} \leftarrow GetVisualTokens(\mathbf{H}^{l_{target}})$   $\triangleright$  Visual tokens from the target layer  
**2. Define anchor (target) and online (prediction)**  
 $\mathbf{z}_{anchor} \leftarrow stop\_gradient(\mathbf{H}_v^{l_{target}})$   
 $\mathbf{p}_{online} \leftarrow f_{pred}(\mathbf{H}_v^{l_{target}})$   
**3. Compute losses**  
 $\mathcal{L}_{LM} \leftarrow LanguageModelLoss$   
 $\mathcal{L}_{PRe} \leftarrow -CosineSimilarity(\mathbf{p}_{online}, \mathbf{z}_{anchor})$   
 $\mathcal{L}_{total} \leftarrow \mathcal{L}_{LM} + \lambda \mathcal{L}_{PRe}$   
**4. Update model parameters**  
 $\mathcal{L}_{total}.backward()$   
 $optimizer.step()$

---

To understand the microscopic mechanism behind the global functional degradation, we analyze the semantic structure at the patch level. Our analysis hinges on quantifying the separability of patch representations by leveraging ground-truth segmentation masks from the COCO-Stuff dataset [8]. For each image and for the sequence of patch representations  $\mathbf{H}_v^l$  at a given layer  $l$ , we compute the following metrics:

First, we measure the Intra-Object Cohesion, which quantifies how tightly clustered the representations of patches belonging to the same object instance are. It is defined as the average cosine similarity between all unique pairs of patches within the same object class  $c$ :

$$Cohesion(\mathbf{H}_v^l) = \mathbb{E}_{c \in \mathcal{C}} [\mathbb{E}_{i, j \in \mathcal{S}_c, i \neq j} [\cos(\mathbf{h}_{v,i}^l, \mathbf{h}_{v,j}^l)]] \quad (3)$$

where  $\mathcal{C}$  is the set of all object classes in the image,  $\cos$  is cosine similarity distance, and  $\mathcal{S}_c$  is the set of patch indices belonging to class  $c$ . A higher cohesion value indicates that the model groups patches of the same object together more effectively.

Second, we measure the Inter-Object Coupling, which quantifies the degree of confusion or similarity between different objects. It is defined as the average cosine similarity between pairs of patches belonging to different object

Table 4. **The effectiveness of the predictive regularization on a 3B-LLM.** \* indicates the visual encoder is frozen during training, without \* indicates the visual encoder is trainable.

Visual Encoder	$\mathcal{L}_{\text{Pre}}$	General & Knowledge					OCR-related		Vision-centric	
		GQA	MMMU	AI2D	MMStar	SQA <sup>I</sup>	VQA <sup>T</sup>	OCRbench	RWQA	MMVP
<i>LLM: Qwen2.5-3B-Instruct</i>										
CLIP-ViT-L/14@336 *	–	61.3	40.7	62.2	41.3	71.9	43.3	318	52.5	31.3
	✓	61.3 <sup>+0.0</sup>	40.9 <sup>+0.2</sup>	62.5 <sup>+0.3</sup>	41.1 <sup>-0.2</sup>	72.5 <sup>+0.6</sup>	43.4 <sup>+0.1</sup>	320 <sup>+2</sup>	55.3 <sup>+2.8</sup>	30.7 <sup>-0.6</sup>
CLIP-ViT-L/14@336	–	61.0	39.6	62.2	42.6	72.7	45.5	359	55.6	38.0
	✓	60.9 <sup>-0.1</sup>	40.0 <sup>+0.4</sup>	62.4 <sup>+0.2</sup>	43.2 <sup>+0.6</sup>	72.9 <sup>+0.2</sup>	45.9 <sup>+0.4</sup>	362 <sup>+3</sup>	55.8 <sup>+0.2</sup>	38.7 <sup>+0.7</sup>

classes  $c$  and  $c'$ :

$$\text{Coupling}(\mathbf{H}_v^l) = \mathbb{E}_{c,c' \in \mathcal{C}, c \neq c'} [\mathbb{E}_{i \in \mathcal{S}_c, k \in \mathcal{S}_{c'}} [\cos(\mathbf{h}_{v,i}^l, \mathbf{h}_{v,k}^l)]] \quad (4)$$

A lower coupling value is desirable, as it signifies better separation between different objects.

Finally, we define the Semantic Contrast Ratio as the ratio of these two metrics. It serves as a holistic measure of patch-level separability, where a higher value is better:

$$\text{Contrast}(\mathbf{H}_v^l) = \frac{\text{Cohesion}(\mathbf{H}_v^l)}{\text{Coupling}(\mathbf{H}_v^l)} \quad (5)$$

These metrics allow us to quantitatively track how the semantic boundaries between and within objects evolve across the layers of the MLLM.

**Details of ‘degradation as a visual sacrifice for language capability’ paragraph.** In Fig. 5 (left), we quantify the Geometric Complexity using the PCA effective dimension. We compute the covariance matrix of the centered features and perform eigenvalue decomposition. The effective dimension is the minimum number of principal components  $k$  required to explain 95% of the total variance:

$$\text{EffectiveDim}(\mathbf{X}) = \min \left\{ k \mid \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^D \lambda_i} \geq 0.95 \right\} \quad (6)$$

where  $\lambda_i$  are the eigenvalues sorted in descending order. A higher effective dimension indicates a more complex and higher-capacity representation.

In Fig. 5 (right), we measure the Linear Feature Redundancy by calculating the mean absolute off-diagonal correlation. This metric quantifies the degree of statistical dependence between feature dimensions. We compute the Pearson correlation matrix  $\mathbf{C} \in \mathbb{R}^{D \times D}$  of the features, and the metric is then the average of the absolute values of its off-diagonal elements:

$$\text{Redundancy}(\mathbf{X}) = \frac{1}{D(D-1)} \sum_{i \neq j} |\mathbf{C}_{ij}| \quad (7)$$

where a lower value indicates better statistical independence (*i.e.*, less redundancy).

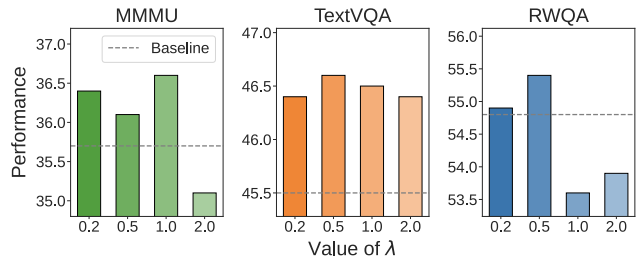


Figure 12. **The impact of  $\lambda$ .** The baseline architecture employed is CLIP-ViT-L/14@336 (frozen) and Vicuna-7B-v1.5.

## 8. More Experiments

**Results on 3B LLMs.** In Tab. 1 of the main paper, we demonstrated the effectiveness of our method across a wide range of 7B-LLM-based architectures. Here, in Tab. 4, we further supplement these findings with experimental results on 3B LLMs. Compared to the baseline, employing PRe yields performance improvements across multiple datasets. These results further underscore the importance of maintaining robust internal visual representations during MLLM’s language-driven next-token prediction training, and prove the universally effectiveness of our PRe.

**Impact of the Regularization Weight  $\lambda$ .** In Fig. 11 of the main paper, we show the impact of hyperparameter  $\lambda$  set to 0.5, 1.0, and 2.0. Here, in Fig. 12, we provide more results. We can observe that when  $\lambda$  is too small (*e.g.*, 0.2), the regularization is too weak to effectively counteract the representation degradation, resulting in only marginal gains over the baseline. As we increase  $\lambda$ , the performance improves, however, as we further increase  $\lambda$  to larger values (*e.g.*, 2.0), the performance begins to decline. This suggests that a moderate regularization strength is optimal, providing a sufficient signal to preserve visual fidelity without unduly interfering with the language modeling task. Based on this analysis, we set  $\lambda = 0.5$  as our default value for all main experiments, as it offers the best trade-off between preserving visual fidelity and maintaining language capability.

**Results using different anchor features.** In Tab. 3 of the main paper, we present results utilizing various sources

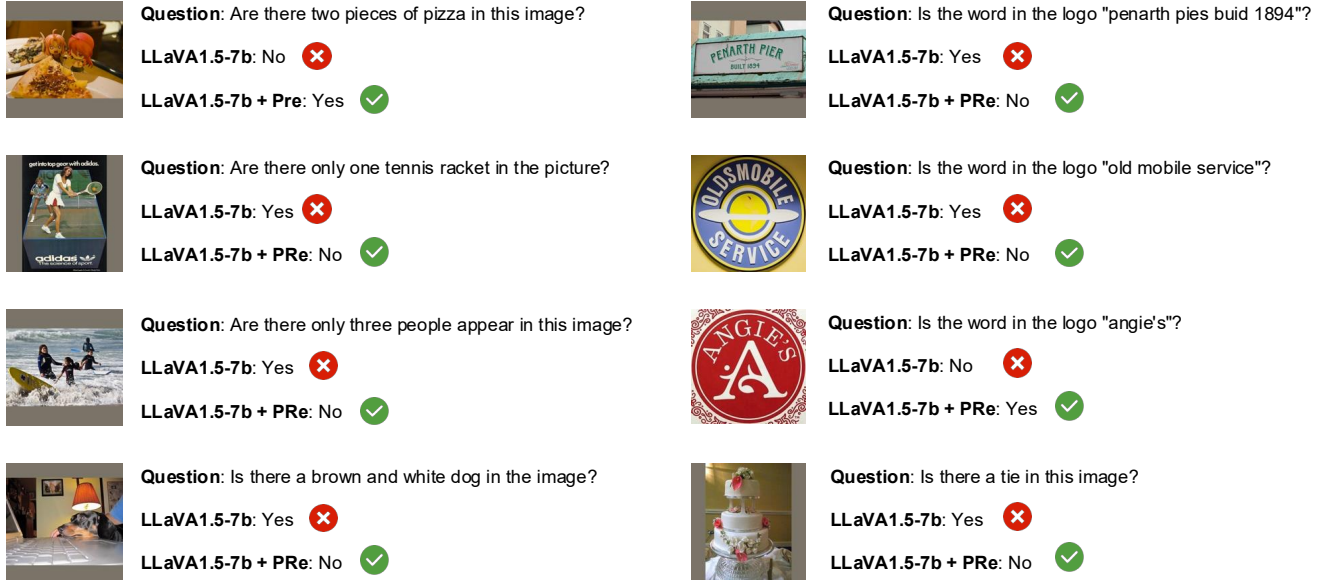


Figure 13. Case studies comparing LLaVA-1.5-7B with and without our PR method.

Table 5. **Results using different anchor features.** When the CLIP is frozen, the ‘Pre-Proj’ setting corresponds to using the penultimate layer features from a pre-trained CLIP model, together with DINOv2, DINOv3, and SAM, which form an ablation using vision foundation models

Anchor Source	GQA	MMMU	VQA <sup>T</sup>	RWQA	MMVP
<i>CLIP (frozen) + Vicuna-7B-v1.5</i>					
Baseline	62.0	35.7	45.5	54.8	20.0
Pre-LLM	62.7	<b>36.1</b>	<b>46.6</b>	<b>55.4</b>	22.0
Pre-Proj	62.7	35.1	46.4	54.4	<b>32.7</b>
DINOv2-vitb14-reg4	<b>62.8</b>	35.9	46.5	54.6	28.7
DINOv3-vitb16	62.6	35.3	46.1	54.9	28.7
SAM-vitb	62.6	34.2	<b>46.6</b>	54.4	26.0

of anchor features. Here, in Tab. 5, we further supplement these findings with results on DINOv3 and SAM. Our observations indicate that employing stronger visual foundation models, such as DINOv3, does not necessarily lead to superior performance improvements. This is likely because the feature spaces of these pure visual foundation models are optimized solely for visual tasks. Forcing the MLLM’s internal multimodal feature space to predict such a pure visual space might introduce conflicting optimization objectives. Consequently, while these approaches offer a slight improvement over the baseline, their performance does not surpass our default Pre-LLM strategy. This is because Pre-LLM’s features, being projections from the projector, strike a better balance between alignment with the LLM’s input space and retaining robust visual capabilities.

**Case studies.** In Fig. 13, we present several case studies demonstrating that the application of PR leads to notable improvements on questions related to counting, OCR,

Table 6. **Computational overhead.**

Metric	clip+qwen2.5-3b	+ PR	Overhead ( $\Delta$ )
Total Training PFLOPs	2.214	2.215	+0.045%
Throughput (samples/s)	18.74	18.58	-0.85%
Total Training Time (s)	$35.5 \times 10^3$	$35.8 \times 10^3$	+0.85%

Table 7. **Results on more datasets.**

	mme <sup>p</sup>	mmb <sup>en</sup>	mmb <sup>cn</sup>	seed <sup>i</sup>	pope(acc/f1)
clip+qwen2.5-3b	1408.3	70.1	66.1	69.9	87.5/86.3
<b>+PR</b>	1434.8	68.3	68.1	70.8	88.7/87.8
clip*+vicuna-7b	1510.7	64.3	46.4	58.6	84.7/85.8
<b>+PR</b>	1528.3	65.7	56.2	66.3	86.9/85.8

color perception, and object existence. This underscores the importance of maintaining robust internal visual features within the MLLM and validates the effectiveness of our method.

**Computational overhead.** As shown in Tab. 6, the additional training cost is negligible, and there is zero overhead during inference since the PR module is discarded after training.

**More results.** The main paper covers 9 diverse and general benchmarks spanning General Knowledge, OCR-related, and Vision-centric tasks. In Tab. 7, we provide additional results on 5 widely-used VQA datasets. These consistent gains across total 14 benchmarks confirm the effectiveness of PR.

**Generalizability to stronger encoders and higher resolutions.** We have demonstrated the robustness of PR across different **vision encoders** (CLIP, SigLIP2), **LLM backbones** (Vicuna, Qwen), **LLM scales** (7b, 3b), and

Table 8. **Results on stronger encoders and higher resolutions.**  
LLM: Qwen2.5-3b-Instruct.

	GQA	AI2D	MMStar	VQA <sup>T</sup>	RWQA
CLIP (Res: 336 <sup>2</sup> )	<b>61.3</b>	62.2	<b>41.3</b>	43.3	52.5
<b>+Pre</b>	<b>61.3</b>	<b>62.5</b>	41.1	<b>43.4</b>	<b>55.3</b>
CLIP (Res: 672 <sup>2</sup> )	60.6	62.1	42.6	45.1	52.7
<b>+Pre</b>	<b>61.0</b>	<b>62.6</b>	<b>44.0</b>	<b>45.9</b>	<b>56.7</b>
Qwen2.5-VL(NaViT)	59.4	60.7	42.5	50.0	50.5
<b>+Pre</b>	<b>60.1</b>	<b>62.2</b>	<b>43.0</b>	<b>53.2</b>	<b>51.9</b>
SigLIP2-NaFlex (NaViT)	60.8	<b>63.1</b>	43.7	42.7	51.4
<b>+Pre</b>	<b>61.4</b>	62.9	<b>44.7</b>	<b>46.4</b>	<b>53.3</b>

**training paradigms** (Frozen, Unfrozen). In Tab. 8, we further validate PRe across **resolutions**, and **stronger encoders**. The consistent gains across all these settings prove the generalizability of PRe to stronger, modern baselines.

## 9. Limitations and Future Work.

In this work, we first identify the issue of visual degradation within MLLMs and subsequently propose Predictive Regularization (PRe) to maintain robust internal visual representations, thereby enhancing performance on vision-language tasks. Our approach is directly inspired by classic and influential works in visual representation learning, such as SimSiam and JEPA. However, a potential limitation of our current PRe strategy is that using internal representations to predict the initial visual representations might inherit biases inherent to the original visual encoder. Furthermore, the broader integration of various visual representation learning methods with MLLM pre-training remains largely unexplored. For instance, the potential of contrastive learning approaches, beyond our current predictive coding-inspired method, is yet to be fully investigated within this context. In future work, we aim to further strengthen the link between advanced visual representation learning techniques and MLLM pre-training. Our goal is to continuously optimize the visual representations within MLLMs and leverage a wider spectrum of self-supervised learning paradigms to foster even more robust and versatile multimodal models.