

RECONSTRUCTIVE VISUAL INSTRUCTION TUNING

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper introduces **re**constructive visual instruction tuning (**ROSS**), a family of Large Multimodal Models (LMMs) that exploit *vision-centric supervision signals*. In contrast to conventional visual instruction tuning approaches that exclusively supervise text outputs, **ROSS** prompts LMMs to supervise *visual outputs* via reconstructing input images. By doing so, it capitalizes on the inherent richness and detail present within input images themselves, which are often lost in pure text supervision. However, producing meaningful feedback from natural images is challenging due to the heavy spatial redundancy of visual signals. To address this issue, **ROSS** employs a denoising objective to reconstruct latent representations of input images, avoiding directly regressing exact raw RGB values. This *intrinsic activation* design inherently encourages LMMs to maintain image detail, thereby enhancing their fine-grained comprehension capabilities and reducing hallucinations. Empirically, **ROSS** *consistently* brings significant improvements across different visual encoders and language models. In comparison with *extrinsic assistance* state-of-the-art alternatives that aggregate multiple visual experts, **ROSS** delivers competitive performance with a single SigLIP visual encoder, demonstrating the efficacy of our *vision-centric supervision* tailored for *visual outputs*. The code will be made publicly available upon acceptance.

1 INTRODUCTION

The success of GPT-style Large Language Models (LLMs) (Radford et al., 2018; 2019; Brown et al., 2020; OpenAI, 2023b; Yang et al., 2024a; Touvron et al., 2023; Chiang et al., 2023; Dubey et al., 2024) has motivated researchers to adapt LLMs to understand multimodal inputs (Liu et al., 2023a; 2024a; Dai et al., 2023; Bai et al., 2023). Notably, visual instruction tuning approaches (Liu et al., 2023a) demonstrate superior performance with cost-efficient training recipes. Some approaches (Chen et al., 2024b; Li et al., 2024c) even surpass GPT-4V(ision) (OpenAI, 2023a) on benchmark evaluations.

Typically, these Large Multimodal Models (LMMs) based on visual instruction tuning adopt a plug-in architecture, as depicted in Figure 1a, where pre-trained vision-language foundation models such as CLIP (Radford et al., 2021) are responsible for projecting images into visual tokens. They serve as prefix tokens for multimodal comprehension. However, this type of design, *i.e.*, *visual encoder* \rightarrow *connector* \rightarrow *LLM* \leftarrow *language instructions*, where “ \leftarrow ” indicates supervision, is primarily LLM-centric: (i) visual comprehension largely depends on vision-to-text alignment and the selected vision models, and (ii) *supervision* derives exclusively from text data. As a result, they exhibit systematic visual shortcomings such as recognizing specific visual patterns (Tong et al., 2024b).

Until very recently, some concurrent works proposed *vision-centric* solutions (Tong et al., 2024a;b). Illustrated in Figure 1b, their solutions leverage *extrinsic assistance* via aggregating several different visual experts. Inspired by the evolution in image recognition, from manually designed visual features (Sánchez & Perronnin, 2011) to learnable deep convolutional models (Krizhevsky et al., 2012), we suggest that *intrinsic activation* offers a more viable path forward. Just as deep models automatically learn hierarchical and abstract features from raw data, we believe *intrinsic activation* methods are similarly more adaptable for multimodal comprehension, reducing reliance on hand-crafted engineering, thereby enhancing both generalization and performance. Therefore, we aim to explore *intrinsic activation* solutions based on the following principles:

1. **Supervise Visual Outputs.** Current LMMs solely supervise text outputs, neglecting a significant amount of visual outputs *unused*. For instance, LLaVA-v1.5 (Liu et al., 2024a) utilizes 576 visual tokens to represent a single 336×336 image, yet their corresponding outputs

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

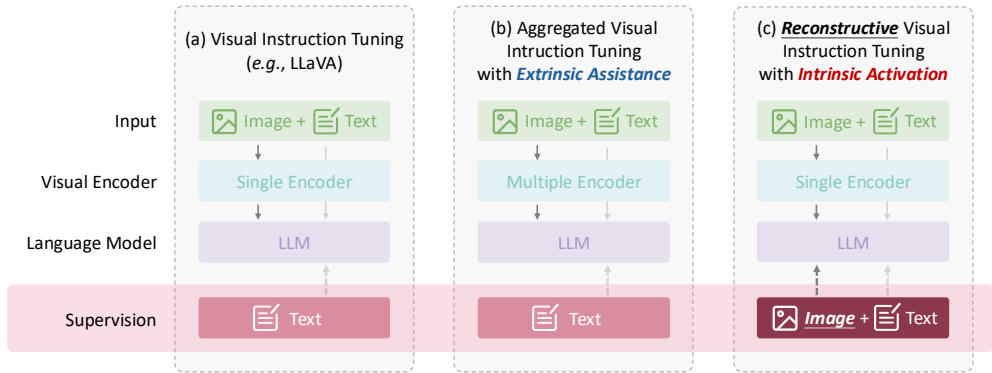


Figure 1: **Conceptual comparison** between different pipelines. (a) Typical visual instruction tuning approaches (Liu et al., 2023a; 2024a) follow a LLM-centric design that solely leverage text supervision. (b) Aggregated visual instruction tuning alternatives (Tong et al., 2024a;b) leverages *extrinsic assistance* via combining several visual experts, requiring a careful selection of visual experts. (c) Our **ROSS**, with a single visual encoder, e.g., CLIP (Radford et al., 2021) and SigLIP (Zhai et al., 2023), designs extra vision-centric reconstructive supervision as *intrinsic activation*. In this way, LMMs are required to preserve every detail of input images, thereby enhancing multimodal comprehension capabilities and reducing hallucinations.

remain *unsupervised*. Intuitively, since input images themselves inherently provide rich and detailed information, we regard LMMs *reconstructing input images* as the supervision of those visual outputs. This approach encourages LMMs to maintain low-level details, thereby enhancing their fine-grained comprehension abilities and reducing hallucinations.

2. **Explore the Optimal Formulation.** Designing this *self-supervised* task effectively is not straightforward. Motivated by the success of *masked* autoencoder (He et al., 2022) compared to its basic version denoising autoencoder (Vincent et al., 2008), we identify handling *heavy spatial redundancy of visual signals* as the underlying key factor. To this end, we formulate our approach as follows: (i) for reconstruction *targets*, instead of raw RGB pixels, we make LMMs reconstruct *latent visual tokens*, and (ii) for reconstruction *objectives*, to avoid directly regressing exact token values, we adopt per-token *denoising*.

To this end, we propose **ROSS**, termed of **reconstructive visual instruction tuning**, which utilizes input images as direct supervision signals illustrated in Figure 1c. Technically, to address the spatial redundancy inherent in natural visual signals (He et al., 2022), we train a small denoising network, which takes high-level visual outputs x as conditions to recover low-level fine-grained visual tokens z , representing an underlying distribution $p(z|x)$. These latent tokens z are derived from a frozen teacher tokenizer such as continuous VAE (Kingma, 2013) and discrete VQGAN (Esser et al., 2021). Unlike *extrinsic assistance* solutions (Tong et al., 2024a;b), our *intrinsic activation* solution naturally maintains a lightweight inference procedure. More importantly, when adapting to new visual domains, our solution avoids a careful choice of new domain-specific experts, e.g., MiDaS-3.0 (Birkel et al., 2023) for understanding depth maps, which is more efficient and easier to implement.

Empirically, **ROSS** achieves top performance across a wide range of multimodal comprehension benchmarks. Notably, our **ROSS** excels in fine-grained vision-centric benchmarks (Tong et al., 2024b; Masry et al., 2022) and hallucination benchmarks (Guan et al., 2024; Li et al., 2023c). To be specific, with a *single* SigLIP (Zhai et al., 2023) as the visual encoder, **ROSS-7B** achieves 57.3 on HallusionBench (Guan et al., 2024) and 54.7 on MMVP (Tong et al., 2024b), significantly outperforms state-of-the-art alternatives with similar model sizes which aggregate several visual experts as *extrinsic assistance*, e.g., Cambrian-1-8B (Tong et al., 2024a). In-depth analysis demonstrates the effectiveness of **ROSS** for directing focus towards visual elements and understanding depth maps. We hope our research will inspire future work in designing supervision signals for large multimodal models.

2 RELATED WORK

Visual Instruction Tuning. Most visual instruction tuning-based LMMs adopt a plug-in architecture (Liu et al., 2023a; 2024a; Bai et al., 2023; Wang et al., 2024a), where a language-supervised

visual encoder (Radford et al., 2021; Zhai et al., 2023) is responsible for extracting visual tokens. A connector is used to map those visual representations into the LLM space, *e.g.*, Resamplers (Alayrac et al., 2022), Q-Formers (Li et al., 2023b; Dai et al., 2023; Bai et al., 2023; Ge et al., 2024a), and MLPs (Liu et al., 2023a; 2024a; Li et al., 2024c; Liu et al., 2024b; Li et al., 2024a). These LLMs usually follow a two-stage training recipe. During the alignment stage, the connector is trained on high-quality caption data. Next, the full model is trained on single-image visual instruction tuning data. However, *only text outputs* are supervised. **ROSS**, on the other hand, introduces novel vision-centric supervision via reconstructing fine-grained visual tokens conditioned on *visual outputs*.

Visual Encoders for LLMs. As the original CLIP (Radford et al., 2021) adopted by conventional visual instruction tuning approaches is trained on noisy image-text pairs, it exhibits specific visual shortcomings, and thus stronger backbones (Fang et al., 2024; Zhai et al., 2023; Chen et al., 2024c) have been introduced to LLMs. Some concurrent works (Tong et al., 2024b;a) leverage *extrinsic assistance*, which further utilizes vision-only self-supervised models (Oquab et al., 2023; Wang et al., 2023a;b;c; He et al., 2022; Caron et al., 2021) and domain experts (Kirillov et al., 2023; Birkl et al., 2023; Rombach et al., 2022). **ROSS**, from a new *intrinsic activation* perspective, aims to catalyze enhanced comprehension through *reconstructing input images* with *no* extra visual experts.

Generative Objectives for LLMs. Another line of work introduces *pre-trained* text-to-image diffusion models (Rombach et al., 2022) to make LLMs capable of both comprehension and *generation* (Dong et al., 2024; Ge et al., 2024a; Sun et al., 2024b; Ge et al., 2024b; Sun et al., 2023). Our **ROSS**, with a totally different motivation, targets to catalyze multimodal comprehension via *reconstruction*. Specifically, conditions are different, where Dong et al. (2024) and Sun et al. (2024b) take outputs corresponding to *learnable queries* as conditions, while our **ROSS** takes outputs corresponding to *visual inputs*. Those methods are *generative* while **ROSS** is *reconstructive*. The detailed pipeline comparison can be found in Appendix C.

3 PRELIMINARIES

Large Multimodal Models. In the literature (Radford et al., 2018; 2019), a θ -parameterized LLM models the canonical *causal* distribution of each *text* token \mathbf{x}_i as $p_\theta(\mathbf{x}) = \prod_{i=1}^T p_\theta(\mathbf{x}_i | \mathbf{x}_{<i})$, where $\{\mathbf{x}_i\}_{i=1}^T$ represents a sequence of text tokens. To make LLMs understand visual contents, typical plug-in style LLMs (Liu et al., 2023a; 2024a) regard a sequence of visual tokens as prefix tokens. Specifically, an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ is first projected into a sequence of visual tokens by a ξ -parameterized visual encoder \mathcal{G}_ξ such as CLIP (Radford et al., 2021) and SigLIP (Zhai et al., 2023), where (H, W) indicates the spatial resolution. Then, a ϕ -parameterized multimodal projector \mathcal{H}_ϕ is utilized to project these visual tokens into the feature space of LLMs. As a result, the canonical causal distribution in a *multimodal* sentence containing an image \mathbf{I} becomes

$$p_\Theta(\mathbf{x}) = \prod_{i=1}^T p_\Theta(\mathbf{x}_i | \mathbf{x}_{<i}, \mathbf{v}), \quad \mathbf{v} = \mathcal{H}_\phi \circ \mathcal{G}_\xi(\mathbf{I}), \quad (1)$$

where $\Theta = \{\theta, \xi, \phi\}$ is the parameters and $\mathbf{v} \in \mathbb{R}^{N \times D}$ indicates the projected visual tokens. N is the number of visual tokens and D indicates the feature channel. The visual encoder \mathcal{G}_ξ could be either frozen (Liu et al., 2023a; 2024a; Tong et al., 2024a) or fine-tuned (Liu et al., 2024b; Bai et al., 2023; Li et al., 2024c; Wang et al., 2024c).

Training Recipes for LLMs. LLMs almost follow a two-stage training recipe (Liu et al., 2023a), *i.e.*, the pre-training stage (or the alignment stage) and the supervised fine-tuning stage (or the instruction tuning stage). The instruction (supervision) comes from languages such as the answers to VQA tasks, maximizing the log-likelihood of *text* outputs:

$$\mathcal{L}_{\text{LLM}}^{\text{text}}(\Theta = \{\theta, \xi, \phi\}, \mathbf{x}, \mathbf{I}) = \frac{-1}{T - N} \sum_{i=N+1}^T \log p_\Theta(\mathbf{x}_i | \mathbf{x}_{<i}, \mathbf{v}), \quad (2)$$

where N represents the number of visual tokens and visual outputs (one input token corresponds to one visual output). From Equation (2), we can tell that only text outputs $\mathbf{x}_{i>N}$ are supervised.

4 ROSS: RECONSTRUCTIVE VISUAL INSTRUCTION TUNING

In this section, we first provide an overview of our reconstructive visual instruction tuning (**ROSS**). Then, we discuss our explorations towards the optimal formulation in the following subsections, with the ultimate goal of *handling spatial redundancy of visual signals* to provide meaningful visual supervision. Our explorations mainly include reconstruction *targets* and the training *objective*.

Overview. Illustrated in Figure 2, the overall philosophy of our **ROSS** is to construct *reconstructive* visual supervision signals on visual outputs $\mathbf{x}_{i \leq N}$. The training objective includes (i) the original next-token prediction on $\mathbf{x}_{i > N}$ shown in the right part of Figure 2, and (ii) another *reconstructive term* in the left part of Figure 2, i.e., $\mathcal{L}_{\text{ROSS}} = \mathcal{L}_{\text{LMM}}^{\text{text}} + \mathcal{L}_{\text{LMM}}^{\text{visual}}$. Specifically, this visual term could be any custom measurements \mathcal{M} between $\mathbf{x}_{i \leq N}$ and specific reconstruction targets of image \mathbf{I} :

$$\begin{aligned} \mathcal{L}_{\text{LMM}}^{\text{visual}}(\Theta = \{\theta, \xi, \phi, \pi\}, \mathbf{x}, \mathbf{I}) \\ = \mathcal{M}(\mathcal{J}_{\pi}(\mathbf{x}_{i \leq N}), \mathcal{F}(\mathbf{I})), \end{aligned} \quad (3)$$

where \mathcal{J}_{π} indicates the π -parameterized post projection that maps the dimensions of visual tokens $\mathbf{x}_{i \leq N}$ to be consistent with the teacher tokenizer \mathcal{F} .

Variants of ROSS. Evidently, different choices of \mathcal{F} and \mathcal{M} contribute to different variants. \mathcal{F} controls the reconstruction *target* while \mathcal{M} defines the *objective*:

1. Towards the *target*, \mathcal{F} can be the patchify operation (Dosovitskiy et al., 2021), resulting in *pixel-level* reconstruction, or pre-trained fine-grained visual tokenizers such as VAE (Kingma, 2013) and VQGAN (Esser et al., 2021), leading to *latent-level* reconstruction. \mathcal{F} could even be vision-only models such as DINOv2 (Oquab et al., 2023), making LMMs learn specific visual patterns from \mathcal{F} , which is also a type of *latent-level* reconstruction.
2. Towards the *objective*, the most straightforward choice of \mathcal{M} is MSE or cosine similarity for *regressing* raw pixel values or latent features, respectively. We also explore the *denoising* objective (Ho et al., 2020) to avoid being overwhelmed by fitting exact values.

We introduce our explorations step by step in the following sections. The ultimate goal of our exploration is to design an appropriate self-supervised reconstructive pre-text task that provides meaningful vision-centric supervision signals to LMMs, where handling the *spatial redundancy* of visual signals (He et al., 2022) becomes the crux.

4.1 ROSS^R: REGRESSING AS RECONSTRUCTIVE VISUAL INSTRUCTION

In this section, we introduce straightforward variants, i.e., *regressing* as reconstructive visual instruction. As shown in Figure 3, depending on the choice of \mathcal{F} , it mainly has three variants: (a) ROSS^R-Pixel, (b) ROSS^R-Latent, and (c) ROSS^R-Latent2Pixel.

Directly Regressing Raw RGB Values. The most straightforward variant is to directly regress raw RGB values illustrated in Figure 3a, called “ROSS^R-Pixel”. Under such a setting, \mathcal{F} is the patchify operation (Dosovitskiy et al., 2021), reshaping the image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ into a sequence of flattened 2D patches $\mathbf{I}_p \in \mathbb{R}^{N \times (3P^2)}$, where (P, P) is the resolution of each image patch and $N = HW/P^2$ indicates the resulting number of patches. \mathcal{J}_{π} can be a simple MLP, mapping the dimension of visual outputs $\mathbf{x}_{i \leq N}$ from D to $3P^2$. The measurement \mathcal{M} is MSE. However, as visual signals suffer from *heavy spatial redundancy* (He et al., 2022), such a design may not provide meaningful supervision to LMMs. An intuitive alternative to avoid directly regressing raw RGB values while still reconstructing the image is to urge LMMs to reconstruct *latent tokens*, introduced as follows.

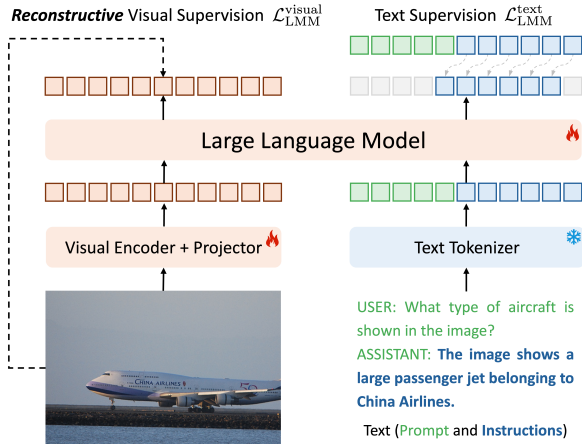


Figure 2: **Overview of ROSS.** Given an input image and the corresponding text to this image, **ROSS** aims to supervise visual outputs by reconstruction.

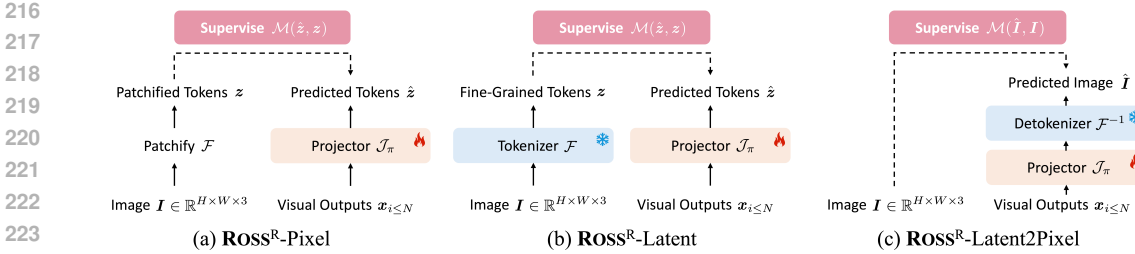


Figure 3: Variants of ROSS^R , where *regression* objectives are either computed on raw RGB values in (a) and (c), or specific latent space determined by \mathcal{F} in (b). We adopt MSE as \mathcal{M} for *pixel* regression in (a) and (c), and cosine-similarity for *latent* regression in (b), respectively.

Regressing Latent Tokens. Illustrated in Figure 3b, ROSS^R -Latent aims to regress fine-grained *latent* tokens extracted by the teacher tokenizer. \mathcal{F} can be models trained with discriminative tasks such as DINOv2 (Oquab et al., 2023) and DEIT-III (Touvron et al., 2022). The *encoder* part of models trained with reconstruction tasks such as VQGAN (Esser et al., 2021) and VAE (Kingma, 2013) are also capable. \mathcal{M} here is the cosine-similarity. Intuitively, the *decoder* part of the latter is able to remap latent tokens into the pixel space. Therefore, supervising in the pixel space via *decoding* becomes another valid variant introduced as follows.

Regressing RGB Values via Decoding. Shown in Figure 3c, ROSS^R -Latent2Pixel requires a *decoder* to project predicted latent tokens \hat{z} into the RGB pixel space, resulting in predicted image \hat{I} . Let \mathcal{F}^{-1} be the *decoder* part of VQGAN (Esser et al., 2021) or VAE (Kingma, 2013), and the *regressive* MSE objective \mathcal{M} is performed on pixel-space. Note that we simply use \mathcal{F}^{-1} to represent the decoding process, which is actually *not* the inverse function of \mathcal{F} mathematically.

Discussion. Recall that we need to find the optimal solution to address the spatial redundancy of natural visual signals, the *target-level* exploration above achieves this goal *partially*, as the *objective* is limited to vanilla regression. To this end, inspired by Ho et al. (2020) and Li et al. (2024e), we further incorporate a novel *denoising* objective in the following section.

4.2 ROSS^D : DENOISING AS RECONSTRUCTIVE VISUAL INSTRUCTION

As an objective for handling *heavy spatial redundancy* to provide meaningful vision-centric supervision signals, denoising is better than vanilla regressing, since the introduction of noise into the training data acts as an implicit form of data augmentation and regularization. The denoising process encourages the model to focus on the underlying data manifold rather than memorizing specific instance values (Chen et al., 2023c; Song & Ermon, 2019; Karras et al., 2022; Yang et al., 2024b).

Technically, as illustrated in Figure 4a, our final ROSS^D takes high-level visual outputs $\mathbf{x}_{i \leq N}$ as conditions to recover *clean* fine-grained tokens \mathbf{z}_0 from *noisy* tokens \mathbf{z}_t . Specifically, clean tokens $\mathbf{z}_0 = \mathcal{F}(I)$ are obtained from the teacher tokenizer \mathcal{F} . By default, we utilize a continuous VAE (Kingma, 2013) regularized by Kullback–Leibler (KL) divergence provided by Rombach et al. (2022), since it is believed to capture sufficient image details. The training procedure of the denoiser \mathcal{J}_π follows a diffusion process (Ho et al., 2020):

$$\mathcal{L}_{\text{LMM}}^{\text{visual}}(\Theta = \{\theta, \xi, \phi, \pi\}, \mathbf{x}, I) = \mathbb{E}_{t, \epsilon} [||\mathcal{J}_\pi(\mathbf{z}_t; \mathbf{x}_{i \leq N}, t) - \epsilon||^2]. \quad (4)$$

The denoiser \mathcal{J}_π actually estimates the conditional expectation $\mathbb{E}[\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) | \mathbf{z}_t]$. More details about the background knowledge of diffusion models can be found in Appendix A.

Architecture of the Denoiser. As conditions $\mathbf{x}_{i \leq N}$ are *causal*, we introduce a self-attention module to model the inter-token dependencies illustrated in Figure 4b. Specifically, the architecture of the denoiser \mathcal{J}_π is a stack of Transformer Encoder blocks (Vaswani et al., 2017) and each block contains three extra projections for conditions $\mathbf{x}_{i \leq N}$, inputs \mathbf{z}_t , and timesteps t , respectively.

Choices of the Teacher Tokenizer. By default, we adopt *latent* denoising and we take a continuous tokenizer provided by Rombach et al. (2022) as \mathcal{F} , since it manages to reconstruct input images with a low rFID (Heusel et al., 2017) and thus it is expected to preserve many low-level details of input images. This extra reconstructive objective, however, is *not* limited to any certain tokenizer \mathcal{F} . Discrete tokenizers such as VQGAN (Esser et al., 2021), and vision self-supervised models

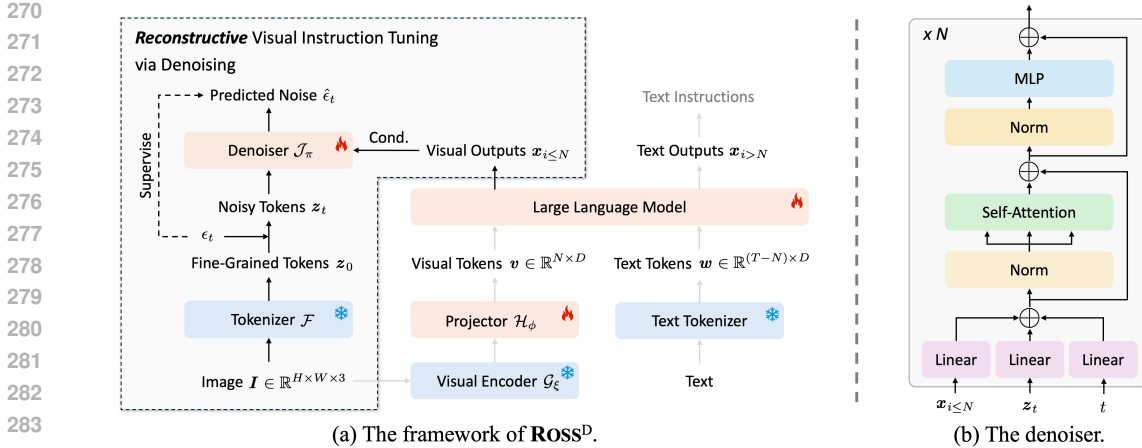


Figure 4: Illustration of (a) the training procedure of **ROSS^D** and (b) the detailed architecture of the denoiser \mathcal{J}_π . (a) **ROSS^D** introduces visual guidance via *denoising fine-grained visual tokens z_0 conditioning on visual outputs $x_{i \leq N}$* . (b) The denoiser takes noisy tokens z_t , current timesteps t , and conditions $x_{i \leq N}$ as inputs and outputs the predicted noise $\hat{\epsilon}_t$. Each denoiser block consists of three linear projection layers and a standard self-attention block (Vaswani et al., 2017).

such as DINOv2 (Oquab et al., 2023), are also qualified to be the tokenizer. Even the patchify operation (Dosovitskiy et al., 2021) is capable, resulting in *pixel* denoising.

5 EXPERIMENTS

5.1 ABLATION STUDY

Implementation Details. All ablation studies are implemented based on LLaVA-v1.5 (Liu et al., 2024a). The visual encoder \mathcal{G}_ξ is CLIP-ViT-L/14@336 (Radford et al., 2021) and the base LLM is Qwen2-7B-Instruct (Yang et al., 2024a). The training data is LLaVA-558K (Liu et al., 2023a) and Cambrian-737K (Tong et al., 2024a) for the pre-training stage and the instruction tuning stage, respectively. We evaluate our each variant of **ROSS** mainly on (i) hallucination: POPE (Li et al., 2023c) and HallusionBench (Guan et al., 2024), (ii) fine-grained comprehension: MMVP (Tong et al., 2024b) and ChartQA (Masry et al., 2022), and (iii) general comprehension: MMBench (Liu et al., 2023b) English dev split. All evaluations are conducted with VLMEvalKit (Duan et al., 2024). Evaluation prompts can be found in Appendix B.

Pixel Regression v.s. Latent Regression. Starting from the visual instruction tuning baseline (Liu et al., 2023a; 2024a), we first explore the effectiveness of using *regression* as the objective for our reconstructive visual instruction tuning. We utilize a continuous VAE (Kingma, 2013) with an encoder-decoder architecture provided by Rombach et al. (2022), where the *encoder* part serves as \mathcal{F} for **ROSS^R-Latent** while the *decoder* part is \mathcal{F}^{-1} for **ROSS^R-Latent2Pixel**. As illustrated in Figure 5, our vision-centric regression supervision outperforms the visual instruction tuning baseline in most cases. Moreover, latent regression performs the best since *regressing raw RGB pixels fails to provide meaningful supervision signals*, regardless of whether utilizing a decoder or not.

Choices of \mathcal{F} . We study the effectiveness across different latent teacher tokenizers \mathcal{F} in Figure 6, including KL-16 provided by Rombach et al. (2022), which is a continuous VAE (Kingma, 2013) with Kullback–Leibler (KL) divergence, self-supervised DINOv2 (Oquab et al., 2023), fully-supervised DEiT-III (Touvron et al., 2022), and language-supervised EVA02CLIP (Fang et al., 2024). Among them, KL-16 is the best choice. One intuitive explanation is that it is expected to preserve the most image details, since it was originally designed to accurately reconstruct input images.

Regression v.s. Denoising. In Figure 7, we study the effectiveness of the *denoising* objective over vanilla regression across different tokenizers, *i.e.*, KL-16 (Rombach et al., 2022) and DINOv2 (Oquab et al., 2023). Notably, even if **ROSS^R-Latent** (KL-16) has already outperformed the visual instruction tuning baseline by a large margin, **ROSS^D** manages to bring significant improvements by replacing regression with denoising. Therefore, *denoising is better at handling visual spatial redundancy*.

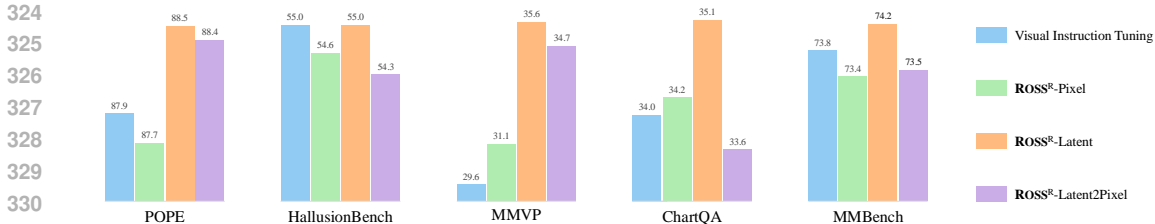


Figure 5: **Pixel Regression v.s. Latent Regression.** The teacher tokenizer \mathcal{F} for ROSS^R-Latent is the *encoder* of a continuous VAE (Kingma, 2013) provided by Rombach et al. (2022), while its *decoder* serves as \mathcal{F}^{-1} for ROSS^R-Latent2Pixel. Our vision-centric reconstructive supervision surpasses the visual instruction tuning baseline in most cases. Among three regression variants, ROSS^R-Latent performs the best, as it avoids explicitly regressing redundant raw RGB values.

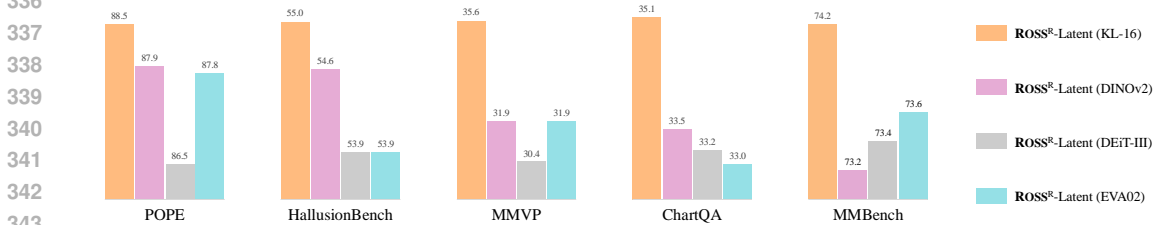


Figure 6: **Choices of the latent teacher tokenizer \mathcal{F} .** KL-16 (Rombach et al., 2022) is the best tokenizer as it is originally used for *reconstruction*, and it is expected to preserve the most image details. Other alternatives are utilized for classification (Touvron et al., 2022), instance-level representation learning (Oquab et al., 2023), and language alignment (Fang et al., 2024), respectively.

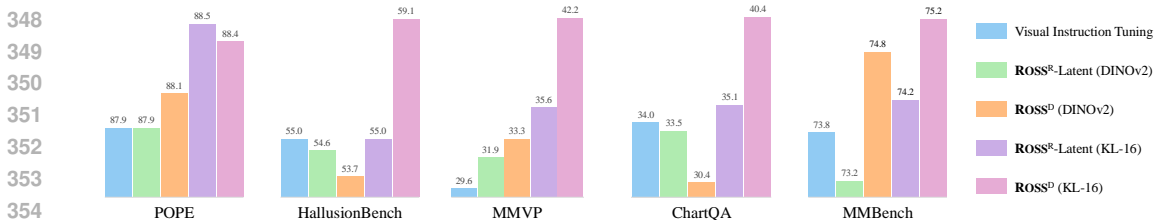


Figure 7: **Regression v.s. Denoising.** With KL-16 as the tokenizer, the denoising objective introduced in Equation (4) brings significant improvements over vanilla regression using MSE as it avoids overfitting exact latent token values, even if ROSS^R-Latent (KL-16) has already outperformed the visual instruction tuning baseline by a large margin.

Finally, we leverage the insights and conclusions from all our previous studies to train our ROSS. Specifically, we regard the optimal formulation ROSS^D (KL-16), *i.e.*, denoising with the KL-16 tokenizer, as our final ROSS. Please refer to Appendix C.1 for ablations on the architecture of the denoiser, continuous tokenizer v.s. discrete tokenizer, and the denoising schedule.

5.2 IN-DEPTH ANALYSIS

Attention Analysis. We compute the attention scores of the *last* token with respect to *all visual tokens* on MMVP (Tong et al., 2024b). Quantitative and qualitative comparisons between the visual instruction tuning baseline (LLaVA) (Liu et al., 2024a) and our ROSS are provided in Table 1 and Figure 8, respectively. Table 1 reveals that the attention scores achieved by ROSS are *significantly higher* than those of LLaVA, indicating that the inclusion of vision-centric reconstructive objective $\mathcal{L}_{LMM}^{visual}$ effectively directs focus towards input images, thereby enhancing the comprehending visual signals. Similarly, Figure 8 demonstrate that the implementation of $\mathcal{L}_{LMM}^{visual}$ enables the alignment of attention closely with the *relevant visual elements* corresponding to the text query.

Table 1: **Quantitative comparison on attention values.** We conduct a T-test (Student, 1908) to compare the *means* and a Mann-Whitney U test (Mann & Whitney, 1947) to compare the *medians* of the two distributions. The mean and median of ROSS are both *significantly higher* than those of LLaVA.

Statistic ($\times 10^{-4}$)	LLaVA	ROSS	P-value
Mean	2.03	2.36	1.27×10^{-7}
25th Percentile	1.50	1.81	–
Median	1.90	2.26	4.39×10^{-9}
75th Percentile	2.42	2.76	–
95th Percentile	3.51	3.69	–

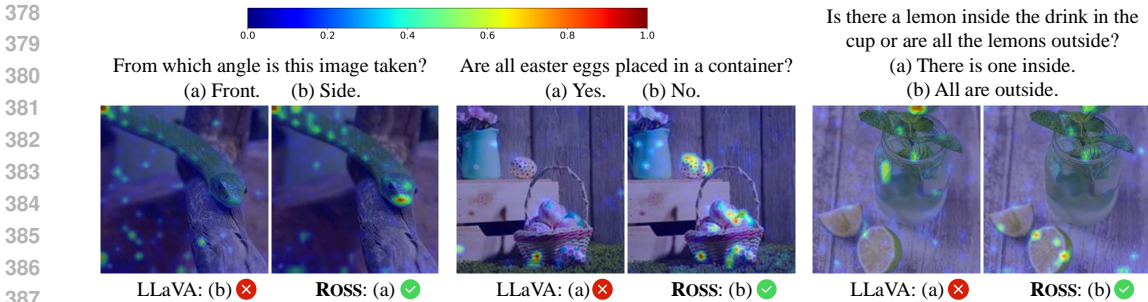


Figure 8: **Qualitative comparison on attention maps on MMVP (Tong et al., 2024b)**, where we keep the *same* LLM and training data. With extra vision-centric supervision signals, **ROSS** urges the model to *focus on specific image contents corresponding to the question with higher attention values*.

Table 2: **Generative v.s. Reconstructive**. Following Sun et al. (2024b) and Dong et al. (2024), we adopt 576 learnable latent tokens to *query* the LMM and utilize the corresponding outputs as conditions to the denoiser for generative cases. Extra 102K caption data from ShareGPT4V (Chen et al., 2023b) is introduced to the original SFT data, facilitating text-to-image creation. Reconstructive objectives boost comprehension while generative alternatives *cannot*.

Method	SFT Data	w/ $\mathcal{L}_{LMM}^{visual}$		Hallucination		Fine-Grained		General
		737K	102K	POPE	Hallu.	MMVP	ChartQA	MMB ^{EN}
Baseline	737K + Caption 102K	–	–	86.2	55.1	32.0	30.9	74.4
Reconstructive	737K + Caption 102K	✓	✓	87.6	58.0	38.7	40.4	75.2
Reconstructive	737K + Caption 102K	–	✓	87.6	56.3	37.3	39.7	74.9
Generative	737K + Creation 102K	–	✓	85.4	52.0	30.0	31.2	73.9

Table 3: **The effectiveness of the vision-centric supervision among various LLMs and visual encoders**, where $\mathcal{L}_{LMM}^{visual}$ manages to bring significant improvements *consistently*.

Language Model	$\mathcal{L}_{LMM}^{visual}$	POPE	Hallu.	MMVP	ChartQA	OCR Bench	MMB ^{EN}
<i>Visual Encoder: CLIP-ViT-L/14@336</i>							
Vicuna-7B-v1.5	–	86.3	52.5	28.0	32.9	339	67.0
	✓	87.2 ↑ 0.9	55.8 ↑ 3.3	36.3 ↑ 8.3	39.8 ↑ 6.9	350 ↑ 11	67.6 ↑ 0.6
Qwen2-7B-Instruct	–	87.9	55.0	29.6	34.0	363	73.8
	✓	88.4 ↑ 0.5	59.1 ↑ 4.1	42.2 ↑ 12.6	40.4 ↑ 6.4	381 ↑ 18	75.2 ↑ 1.4
<i>Visual Encoder: SigLIP-ViT-SO400M/14@384</i>							
Vicuna-7B-v1.5	–	86.0	50.4	27.3	36.2	354	64.5
	✓	86.8 ↑ 0.8	53.2 ↑ 2.8	38.0 ↑ 10.7	41.6 ↑ 5.4	365 ↑ 11	65.7 ↑ 1.2
Qwen2-7B-Instruct	–	88.5	57.3	40.7	44.4	432	76.3
	✓	88.7 ↑ 0.2	58.2 ↑ 0.9	49.3 ↑ 8.6	46.3 ↑ 1.9	448 ↑ 16	76.9 ↑ 0.6

Generative v.s. Reconstructive. We ablate the effectiveness of reconstruction over generation in Table 2. Similar to Sun et al. (2024b) and Dong et al. (2024), for generative cases, we adopt 576 learnable latent tokens to *query* the LMM and utilize the corresponding outputs as conditions to the denoiser. The detailed pipeline of these two methods can be found at Figure 11 in Appendix C.1. However, generative methods require *specific creation data* and can *not* be naively implemented on the original SFT data. To build creation data, we utilize GPT-4o to transfer 102K *caption* into *text-to-image* creation data from ShareGPT4V (Chen et al., 2023b) and combine them with the original SFT data. From Table 2, we can tell that reconstructive objectives boost comprehension while generative alternatives *cannot*. An intuitive explanation of this evidence can be found in Appendix C.1.

ROSS with Different LLMs and Visual Encoders. To demonstrate the effectiveness of our vision-centric supervision $\mathcal{L}_{LMM}^{visual}$ adopted by our ROSS, we conduct systematic experiments across different base LLMs and visual encoders. From Table 3, ROSS contributes to significant improvements *consistently*, especially on fine-grained comprehension benchmarks, *i.e.*, MMVP (Tong et al., 2024b) and ChartQA (Masry et al., 2022). Extended experiments on more representative benchmarks can be found at Table 12 in Appendix C.1.

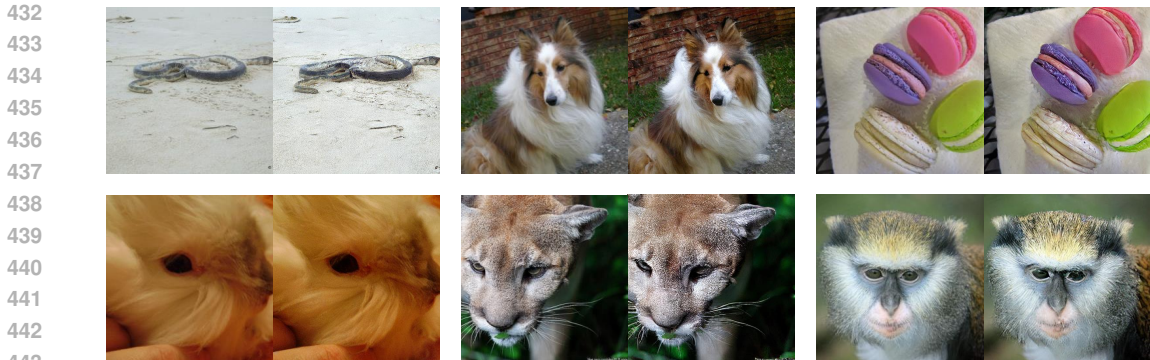


Figure 9: **Reconstruction results** on ImageNet-1K (Deng et al., 2009) validation set. For each tuple, we show the input image (left) and the reconstructed image (right). Reasonable reconstruction results demonstrate that *high-level features of ROSS-7B can be projected back into the pixel space*.

Table 4: **Comparison to state-of-the-art LMMs**. A mixture of 2M caption data and 1.2M instruction tuning data are utilized for pre-training and fine-tuning, respectively. Our model outperforms them in most of the settings. We evaluate these models on: POPE (Li et al., 2023c) averaged accuracy, Hallu.: HallusionBench (Guan et al., 2024) average accuracy, MMB^{EN}: MMBench (Liu et al., 2023b) English dev split, MMB^{CN}: MMBench (Liu et al., 2023b) Chinese dev split, SEED^I: SEED-Bench-1 (Li et al., 2023a) with image accuracy, MMMU (Yue et al., 2024) validation split, MMVP (Tong et al., 2024b), GQA (Hudson & Manning, 2019) test-dev-balanced split, and AI2D (Hiippala et al., 2021) test split. †We evaluate the official checkpoint/api using VLMEvalKit (Duan et al., 2024).

Model	POPE	Hallu.	MMB ^{EN}	MMB ^{CN}	SEED ^I	MMMU	MMVP	GQA	AI2D
GPT-4V-1106 (OpenAI, 2023a)	75.4	65.8 [†]	75.8	75.1 [†]	71.6	53.8	50.0	36.8	78.2
Gemini-1.5 Pro (Team et al., 2023)	–	–	73.6	–	70.7	47.9	–	–	–
MM-1-8B (McKinzie et al., 2024)	86.6	–	72.3	–	69.9	37.0	–	72.6	–
Mini-Gemini-8B (Li et al., 2024f)	–	–	72.7	–	73.2	37.3	18.7	64.5	73.5
DeepSeek-VL-7B (Lu et al., 2024)	85.8 [†]	44.1 [†]	73.2	72.8	70.4	36.6	–	–	64.9 [†]
Cambrian-1-8B (Tong et al., 2024a)	87.4 [†]	48.7 [†]	75.9	68.9 [†]	74.7	42.7	51.3	64.6	73.0
ROSS-7B	89.2	57.3	79.0	76.1	73.0	43.4	54.7	65.5	76.8
<i>Base LLM: Vicuna-7B-v1.5</i>									
LLaVA-v1.5-7B [†] (Liu et al., 2024a)	86.2	47.5	65.5	58.5	66.0	34.4	20.0	62.0	55.4
LLaVA-v1.6-7B [†] (Liu et al., 2024b)	86.5	35.8	67.4	60.1	70.2	35.8	37.3	64.2	67.1
ROSS-7B_{vicuna}	88.2	55.2	67.7	61.3	67.6	36.9	39.3	63.7	69.3
<i>Base LLM: Vicuna-13B-v1.5</i>									
LLaVA-v1.5-13B [†] (Liu et al., 2024a)	82.5	44.9	68.8	63.6	68.2	36.6	32.0	63.3	60.8
LLaVA-v1.6-13B [†] (Liu et al., 2024b)	86.2	36.7	70.0	64.1	71.9	36.2	35.3	65.4	72.4
Mini-Gemini-13B (Li et al., 2024f)	–	–	68.6	–	73.2	37.3	19.3	63.7	70.1
Cambrian-1-13B (Tong et al., 2024a)	85.7 [†]	54.0 [†]	75.7	65.9 [†]	74.4	40.0	41.3	64.3	73.6
ROSS-13B_{vicuna}	88.7	56.4	73.6	67.4	71.1	41.3	44.7	65.2	73.8

Reconstruction Results. We fine-tune the denoiser to recover latent tokens from a frozen KL-16 provided by Rombach et al. (2022) conditioned on ROSS-7B features on ImageNet-1K (Deng et al., 2009) for *only five epochs*, where the denoiser manages to produce reasonable reconstruction results as illustrated in Figure 9. This interesting finding demonstrates that high-level ROSS-7B features *actually contain image details*. We hope this finding will inspire future work.

Computational Overhead. The denoising process introduces a negligible increase in training time ($\approx 10\%$ compared to the baseline), while the benefits outweigh the minor additional costs. Please refer to Table 10 in Appendix B for details.

5.3 COMPARISON WITH STATE-OF-THE-ARTS

ROSS utilizes a *single* SigLIP-ViT-SO400M/14@384 (Zhai et al., 2023) as the visual encoder. ROSS-7B utilizes Qwen2-7B-Instruct (Yang et al., 2024a) and ROSS-13B_{vicuna} adopts Vicuna-13B-v1.5 (Chiang et al., 2023) as the base LLM. The implementation almost follows LLaVA-v1.5 (Liu et al., 2024a) *without* the high-resolution image-slicing technique (Liu et al., 2024b). Thus, our primary compar-

Table 5: **Transfer learning** on SpatialBench (Cai et al., 2024). “RGB” indicates using only RGB images for testing, while “RGB + D” represents taking depth maps as extra inputs. The performance of GPT-4o is obtained from Cai et al. (2024). *LMMs can better comprehend depth maps with $\mathcal{L}_{LMM}^{visual}$* .

Method	Test Inputs	$\mathcal{L}_{LMM}^{visual}$	MiDaS	Size	Reaching	Position	Existence	Counting	Average
LLaVA	RGB	–	–	20.0	51.7	58.8	70.0	74.6	55.0
	RGB + D	–	–	21.7 \uparrow 1.7	45.0 \downarrow 6.7	58.8 – 0.0	65.0 \downarrow 5.0	77.7 \uparrow 3.1	53.6 \downarrow 1.4
	RGB	–	✓	21.7	60.0	64.7	80.0	84.1	62.1
	RGB + D	–	✓	21.7 – 0.0	51.7 \downarrow 8.3	70.6 \uparrow 5.9	65.0 \downarrow 15.0	91.1 \uparrow 7.0	60.0 \downarrow 2.1
ROSS	RGB	✓	–	25.0	53.3	64.7	70.0	75.3	57.7
	RGB + D	✓	–	28.3 \uparrow 3.3	65.0 \uparrow 11.7	67.6 \uparrow 2.9	85.0 \uparrow 15.0	84.6 \uparrow 8.7	66.1 \uparrow 8.4
GPT-4o	RGB	–	–	43.3	51.7	70.6	85.0	84.5	67.0
	RGB + D	–	–	40.0 \downarrow 3.3	51.7 – 0.0	61.8 \downarrow 8.8	90.0 \uparrow 5.0	85.2 \uparrow 0.7	65.7 \downarrow 1.3

ison of **ROSS** with alternative methods focuses on benchmarks that do *not* require exceptionally high-resolution inputs. We use a mixture of 2M caption data for the pre-training stage, which consists of 1246K from ShareGPT4V (Chen et al., 2023b) and 707K from ALLaVA (Chen et al., 2024a). The instruction tuning data is a mixture of Cambrian-737K (Tong et al., 2024a) and SMR-473K (Zhang et al., 2024). We further incorporate our **ROSS** with the “anyres” technique (Liu et al., 2024b) and compare with others on high-resolution benchmarks at Table 13 in Appendix C.2.

Illustrated in Table 4, we compare our **ROSS** with both private models (OpenAI, 2023a; Team et al., 2023; McKinzie et al., 2024) and open-sourced alternatives (Liu et al., 2024a;b; Tong et al., 2024a; Li et al., 2024f; Lu et al., 2024). The previous open-source state-of-the-art Cambrian-1 (Tong et al., 2024a) leverages *extrinsic assistance* that aggregates CLIP (Radford et al., 2021), SigLIP (Zhai et al., 2023), DINOv2 (Oquab et al., 2023), and ConvNext (Liu et al., 2022). On the other hand, our **ROSS** stands for *intrinsic activation*. With only a single SigLIP (Zhai et al., 2023) model as the visual encoder, our **ROSS** surpasses Cambrian-1 (Tong et al., 2024a), under most cases, *without* a careful choice of the visual experts and naturally maintains a lightweight inference procedure. **ROSS** is also data-efficient compared with Cambrian-1 (Tong et al., 2024a), since it requires 7M instruction tuning data. Notably, **ROSS-7B** even surpasses GPT-4V-1106 and Gemini-1.5 Pro on several benchmarks such as POPE (Li et al., 2023c), MMBench (Liu et al., 2023b), and MMVP (Tong et al., 2024b).

5.4 APPLICATIONS

Transfer Learning on Understanding Depth Maps. We further evaluate the transfer learning capability of our **ROSS** on SpatialBench (Cai et al., 2024), which requires the model to understand *depth maps*. We compare our **ROSS** with the visual instruction tuning baseline, with the *same* training data and model architecture. Also, we compare the effectiveness of the *extrinsic assistance* solution, *i.e.*, combining a depth expert MiDaS-3.0 (Birkl et al., 2023) to visual instruction tuning, with our *intrinsic activation* solution. Specifically, the pre-training data is LLaVA-558K (Liu et al., 2023a) and the fine-tuning data is SpatialQA-853K (Cai et al., 2024), where each conversation contains the RGB image and the *depth* maps extracted by ZoeDepth (Bhat et al., 2023). The visual encoder is CLIP-ViT-L/14@336 (Radford et al., 2021) and the base LLM is Qwen2-7B-Instruct (Yang et al., 2024a). As demonstrated in Table 5, our **ROSS** manages to make use of the extra depth map, as consistent and significant improvements are observed when taking “RGB + D” inputs for testing. Extrinsic assistance approaches *cannot* take advantage of extra depth maps when testing. Even GPT-4o *cannot* fully understand depth maps. *Qualitative results can be found at Figure 16 in Appendix C.4.*

6 CONCLUSION

This paper introduces reconstructive visual instruction tuning (**ROSS**), leveraging a vision-centric *reconstructive* objective to supervise visual outputs. To avoid being overwhelmed by heavily redundant raw RGB values, we train a denoiser to recover clean latent visual representations conditioning on visual outputs. Experimentally, the proposed objective indeed brings enhanced comprehension capabilities and reduced hallucinations. **ROSS** outperforms the state-of-the-art under most cases with only a *single* SigLIP (Zhai et al., 2023) as the visual encoder. *The in-depth analysis demonstrates that high-level features from ROSS-7B actually contain sufficient details for low-level image reconstruction. This finding reveals the possibility of equipping comprehension LMMs with the ability of naive generation without the help of generation experts such as Stable Diffusion (Rombach et al., 2022).*

REFERENCES

- 540
541
542 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur
543 Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot
544 learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- 545 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and
546 Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint*
547 *arXiv:2308.12966*, 2023.
- 548 Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer
549 by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023.
- 550 Reiner Birkel, Diana Wofk, and Matthias Müller. Midas v3. 1—a model zoo for robust monocular relative depth
551 estimation. *arXiv preprint arXiv:2307.14460*, 2023.
- 552 Ali Furkan Biten, Ron Litman, Yusheng Xie, Srikanth Appalaraju, and R Manmatha. Latr: Layout-aware
553 transformer for scene-text vqa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
554 *Recognition (CVPR)*, 2022.
- 555 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
556 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
557 *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 558 Wenxiao Cai, Yaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and Bo Zhao.
559 Spatialbot: Precise spatial understanding with vision language models. *arXiv preprint arXiv:2406.13642*,
560 2024.
- 561 Jie Cao and Jing Xiao. An augmented benchmark dataset for geometric question answering through dual parallel
562 text encoding. In *Proceedings of the 29th International Conference on Computational Linguistics (COLING)*,
563 2022.
- 564 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin.
565 Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International*
566 *Conference on Computer Vision (ICCV)*, 2021.
- 567 Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong
568 Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite
569 vision-language model. *arXiv preprint arXiv:2402.11684*, 2024a.
- 570 Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric P Xing, and Liang Lin. Geoqa:
571 A geometric question answering benchmark towards multimodal numerical reasoning. *arXiv preprint*
572 *arXiv:2105.14517*, 2021.
- 573 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok,
574 Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image
575 synthesis. *arXiv preprint arXiv:2310.00426*, 2023a.
- 576 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v:
577 Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023b.
- 578 Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution
579 recovery of diffusion models on low-dimensional data. In *International Conference on Machine Learning*
580 *(ICML)*, 2023c.
- 581 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu,
582 Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models
583 with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024b.
- 584 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou
585 Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic
586 tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
587 2024c.
- 588 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
589 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source
590 chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.

- 594 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li,
595 Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction
596 tuning. *arXiv preprint arXiv:2305.06500*, 2023.
- 597 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
598 image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
599 *(CVPR)*, 2009.
- 600 Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu
601 Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. In *International*
602 *Conference on Learning Representations (ICLR)*, 2024.
- 603 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
604 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words:
605 Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*,
606 2021.
- 607 Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang,
608 Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models.
609 *arXiv preprint arXiv:2407.11691*, 2024.
- 610 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
611 Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint*
612 *arXiv:2407.21783*, 2024.
- 613 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In
614 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- 615 Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02: A visual
616 representation for neon genesis. *Image and Vision Computing*, 2024.
- 617 Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and
618 draw with seed tokenizer. In *International Conference on Learning Representations (ICLR)*, 2024a.
- 619 Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan.
620 Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint*
621 *arXiv:2404.14396*, 2024b.
- 622 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron
623 Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing*
624 *Systems (NeurIPS)*, 2014.
- 625 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter:
626 Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE/CVF*
627 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- 628 Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen,
629 Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language
630 hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference*
631 *on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- 632 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders
633 are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
634 *Recognition (CVPR)*, 2022.
- 635 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
636 Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 637 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free
638 evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- 639 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by
640 a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing*
641 *systems*, 30, 2017.
- 642 Tuomo Hiippala, Malihe Alikhani, Jonas Haverinen, Timo Kalliokoski, Evanfiya Logacheva, Serafina Orekhova,
643 Aino Tuomainen, Matthew Stone, and John A Bateman. Ai2d-rst: A multimodal corpus of 1000 primary
644 school science diagrams. *Language Resources and Evaluation*, 55:661–688, 2021.

- 648 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
649 *Information Processing Systems (NeurIPS)*, 2020.
- 650
- 651 Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and
652 compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
653 *Pattern Recognition (CVPR)*, 2019.
- 654 Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations
655 via question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
656 *Recognition (CVPR)*, 2018.
- 657 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based
658 generative models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 659
- 660 Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in
661 photographs of natural scenes. In *Proceedings of the Conference on Empirical Methods in Natural Language*
662 *Processing (EMNLP)*, 2014.
- 663 Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram
664 is worth a dozen images. In *European Conference on Computer Vision (ECCV)*, 2016.
- 665
- 666 Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi.
667 Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension.
668 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- 669
- 670 Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 671
- 672 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer
673 Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF*
674 *International Conference on Computer Vision (ICCV)*, 2023.
- 675
- 676 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis
677 Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using
678 crowdsourced dense image annotations. *Proceedings of the IEEE/CVF International Conference on Computer*
679 *Vision (ICCV)*, 2017.
- 680
- 681 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural
682 networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- 683
- 684 Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and
685 Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild, May 2024a. URL
686 <https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms/>.
- 687
- 688 Bo Li, Peiyuan Zhang, Kaichen Zhang, Fanyi Pu, Xinrun Du, Yuhao Dong, Haotian Liu, Yuanhan Zhang,
689 Ge Zhang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Accelerating the development of large multimodal
690 models, 2024b. URL <https://github.com/EvolvingLMMS-Lab/lmms-eval>.
- 691
- 692 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu,
693 and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024c.
- 694
- 695 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking
696 multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023a.
- 697
- 698 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training
699 with frozen image encoders and large language models. In *International Conference on Machine Learning*
700 *(ICML)*, 2023b.
- 701
- 702 Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal
703 arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv preprint*
704 *arXiv:2403.00231*, 2024d.
- 705
- 706 Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without
707 vector quantization. *arXiv preprint arXiv:2406.11838*, 2024e.
- 708
- 709 Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and
710 Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint*
711 *arXiv:2403.18814*, 2024f.

- 702 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucina-
703 tion in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023c.
- 704
- 705 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural*
706 *Information Processing Systems (NeurIPS)*, 2023a.
- 707 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In
708 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024a.
- 709
- 710 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next:
711 Improved reasoning, ocr, and world knowledge, January 2024b. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- 712
- 713 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang,
714 Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint*
715 *arXiv:2307.06281*, 2023b.
- 716 Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On
717 the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023c.
- 718
- 719 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet
720 for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
721 *(CVPR)*, 2022.
- 722
- 723 Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu
724 Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint*
725 *arXiv:2403.05525*, 2024.
- 726
- 727 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps:
728 Interpretable geometry problem solving with formal language and symbolic reasoning. *arXiv preprint*
729 *arXiv:2105.04165*, 2021.
- 730
- 731 Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and
732 Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning.
733 *arXiv preprint arXiv:2209.14610*, 2022.
- 734
- 735 Henry B Mann and Donald R Whitney. On a test of whether one of two random variables is stochastically larger
736 than the other. *The annals of mathematical statistics*, pp. 50–60, 1947.
- 737
- 738 Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question
739 answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Conference on*
740 *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 741
- 742 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for
743 question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- 744
- 745 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In
746 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021.
- 747
- 748 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi
749 Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights from multimodal
750 llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- 751
- 752 Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question
753 answering by reading text in images. In *International Conference on Document Analysis and Recognition*
754 *(ICDAR)*, 2019.
- 755
- 756 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
757 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided
758 diffusion models. In *International Conference on Machine Learning (ICML)*, pp. 16784–16804. PMLR, 2022.
- 759
- 760 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *arXiv*
761 *preprint arXiv:1711.00937*, 2017.
- 762
- 763 OpenAI. GPT-4V(ision) System Card, 2023a. URL [https://cdn.openai.com/papers/GPTV_](https://cdn.openai.com/papers/GPTV_System_Card.pdf)
764 [System_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf).
- 765
- 766 R OpenAI. Gpt-4 technical report. *arXiv*, pp. 2303–08774, 2023b.

- 756 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
757 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without
758 supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 759 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by
760 generative pre-training. 2018.
- 762 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are
763 unsupervised multitask learners. *OpenAI blog*, 2019.
- 764 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
765 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language
766 supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- 767 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image
768 synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
769 Pattern Recognition (CVPR)*, 2022.
- 771 Tanik Saikh, Tirthankar Ghosal, Amish Mittal, Asif Ekbal, and Pushpak Bhattacharyya. Scienceqa: A novel
772 resource for question answering on scholarly articles. *International Journal on Digital Libraries*, 23(3):
773 289–301, 2022.
- 774 Jorge Sánchez and Florent Perronnin. High-dimensional signature compression for large-scale image classifi-
775 cation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
776 2011.
- 777 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo
778 Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for
779 training next generation image-text models. *Advances in Neural Information Processing Systems (NeurIPS)*,
780 35:25278–25294, 2022.
- 781 Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A
782 benchmark for visual question answering using world knowledge. In *European Conference on Computer
783 Vision (ECCV)*, 2022.
- 784 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv
785 Batra. Grad-cam: visual explanations from deep networks via gradient-based localization. *International
786 Journal of Computer Vision (IJCV)*, 128:336–359, 2020.
- 787 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus
788 Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer
789 Vision and Pattern Recognition (CVPR)*, 2019.
- 791 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances
792 in Neural Information Processing Systems (NeurIPS)*, 2019.
- 793 Student. The probable error of a mean. *Biometrika*, pp. 1–25, 1908.
- 794 Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive
795 model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024a.
- 797 Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu,
798 Tiejun Huang, and Xinlong Wang. Generative pretraining in multimodality. *arXiv preprint arXiv:2307.05222*,
799 2023.
- 800 Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Yueze Wang, Yongming Rao, Jingjing Liu,
801 Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *Proceedings of
802 the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024b.
- 804 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
805 Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models.
806 *arXiv preprint arXiv:2312.11805*, 2023.
- 807 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.github.
808 io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 809 ShareGPT Team. Sharegpt, 2023. URL <https://sharegpt.com>.

- 810 Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan
811 Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration
812 of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024a.
- 813 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring
814 the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision
815 and Pattern Recognition (CVPR)*, 2024b.
- 816 Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit III: Revenge of the vit. In *European Conference on
817 Computer Vision (ECCV)*, 2022.
- 818 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Anjad Almahairi, Yasmine Babaei, Nikolay Bashlykov,
819 Soumya Batra, Prajjwal Bhargava, Shrubti Bhosale, et al. LLaMA 2: Open foundation and fine-tuned chat
820 models. *arXiv preprint arXiv:2307.09288*, 2023.
- 821 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser,
822 and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*,
823 2017.
- 824 Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing
825 robust features with denoising autoencoders. In *International Conference on Machine Learning (ICML)*,
826 2008.
- 827 Haochen Wang, Junsong Fan, Yuxi Wang, Kaiyou Song, Tiancai Wang, Xiangyu Zhang, and Zhaoxiang Zhang.
828 Bootstrap masked visual modeling via hard patches mining. *arXiv preprint arXiv:2312.13714*, 2023a.
- 829 Haochen Wang, Junsong Fan, Yuxi Wang, Kaiyou Song, Tong Wang, and Zhaoxiang Zhang. Droppos: Pre-
830 training vision transformers by reconstructing dropped positions. *Advances in Neural Information Processing
831 Systems (NeurIPS)*, 2023b.
- 832 Haochen Wang, Kaiyou Song, Junsong Fan, Yuxi Wang, Jin Xie, and Zhaoxiang Zhang. Hard patches mining
833 for masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
834 Recognition (CVPR)*, 2023c.
- 835 Jiacong Wang, Bohong Wu, Haiyong Jiang, Zhou Xun, Xin Xiao, Haoyuan Guo, and Jun Xiao. World to code:
836 Multi-modal data generation via self-instructed compositional captioning and filtering. In *Proceedings of the
837 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4608–4623, 2024a.
- 838 Ke Wang, Juntong Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal
839 mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*, 2024b.
- 840 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang,
841 Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution.
842 *arXiv preprint arXiv:2409.12191*, 2024c.
- 843 xAI. Grok. 2024.
- 844 Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion
845 model for molecular conformation generation. In *International Conference on Learning Representations
846 (ICLR)*, 2022.
- 847 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li,
848 Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024a.
- 849 Ruofeng Yang, Zhijie Wang, Bo Jiang, and Shuai Li. The convergence of variance exploding diffusion models
850 under the manifold hypothesis, 2024b. URL <https://openreview.net/forum?id=tD4NOxYTFg>.
- 851 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang,
852 Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning
853 benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
854 Recognition (CVPR)*, 2024.
- 855 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really
856 finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- 857 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image
858 pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.

864 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness
865 of deep features as a perceptual metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
866 *Pattern Recognition (CVPR)*, 2018.

867 Yi-Fan Zhang, Qingsong Wen, Chaoyou Fu, Xue Wang, Zhang Zhang, Liang Wang, and Rong Jin. Beyond
868 llava-hd: Diving into high-resolution large multimodal models. *arXiv preprint arXiv:2406.08487*, 2024.

869 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou.
870 Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

871 Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu,
872 Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale corpus of
873 images interleaved with text. *Advances in Neural Information Processing Systems*, 36, 2023.

874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

APPENDIX

A LATENT DIFFUSION MODELS

Given a set of clean latent tokens \mathbf{z}_0 drawn from $p(\mathbf{z})$, the forward diffusion process is a Markov chain that gradually adds random Gaussian noise to the original sample:

$$q(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}\mathbf{z}_{t-1}, \beta_t\mathbf{I}), \quad (5)$$

where $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the Gaussian distribution, and t indicates discrete timesteps. $\beta_t \in (0, 1)$ indicates a pre-defined time-dependent variance schedule. According to Ho et al. (2020), to admit sampling \mathbf{z}_t at an arbitrary timestep t directly from \mathbf{z}_0 , this transition can be reformulated as

$$\begin{aligned} q(\mathbf{z}_t|\mathbf{z}_0) &= \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{z}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \\ \mathbf{z}_t &= \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \end{aligned} \quad (6)$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. A latent diffusion model learns to *reverse* this progressive noise addition process for latent tokens. Specifically, to iteratively generate clean tokens \mathbf{z}_0 from pure noise \mathbf{z}_T conditioned on \mathcal{C} , we need to reverse the forward process by

$$\mathbf{z}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{z}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\pi(\mathbf{z}_t; \mathcal{C}, t) \right) + \sigma_t \boldsymbol{\epsilon}, \quad (7)$$

where a π -parameterized neural network $\boldsymbol{\epsilon}_\pi$ is trained to predict the added noise during the forward process. σ_t indicates the posterior noise variance. The training objective of $\boldsymbol{\epsilon}_\pi$ is

$$\mathcal{L}(\pi, \mathbf{z}_0) = \mathbb{E}_{t, \boldsymbol{\epsilon}} [\|\boldsymbol{\epsilon}_\pi(\sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}; \mathcal{C}, t) - \boldsymbol{\epsilon}\|^2]. \quad (8)$$

B IMPLEMENTATION DETAILS

Table 6: **Hyperparameters of ROSS.** We obtain most of the configurations from Liu et al. (2024a).

Config	Stage I	Stage II
Trainable parts	projector + denoiser	projector + LLM + denoiser
Frozen parts	visual encoder + LLM + teacher tokenizer	visual encoder + teacher tokenizer
Global batch size	256	128
Batch size per GPU	16	4
Accumulated steps	2	4
DeepSpeed zero stage	2	3
Learning rate	1×10^{-3}	2×10^{-5}
Learning rate schedule	warmup + cosine decay	
Warmup ratio	0.03	
Weight decay	0	
Epoch	1	
Optimizer	AdamW	
Precision	bf16	

Table 7: **Details of the instruction tuning dataset** provided by Tong et al. (2024a).

Dataset	# Samples
LLaVA (Liu et al., 2023a)	158K
ShareGPT (Team, 2023)	40K
VQAv2 (Goyal et al., 2017)	83K
GQA (Hudson & Manning, 2019)	72.1K
OKVQA (Marino et al., 2019)	9K
OCRvQA (Mishra et al., 2019)	80K
A-OKVQA (Schwenk et al., 2022)	50K
TextVQA (Singh et al., 2019)	21.9K
RefCOCO (Kazemzadeh et al., 2014)	30K
VG (Krishna et al., 2017)	86.4K
DVQA (Kafle et al., 2018)	13K
DocVQA (Mathew et al., 2021)	15K
ChartQA (Masry et al., 2022)	28.1K
AI2 Diagrams (Kembhavi et al., 2016)	15.5K

Table 8: **Details of the instruction tuning dataset** provided by Zhang et al. (2024).

Dataset	# Samples
SciencQA (Saikh et al., 2022)	9K
TextbookQA (Kembhavi et al., 2017)	9.5K
AI2 Diagrams (Kembhavi et al., 2016)	12.4K
ChartQA (Masry et al., 2022)	28.3K
DVQA (Kafle et al., 2018)	200K
ArxivQA (Li et al., 2024d)	100K
GeoQA3 (Chen et al., 2021)	5K
Geometry3K (Lu et al., 2021)	2.1K
GeoQA+ (Cao & Xiao, 2022)	72.3K
MathVision (Wang et al., 2024b)	2.7K
TabMWP (Lu et al., 2022)	30.7K

Table 9: **Summary of the evaluation benchmarks.** Prompts are mostly borrowed from VLMEvalKit (Duan et al., 2024) and lmms-eval (Li et al., 2024b).

Benchmark	Response formatting prompts
POPE (Li et al., 2023c)	–
HallusionBench (Guan et al., 2024)	Answer the question using a single word or phrase.
MMBench (Liu et al., 2023b)	Answer with the option’s letter from the given choices directly.
SEED-Bench (Li et al., 2023a)	Answer with the option’s letter from the given choices directly.
MMMU (Yue et al., 2024)	Answer with the option’s letter from the given choices directly.
MMVP (Tong et al., 2024b)	Answer with the option’s letter from the given choices directly.
AI2D (Hiippala et al., 2021)	Answer with the option’s letter from the given choices directly.
RealWorldQA (xAI, 2024)	Answer with the option’s letter from the given choices directly.
GQA (Hudson & Manning, 2019)	Answer the question using a single word or phrase.
ChartQA (Masry et al., 2022)	Answer the question using a single word or phrase.
OCRBench (Liu et al., 2023c)	Answer the question using a single word or phrase.
DocVQA (Mathew et al., 2021)	Answer the question using a single word or phrase.
InfoVQA (Biten et al., 2022)	Answer the question using a single word or phrase.
TextVQA (Singh et al., 2019)	Answer the question using a single word or phrase.

Table 10: **Comparisons on computational costs during the instruction tuning stage** with Cambrian-737K (Tong et al., 2024a), where evaluations are conducted using 8 A100 GPUs with a global batch size of 128. Due to the limited GPU memory, we accumulate 4 gradient steps and the batch size per GPU is 4. The whole stage requires 5757 training steps. GPU memories are averaged over 8 GPUs with DeepSpeed Zero 3.

Vision	Base LLM	$\mathcal{L}_{LMM}^{visual}$	Trainable Parameters	Speed (s/iter)	Time	GPU Memory
CLIP-L/336	Qwen2-7B-Instruct	–	7.63 B	8.31	13h 17min	45.34 G
CLIP-L/336	Qwen2-7B-Instruct	✓	7.68 B	9.02 (1.09×)	14h 25min	46.62 G (1.03×)
CLIP-L/336	Vicuna-13B-v1.5	–	13.05 B	13.33	21h 19min	48.62 G
CLIP-L/336	Vicuna-13B-v1.5	✓	13.11 B	14.69 (1.10×)	23h 30min	49.07 G (1.01×)
SigLIP-L/384	Qwen2-7B-Instruct	–	7.63 B	8.77	14h 1min	47.08 G
SigLIP-L/384	Qwen2-7B-Instruct	✓	7.68 B	9.48 (1.08×)	15h 9min	52.07 G (1.11×)
SigLIP-L/384	Vicuna-13B-v1.5	–	13.05 B	14.22	22h 44min	48.80 G
SigLIP-L/384	Vicuna-13B-v1.5	✓	13.11 B	15.32 (1.08×)	24h 30min	52.68 G (1.08×)

Hyperparameters. The hyperparameters of **ROSS** are provided in Table 6. We simply borrow most configurations from LLaVA-v1.5 (Liu et al., 2024a) without further tuning, as we find it works well with our **ROSS**, even if we adopt SigLIP (Zhai et al., 2023) and Qwen2 (Yang et al., 2024a) while the original LLaVA-v1.5 (Liu et al., 2024a) utilized CLIP (Radford et al., 2021) and Vicuna-v1.5 (Chiang et al., 2023). As SigLIP represents a single 384×384 image with 729 tokens and the downsampling ratio of the teacher tokenizer KL-16 (Rombach et al., 2022) is 16, we set the input resolution of the teacher tokenizer as $432 = \sqrt{729} \times 16$ to produce 729 fine-grained tokens as denoising targets.

Instruction Tuning Data. When comparing with state-of-art LMMs in Table 4, our **ROSS** is trained on approximately 1.2M instruction tuning data, which is a mixture of Cambrian-737K (Tong et al., 2024a) and SMR-473K (Zhang et al., 2024). Details of these two instruction tuning datasets are listed in Table 7 and Table 8, respectively. There might be some overlap but we simply concat these two datasets as it is already empirically effective.

Evaluation Prompts. We provide a thorough examination of all evaluation benchmarks utilized in this paper in Table 9. Notably, for MMVP (Tong et al., 2024b), which is not officially supported by VLMEvalKit (Duan et al., 2024), we follow Cambrian-1 (Tong et al., 2024a) to reformat the original question into a multiple-choice format and compute the accuracy using exact matching.

Computational Costs. As demonstrated in Table 10, the denoising process introduces a negligible increase in training time ($\approx 10\%$ compared to the baseline), while the benefits outweigh the minor additional costs.

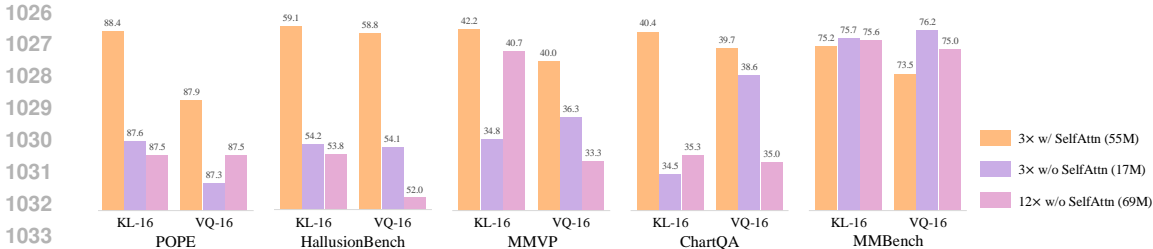


Figure 10: **The architecture of the denoiser and the choice of fine-grained tokenizer.** The self-attention module illustrated in Figure 4b is crucial since orange bars consistently outperform others on hallucination and fine-grained comprehension benchmarks, while maintaining similar performances on the general understanding benchmark. KL-16 provided by Rombach et al. (2022) is better than VQ-16 provided by Sun et al. (2024a), as quantization may lead to information loss.

Table 11: **Ablations on different schedules of β .** ROSS consistently improves the baseline, demonstrating its robustness to the denoising schedule.

Schedule of β	POPE	Hallu.	MMVP	ChartQA	MMB ^{EN}
–	87.9	55.0	29.6	34.0	73.8
Linear (Ho et al., 2020)	88.1 \uparrow 0.2	57.3 \uparrow 2.3	42.0 \uparrow 12.4	39.2 \uparrow 5.2	75.1 \uparrow 1.3
Scaled Linear (Rombach et al., 2022)	88.4 \uparrow 0.5	58.3 \uparrow 3.3	40.0 \uparrow 10.4	40.7 \uparrow 6.7	75.3 \uparrow 1.5
GLIDE Softmax (Nichol et al., 2022)	88.4 \uparrow 0.5	59.1 \uparrow 4.1	42.2 \uparrow 12.6	40.4 \uparrow 6.4	75.2 \uparrow 1.4
GeoDiff Sigmoid (Xu et al., 2022)	88.2 \uparrow 0.3	57.7 \uparrow 2.7	41.3 \uparrow 11.7	38.9 \uparrow 4.9	75.5 \uparrow 1.7

C MORE EXPERIMENTS

C.1 MORE ABLATIONS

KL-16 v.s. VQ-16. Our default tokenizer is a continuous VAE (Kingma, 2013) with Kullback-Leibler (KL) divergence trained by Rombach et al. (2022). We further conduct experiments with a *discrete* tokenizer provided by Sun et al. (2024a), which is a VQGAN (Esser et al., 2021), i.e., VQVAE (Oord et al., 2017) with additional perceptual loss (Zhang et al., 2018) and adversarial loss (Goodfellow et al., 2014). *KL-16 outperforms VQ-16.* One intuitive explanation is that KL-16 preserves more low-level details than VQ-16 since quantization may lead to information loss. Moreover, quantitatively, on ImageNet (Deng et al., 2009) 256 \times 256 validation set, KL-16 achieves 0.87 rFID (Heusel et al., 2017) while the rFID (Heusel et al., 2017) of VQ-16 is 2.19.

Architecture of the Denoiser. Illustrated in Figure 10, *the self-attention module is crucial*, as original visual outputs $x_{i \leq N}$ are actually *causal* and we need to model inter-token discrepancy via self-attention. The number of trainable parameters is *not* the crux.

Schedule of β . We study the effectiveness of different schedules of β in Table 11. From the table, we can tell that even with different schedules of β , **ROSS consistently improves the baseline**, demonstrating its robustness to the denoising schedule.

Generative v.s. Reconstructive. We offer a detailed pipeline comparison in Figure 11. Experimental results have already been provided in Table 2. The implementation of generative methods is similar to Sun et al. (2024b) and Dong et al. (2024), where we adopt 576 learnable queries as inputs and take the corresponding outputs as conditions for the denoiser.

We hypothesize that the underlying reason for the lower performance of generative methods in comprehension tasks is *the weak correspondence between inputs and supervision* under generative settings, which typically arises from both the (1) data and the (2) design of these methods.

(1) Typical generative methods that explore the synergy of comprehension and generation, usually leverage image generation conditioned on text instructions on (i) *text-to-image datasets* or (ii) *interleaved datasets* as extra supervision. However, (i) text-to-image datasets are typically designed to generate *high-aesthetic* samples rather than text-aligned ones, and (ii) interleaved datasets aim to enable few-shot learning via interleaving independent supervised examples, where reasoning becomes more important than alignment. Therefore, there exists a clear disconnect where the supervision

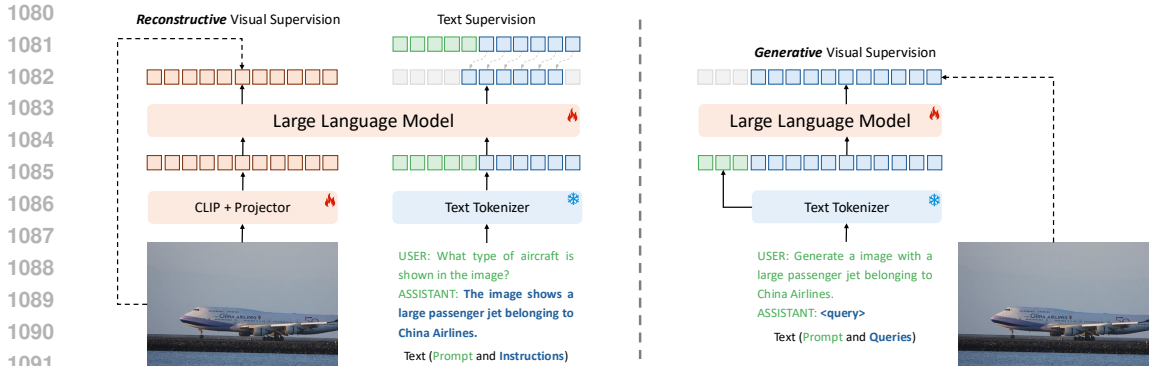


Figure 11: **Pipeline comparison between reconstructive and generative.** The reconstructive objective (left) does *not* require specific data formulations and can be easily combined with current visual instruction tuning data. However, the generative objective (right) needs specific *text-to-image* creation data, which could be converted by image-to-text caption data.

Table 12: **Extended ablations** on The effectiveness of the vision-centric supervision $\mathcal{L}_{\text{LLM}}^{\text{visual}}$ among various LLMs and visual encoders. Pre-training data is LLaVA-558K (Liu et al., 2023a) and instruction tuning data is Cambrian-737K (Tong et al., 2024a). Evaluations of POPE (Li et al., 2023c), HallusionBench (Guan et al., 2024), MMBench (Liu et al., 2023b), SEED-Bench-1 (Li et al., 2023a), MMMU (Yue et al., 2024), MMVP (Tong et al., 2024b), AI2D (Hiippala et al., 2021), OCRBench (Liu et al., 2023c), and RealWorldQA (xAI, 2024) are conducted with VLMEvalKit (Duan et al., 2024), while evaluations of ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), InfoVQA (Biten et al., 2022), and TextVQA (Singh et al., 2019) are conducted with Imms-eval (Li et al., 2024b).

Benchmark	CLIP-ViT-L/14@336				SigLIP-ViT-SO400M/14@384			
	Vicuna-7B-v1.5		Qwen2-7B-Instruct		Vicuna-7B-v1.5		Qwen2-7B-Instruct	
	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS
POPE _{acc}	86.3	87.2 ↑ 0.9	87.9	88.4 ↑ 0.5	86.0	87.7 ↑ 1.7	88.5	88.7 ↑ 0.2
HallusionBench _{aAcc}	52.5	55.8 ↑ 3.3	55.0	59.1 ↑ 4.1	50.4	53.8 ↑ 3.4	57.3	58.2 ↑ 0.9
MMBench-EN _{dev}	67.0	67.6 ↑ 0.6	73.8	75.2 ↑ 1.4	64.5	69.2 ↑ 4.7	76.3	76.9 ↑ 0.6
MMBench-CN _{dev}	60.0	59.8 ↓ 0.2	72.9	73.7 ↑ 0.8	63.1	63.4 ↑ 0.3	75.7	76.3 ↑ 0.7
SEED _{img}	66.7	66.4 ↓ 0.3	70.3	70.7 ↑ 0.4	68.2	69.0 ↑ 0.8	72.3	72.1 ↓ 0.2
MMMU _{dev}	30.0	34.0 ↑ 4.0	44.0	45.3 ↑ 1.3	33.3	38.0 ↑ 4.7	38.7	41.3 ↑ 2.6
MMMU _{val}	35.3	36.0 ↑ 0.7	41.9	42.6 ↑ 0.7	34.2	35.4 ↑ 1.2	41.8	43.8 ↑ 2.0
MMVP	28.0	36.3 ↑ 8.3	29.6	42.2 ↑ 12.6	27.3	38.0 ↑ 10.7	40.7	49.3 ↑ 8.6
AI2D _{test}	61.2	61.4 ↑ 0.2	71.9	73.3 ↑ 1.4	62.6	62.4 ↓ 0.2	74.0	74.5 ↑ 0.5
ChartQA _{test}	32.9	39.8 ↑ 6.9	36.2	41.6 ↑ 5.4	34.0	48.2 ↑ 14.2	44.4	46.9 ↑ 2.5
DocVQA _{val}	33.4	41.6 ↑ 8.2	31.1	44.7 ↑ 13.6	40.4	40.7 ↑ 0.3	39.2	39.3 ↑ 0.1
InfoVQA _{val}	21.2	26.4 ↑ 5.2	22.1	39.3 ↑ 16.2	22.8	23.3 ↑ 0.5	24.0	25.1 ↑ 1.1
TextVQA _{val}	55.7	58.7 ↑ 3.0	52.0	54.1 ↑ 2.1	60.5	62.6 ↑ 2.1	56.3	57.5 ↑ 1.2
OCRBench	339	350 ↑ 11	363	381 ↑ 18	354	365 ↑ 11	432	448 ↑ 16
RealWorldQA	52.7	53.2 ↑ 0.5	56.7	57.4 ↑ 0.7	55.0	57.1 ↑ 2.1	57.9	59.1 ↑ 1.2
Average	47.8	50.6 ↑ 2.8	52.1	56.4 ↑ 4.3	49.2	52.4 ↑ 3.2	55.4	56.9 ↑ 1.5

(image) has little to do with the input (text instruction). For example, the CLIP-Score (Hessel et al., 2021), which measures the similarity between text and images, is only 0.3043 for the LAION-Art dataset (Schuhmann et al., 2022) and 0.2842 for the MMC4 dataset (Zhu et al., 2023), indicating that the supervision signals in these datasets are *not* well-suited for tasks requiring strong text-image alignment.

(2) Even when we attempt to ensure image-text alignment by converting aligned caption data into creation data for supervision, the results demonstrated in Table 2 remain unsatisfactory. This suggests that the *design of generative objectives itself does not inherently require a strong correspondence* between inputs and supervision targets.

Table 13: **Comparison to state-of-the-art LLMs on benchmarks requires high-resolution inputs.** We evaluate models on: ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021) val set, InfoVQA (Biten et al., 2022) val set, TextVQA (Singh et al., 2019) val set, OCRBench (Liu et al., 2023c), and RealWorldQA (xAI, 2024). †We evaluate the official checkpoint.

Model	ChartQA	DocVQA	InfoVQA	TextVQA	OCRBench	RealWorldQA
GPT-4V-1106 (OpenAI, 2023a)	78.5	88.4	–	78.0	645	61.4
Gemini-1.5 Pro (Team et al., 2023)	81.3	86.5	–	78.1	–	67.5
Grok-1.5 (xAI, 2024)	76.1	85.6	–	78.1	–	68.7
LLaVA-v1.5-7B [†] (Liu et al., 2024a)	18.2	28.1	25.7	58.2	317	54.9
LLaVA-v1.6-7B [†] (Liu et al., 2024b)	65.5	74.4	37.1	64.8	532	57.6
Cambrian-1-8B (Tong et al., 2024a)	73.3	77.8	–	71.7	624	64.2
ROSS-7B_{anyres}	76.9	81.8	50.5	72.2	607	66.2

Table 14: **Evaluations on language performance.** We evaluate multi-modal benchmarks that mainly require general knowledge following Tong et al. (2024a). Furthermore, we incorporate representative language benchmarks, including general understanding on MMLU (Hendrycks et al., 2020) and HellaSwag (Zellers et al., 2019), and instruction-following on IFEval (Zhou et al., 2023). ROSS does *not* harm language capabilities as it brings improvements in most cases.

Benchmark	CLIP-ViT-L/14@336				SigLIP-ViT-SO400M/14@384			
	Vicuna-7B-v1.5		Qwen2-7B-Instruct		Vicuna-7B-v1.5		Qwen2-7B-Instruct	
	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS
<i>Vision-Language Benchmarks on Knowledge</i>								
ScienceQA _{test}	68.5	69.0 ↑ 0.5	76.5	77.4 ↑ 0.9	69.6	71.3 ↑ 1.7	78.3	78.5 ↑ 0.2
MMMU _{dev}	30.0	34.0 ↑ 4.0	44.0	45.3 ↑ 1.3	33.3	38.0 ↑ 4.7	38.7	41.3 ↑ 2.6
MMMU _{val}	35.3	36.0 ↑ 0.7	41.9	42.6 ↑ 0.7	34.2	35.4 ↑ 1.2	41.8	43.8 ↑ 2.0
AI2D _{test}	61.2	61.4 ↑ 0.2	71.9	73.3 ↑ 1.4	62.6	62.4 ↓ 0.2	74.0	74.5 ↑ 0.5
<i>Language Benchmarks</i>								
MMLU	26.5	27.4 ↑ 0.9	57.1	60.7 ↑ 3.6	26.0	25.9 ↓ 0.1	60.9	61.0 ↑ 0.1
HellaSwag _{acc-norm}	27.0	26.9 ↓ 0.1	46.4	46.2 ↓ 0.2	27.1	27.0 ↓ 0.1	45.5	46.6 ↑ 1.1
IFEval _{strict-inst}	41.2	44.6 ↑ 3.4	47.1	49.2 ↑ 2.1	43.6	43.8 ↑ 0.2	47.8	48.1 ↑ 0.3
IFEval _{strict-prompt}	28.7	35.3 ↑ 6.7	35.1	37.0 ↑ 1.9	32.5	33.1 ↑ 0.6	35.3	36.2 ↑ 0.9
Average	39.8	41.8 ↑ 2.0	52.5	54.0 ↑ 1.5	41.1	42.1 ↑ 1.0	52.8	53.8 ↑ 1.0

In contrast, reconstructive methods like Ross leverage the original input images as auxiliary supervision, ensuring a strong and direct correspondence, which is crucial for tasks requiring accurate comprehension and interpretation of multimodal data, leading to significantly improved performance.

Extended Ablations on Different LLMs and Visual Encoders. We extend the ablation in Table 3 by incorporating more benchmarks, providing a more balanced and representative distribution of tasks. Empirical results in Table 12 demonstrate that our proposed vision-centric supervision utilized by ROSS leads to significant improvements in most cases. Moreover, we found ROSS contributes more significant improvements over fine-grained comprehension datasets, such as HallusionBench (Guan et al., 2024), MMVP (Tong et al., 2024b), ChartQA (Masry et al., 2022), and OCRBench (Liu et al., 2023c).

C.2 COMPARISON ON HIGH-RESOLUTION BENCHMARKS

We incorporate the “anyres” technique proposed by LLaVA-v1.6 (Liu et al., 2024b) into our ROSS. Specifically, for each image, we employ a grid configuration of $384 \times \{2 \times 2, 1 \times \{2, 3, 4\}, \{2, 3, 4\} \times 1\}$ to identify the input resolution, resulting in a maximum of $5 \times 729 = 3,645$ visual tokens. Each 384×384 crop is required to reconstruct the original input via the denoising objective proposed by ROSS. In Table 13, our ROSS-7B_{anyres} surpasses LLaVA-v1.6-7B (Liu et al., 2024b) and Cambrian-1-8B (Tong et al., 2024a) under most cases. These results indicate that ROSS not only performs well at lower resolutions but also maintains its competitive edge at higher resolutions, making it a robust and versatile method.

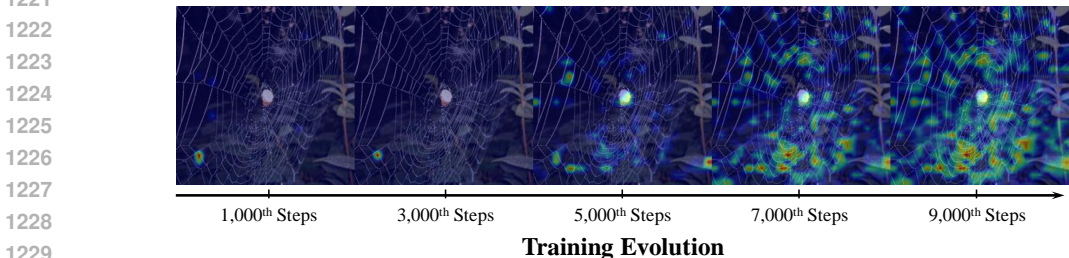
1188 Table 15: **Model scaling of ROSS.** We take Qwen2.5 series (Team, 2024) as the base language
 1189 model and CLIP-ViT-L/14@336 (Radford et al., 2021) as the visual encoder. Pre-training data is
 1190 LLaVA-558K (Liu et al., 2023a) and the instruction tuning data is LLaVA-665K (Liu et al., 2024a).
 1191 **ROSS brings improvements over the baseline across different model sizes in most cases.**

Benchmark	0.5B		1.5B		3B		7B	
	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS	LLaVA	ROSS
POPE _{acc}	50.0	60.4 ↑ 10.4	85.3	87.9 ↑ 2.4	87.3	88.1 ↑ 0.8	87.9	88.4 ↑ 0.5
HallusionBench _{acc}	45.8	48.0 ↑ 2.2	48.7	49.6 ↑ 0.9	52.2	52.2 - 0.0	48.7	53.7 ↑ 5.0
MMBench-EN _{dev}	55.2	60.4 ↑ 5.2	67.5	68.2 ↑ 1.7	70.6	71.4 ↑ 0.8	75.0	75.7 ↑ 0.7
MMBench-CN _{dev}	45.6	48.9 ↑ 3.3	62.4	63.9 ↑ 1.5	68.0	69.1 ↑ 1.1	73.6	73.5 ↓ 0.1
SEED _{img}	55.8	55.6 ↓ 0.2	66.3	66.8 ↑ 0.5	68.2	68.4 ↑ 0.2	70.6	71.0 ↑ 0.4
OCRBench	229	248 ↑ 19	291	298 ↑ 7	313	308 ↓ 5	334	358 ↑ 24
MMMU _{dev}	35.2	36.0 ↑ 0.8	44.7	45.0 ↑ 0.3	48.7	49.0 ↑ 0.3	48.0	48.0 - 0.0
MMMU _{val}	38.0	40.3 ↑ 1.7	41.8	43.6 ↑ 1.8	41.6	42.7 ↑ 1.1	47.3	48.0 ↑ 0.7
AI2D _{test}	45.3	46.0 ↑ 0.7	59.0	59.5 ↑ 0.5	62.9	63.2 ↑ 0.3	68.3	68.5 ↑ 0.2
RealWorldQA	45.1	46.4 ↑ 1.3	50.5	53.5 ↑ 3.0	55.7	57.9 ↑ 2.2	59.5	59.9 ↑ 0.4
Average	43.9	46.7 ↑ 2.8	55.3	56.8 ↑ 1.5	58.9	59.3 ↑ 0.4	61.2	62.3 ↑ 1.1

1205 Table 16: **Data scaling of ROSS.** We take Qwen2-7B-Instruct (Yang et al., 2024a) as the base
 1206 language model and CLIP-ViT-L/14@336 (Radford et al., 2021) as the visual encoder. **ROSS**
 1207 **consistently brings significant improvements as the training data scale increases.**

PT	SFT	$\mathcal{L}_{LMM}^{visual}$	POPE	Hallu.	ChartQA	OCRBench	MMB ^{EN}	AI2D
558K	737K	-	87.9	55.0	34.0	363	73.8	72.4
		✓	88.4 ↑ 0.5	59.1 ↑ 4.1	40.4 ↑ 6.4	380 ↑ 17	75.2 ↑ 1.4	73.3 ↑ 0.9
558K	1.2M	-	88.5	57.3	37.0	389	75.7	74.5
		✓	88.8 ↑ 0.3	57.8 ↑ 0.5	42.0 ↑ 5.0	392 ↑ 3	76.8 ↑ 1.1	74.7 ↑ 0.2
2M	737K	-	88.1	55.6	37.3	384	76.2	72.3
		✓	88.3 ↑ 0.2	56.2 ↑ 0.6	41.9 ↑ 4.5	398 ↑ 14	77.0 ↑ 0.8	73.4 ↑ 1.1
2M	1.2M	-	88.5	53.8	41.2	388	76.5	73.9
		✓	88.9 ↑ 0.4	57.3 ↑ 2.5	43.2 ↑ 2.0	405 ↑ 17	78.0 ↑ 1.5	74.1 ↑ 0.2

1219 From the camera’s perspective, is the spider web very dense or relatively sparse?
 1220 (a) Very dense (b) Relatively sparse



1230 Figure 12: **Qualitative comparison using GradCAM (Selvaraju et al., 2020) on MMVP (Tong**
 1231 **et al., 2024b).** We visualize the gradient of the second-to-last block of the LMM using the option of
 1232 the ground-truth answer as the target class. Equipped with our proposed vision-centric supervision
 1233 signals, **ROSS** provides more reasonable gradients and urges LMMs to focus on relevant regions (the
 1234 spider web) as the training evolves.

1235 **C.3 MORE ANALYSIS**

1237 **Language Capabilities.** One possible concern of **ROSS** is that this type of low-level reconstruction
 1238 may harm the high-level language capabilities. To investigate this issue, we evaluate multi-modal
 1239 benchmarks that mainly require general knowledge following (Tong et al., 2024a), including ScienceQA
 1240 (Saikh et al., 2022), MMMU (Yue et al., 2024), and AI2D (Hiippala et al., 2021). Fur-
 1241 thermore, we incorporate representative language benchmarks, including general understanding on
 MMLU (Hendrycks et al., 2020) and HellaSwag (Zellers et al., 2019), and instruction-following

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Is the right orange circle smaller than the left orange circle?

LLaVA-v1.5-7B	Yes, the right orange circle is smaller than the left orange circle.	LLaVA-v1.5-7B	No , New York is not the only U.S. state that shares a border with Lake Huron. The image shows that the Great Lakes region includes the states of Michigan, Ohio, Indiana, Illinois, Wisconsin, and Minnesota, all of which share borders with Lake Huron.
Cambrian-1-8B	Yes, the right orange circle is smaller than the left orange circle.	Cambrian-1-8B	Based on the image, New York is not the only U.S. state that shares a border with Lake Huron.
GPT-4V	The orange circles in both the left and right images are actually the same size . This is a visual illusion where the surrounding circles influence the perceived size of the central orange circle. In the left image, the smaller blue circles make the orange circle appear larger, whereas in the right image, the larger blue circles make the orange circle appear smaller.	GPT-4V	No , the image does not show New York as bordering Lake Huron. It suggests that New York borders lakes Ontario and Erie, but Lake Huron does not share a border with New York.
ROSS-7B	No, the right orange circle is larger than the left orange circle.	ROSS-7B	Yes, New York is the only U.S. state that shares a border with Lake Huron.

Figure 13: Qualitative comparisons on HallusionBench (Guan et al., 2024).

What is the name of the place shown?

Which image shows the highest contrast?

What direction is Italy in the Mediterranean Sea?

LLaVA-v1.5-7B	A	LLaVA-v1.5-7B	B	LLaVA-v1.5-7B	C
Cambrian-1-8B	A	Cambrian-1-8B	D	Cambrian-1-8B	B
GPT-4V	C. New York	GPT-4V	D. Down right	GPT-4V	C. West
ROSS-7B	D	ROSS-7B	A	ROSS-7B	D

Figure 14: Qualitative comparisons on MMbench (Guan et al., 2024) English dev split.

on IFEval (Zhou et al., 2023). Empirical results in Table 14 demonstrate that ROSS does *not* harm language capabilities as it brings improvements in most cases.

Model Scaling Properties. To study the stability and scalability of ROSS across different model sizes, we use the Qwen2.5 series (Team, 2024) with varying sizes as the base language model while keeping the CLIP-ViT-L/14@336 (Radford et al., 2021) as the visual encoder. The pre-training data is LLaVA-558K (Liu et al., 2023a), and the instruction tuning data is LLaVA-665K (Liu et al., 2024a). The results, shown in Table 15, demonstrate that ROSS *brings improvements over the baseline (LLaVA) across different model sizes in most cases.*

Data Scaling Properties. To study the impact of the training data scale, we used Qwen2-7B-Instruct (Yang et al., 2024a) as the base language model and CLIP-ViT-L/14@336 (Radford et al., 2021) as the visual encoder. We compared the performance of ROSS and the baseline under different scales of training data. Table 16 demonstrates that ROSS *consistently brings significant improvements as the training data scale increases.*

Gradient Analysis. To better explain the reasoning behind how the vision-centric supervision enables the model to focus on relevant areas of the image during VQA tasks, we provide qualitative comparison using GradCAM (Selvaraju et al., 2020) on MMVP (Tong et al., 2024b) in Figure 12, since GradCAM helps in understanding which parts of the image the model is focusing on, making the model’s decision-making process more transparent. In our analysis, we visualize the gradient of the second-to-last block of the LMM, regarding the option of the ground-truth answer as the target class. Specifically in this case, where the providing question is about the spider web, our proposed



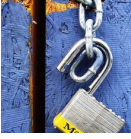



1296		From which angle is this image taken?		Is the butterfly's abdomen visible in the image?		Is the lock locked or unlocked?
1297		A. Front		A. Yes		A. Locked
1298		B. Side		B. No		B. Unlocked
1299						
1300						
1301	LLaVA-v1.5-7B	B	LLaVA-v1.5-7B	A	LLaVA-v1.5-7B	A
1302	Cambrian-1-8B	B	Cambrian-1-8B	A	Cambrian-1-8B	A
1303	ROSS-7B	A	ROSS-7B	B	ROSS-7B	B
1304		Are the ears of the dog erect or drooping?		In this image, how many eyes can you see on the animal?		Can you see the side windows of the vehicles?
1305		A. Erect		A. One		A. Yes
1306		B. Drooping		B. Two		B. No
1307						
1308						
1309	LLaVA-v1.5-7B	A	LLaVA-v1.5-7B	B	LLaVA-v1.5-7B	A
1310	Cambrian-1-8B	A	Cambrian-1-8B	B	Cambrian-1-8B	A
1311	ROSS-7B	B	ROSS-7B	A	ROSS-7B	B

Figure 15: Qualitative comparisons on MMVP (Tong et al., 2024b).




1315		In real world, which dog is smaller in size?		What is the positional relationship between the group of people with the flag and the black car?		Has the man touched the elephant?
1316		A. The dog closer to the camera.		A. Behind the black car.		LLaVA
1317		B. The dog further to the camera.		B. Left of the black car.		LLaVA (w/ MiDaS)
1318		C. They seem to be equally large.		C. Right of the black car.		LLaVA (w/ MiDaS)
1319		D. It can not be decided from the image because information given is not enough.		D. In front of the black car.		ROSS
1320						Yes.
1321						Yes.
1322						No.
1323	LLaVA	A	LLaVA	D	LLaVA	Yes.
1324	LLaVA (w/ MiDaS)	A	LLaVA (w/ MiDaS)	A	LLaVA (w/ MiDaS)	Yes.
1325	ROSS	B	ROSS	D	ROSS	No.

Figure 16: Qualitative comparisons on SpatialBench (Cai et al., 2024). We take RGB + D inputs when testing. Notably, the extra depth expert MiDaS-3.0 (Birkl et al., 2023) sometimes harms comprehension (see the second example).

vision-centric supervision signals provide more reasonable gradients and urge LMMs to focus on relevant regions, *i.e.*, the spider web, as the training evolves.

C.4 QUALITATIVE COMPARISONS

We provide sufficient qualitative comparisons in Figure 13, Figure 14, Figure 15, and Figure 16 on HallusionBench (Guan et al., 2024), MMBench (Liu et al., 2023b) English dev split, MMVP (Tong et al., 2024b), and SpatialBench (Cai et al., 2024), respectively. In Figure 13, Figure 14, and Figure 15, we compare our **ROSS-7B** with the instruction tuning baseline LLaVA-v1.5-7B (Liu et al., 2024a), the state-of-the-art open-source method using extrinsic assistance Cambrian-1-8B (Tong et al., 2024a), and GPT-4V (OpenAI, 2023a).

As demonstrated in Figure 13, where we highlight the *wrong* parts of each prediction in red, our **ROSS** manages to correctly answer the question with reduced hallucinations even when GPT-4V fails. Cambrian-1 (Tong et al., 2024a) even fails to follow the instructions in the second example. This could be because a super huge SFT data (7M) may harm the instruction-following abilities of LMMs. Qualitative results shown in Figure 14 demonstrate both enhanced reasoning abilities (the first example), low-level comprehension capabilities (the second example), and spatial understanding skills (the third example). Figure 15 illustrates that our **ROSS** is good at recognizing various visual patterns, implying that the introduced reconstructive vision-centric objective indeed makes up the visual shortcomings of the original visual encoder.

1350 Figure 16 provides qualitative results on SpatialBench (Cai et al., 2024). The extra depth understand-
1351 ing visual expert, *i.e.*, MiDaS (Birkl et al., 2023), fails to help LMMs understand depth maps both
1352 quantitatively in Table 5 and qualitatively in Figure 16.
1353

1354 D DISCUSSION

1355 One limitation is that **ROSS** does not have generation capabilities, since **ROSS** is designed for
1356 enhanced multimodal comprehension, without the need to generate photorealistic aesthetic images.
1357 Furthermore, the gap in *training data* between comprehension and generation methods also matters.
1358 For instance, PixArt- α (Chen et al., 2023a), which is one of the most *efficient* text-to-image models,
1359 was trained on nearly 400M images to model the pixel discrepancy just in the *first* training stage.
1360 By contrast, our **ROSS** is only trained on nearly 3M images for one epoch. Future topics include
1361 achieving photorealistic text-to-image generation via incorporating more training samples.
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403