

GRANViT: A FINE-GRAINED VISION MODEL FOR AUTOREGRESSIVE MULTIMODAL LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Vision encoders are indispensable for allowing impressive performance of Multimodal Large Language Models (MLLMs) in vision–language tasks such as visual question answering and reasoning. However, existing vision encoders focus on global image representations but overlook fine-grained regional analysis. They are limited in fine-grained perception due to the scarcity of fine-grained annotated data and the lack of a fine-grained pre-training paradigm. In this paper, we propose GranViT, a novel Vision Transformer that integrates fine-grained feature extraction with semantic alignment to Large Language Models (LLMs) via region-level autoregressive training. We first construct *Gran-29M*, a dataset comprising 29 million natural and OCR images paired with over 180 million high-quality region-level annotations, to enable large-scale fine-grained pretraining. Consequently, we develop a pretraining-adaptation framework along with a self-distillation mechanism to train fine-grained GranViT on *Gran-29M*. We sufficiently exploit the fine-grained annotations from *Gran-29M* to resort to bounding-box-to-caption regression to enhance localized visual representation of the vision encoder in the pretraining and caption-to-bounding-box regression to improve vision feature utilization and localization for LLM in the adaptation. We further incorporate a self-distillation mechanism that imposes explicit localization constraints on the vision encoder to strengthen its regional reasoning capability. Extensive experiments show that GranViT surpasses existing vision encoders and attains strong transferability to varying LLMs. Remarkably, it achieves state-of-the-art results on fine-grained recognition, multimodal VQA, and OCR understanding.

1 INTRODUCTION

Multimodal Large Language Models (MLLMs) have stimulated substantially growing research interests and efforts in recent years (Wang et al., 2024; Bai et al., 2025; Dong et al., 2025; Zhu et al., 2025; Wang et al., 2025; Li et al., 2025). Existing architectures for MLLMs usually consist of a pretrained vision encoder that extracts visual information and a projection module that maps visual information to visual tokens for image understanding and reasoning with Large Language Models (LLMs). Projection modules such as multilayer perceptrons (MLPs) (Liu et al., 2023c) or Q-Formers (Li et al., 2023b) bridge visual features to the semantic space of LLMs, whereas vision encoders are primarily for the ability of capturing visual information for MLLMs.

Vision Transformers (ViTs) (Dosovitskiy et al., 2020) and their variants (Liu et al., 2021; Ravi et al., 2024) have been widely adopted as vision encoders in MLLMs due to their exceptional capabilities in visual feature extraction and scalability (Dosovitskiy et al., 2020; Carion et al., 2020; Kirillov et al., 2023; Ravi et al., 2024). Existing ViTs are usually trained to align visual representations with textual semantics. Contrastive Language–Image Pre-training (CLIP) (Radford et al., 2021; Zhai et al., 2023; Tschannen et al., 2025; Shi et al., 2024) is one prevailing paradigm that projects images and texts into a learned shared embedding space to aggregate matched image-text pairs, and non-matching pairs are discriminated to preserve the semantic relationship. Another popular alternative is autoregressive modeling (Chen et al., 2024c; Fini et al., 2025; Tschannen et al., 2025) that directly maps visual features into the textual space by connecting a cascaded vision encoder and textual decoder. It allows superior alignment with textual space but could sacrifice the discrimination

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

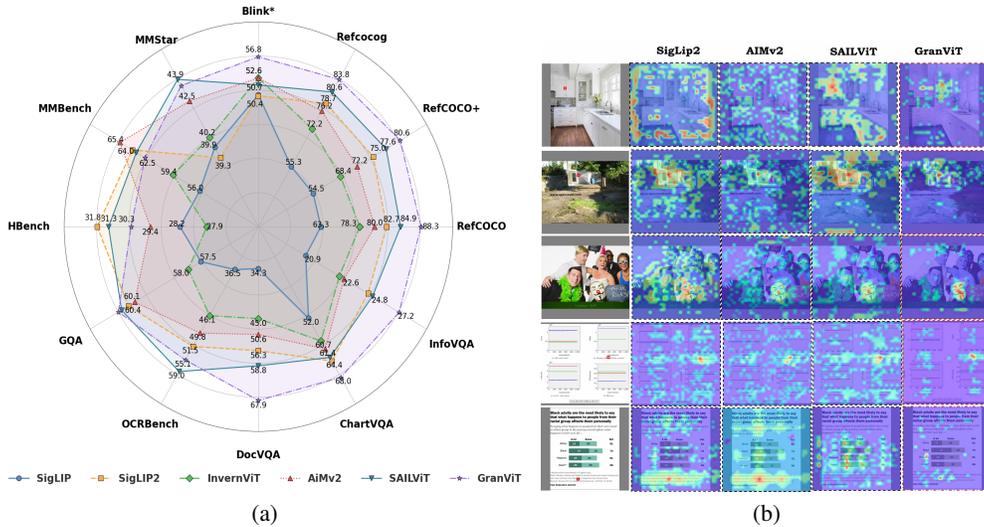


Figure 1: (a) Compared to existing vision encoders, GranViT demonstrates outstanding performance across fine-grained natural image and OCR understanding. HBench denotes HallusionBench. (b) Attention visualization of existing vision encoders according to the query token. The small red rectangle indicates the query token. Best viewed with zoom in.

ability of visual features. Nevertheless, these approaches over-emphasize image-level global feature extraction but neglect essential fine-grained details required for multimodal understanding.

To address the limitation, in this paper, we investigate integrating fine-grained localization capabilities into the vision encoder within an LLM cascade architecture. It is non-trivial to address the following two challenges, *i.e.*, i) Data scarcity: scarcity of high-quality datasets with fine-grained annotations, and ii) Fine-grained pre-training: lack of a dedicated framework to train fine-grained vision encoders that effectively align with LLMs.

i) Data scarcity. We build a high-quality annotated dataset termed *Gran-29M* that contains 29 million natural and OCR images with image-level annotations along with 183 million region-level annotations. Specifically, we leverage the UMG-41M (Shi et al., 2024), FLICKR30k (Young et al., 2014), and LAION (Schuhmann et al., 2022) datasets to collect natural images of varying scales and diversity and generate image-level and region-level annotations (e.g., bounding boxes) using ViTDet (Li et al., 2022) and Qwen2.5-VL (Bai et al., 2025). Moreover, we consolidate publicly available OCR datasets (Li et al., 2025; 2024a) and utilize PaddleOCR (Cui et al., 2025) for localized text detection and bounding box prediction. *Gran-29M* is achieved with rigorous filtering based on bounding box aspect ratio, area, quantity, and image resolution.

ii) Fine-grained pre-training and adaptation. We propose a novel pretraining-adaptation framework to improve fine-grained understanding of natural and OCR images beyond enhancing overall perception. In the pretraining, the proposed framework optimizes the vision encoder with bounding-box-to-caption (*Bbox2Caption*) regression for fine-grained feature extraction. Additionally, we develop localized self-distillation to optimize the vision encoder and explicitly augment its ability to extract fine-grained features. As for adaptation, the LLM is tunable for fine-grained vision feature localization with caption-to-bounding-box (*Caption2Bbox*) regression.

To validate the effectiveness of GranViT, we perform comprehensive performance comparisons and extensive visualizations after downstream supervised fine-tuning (SFT) (Li et al., 2024b), including visual question answering, visual grounding, and OCR understanding. Fig. 1 shows that GranViT achieves state-of-the-art performance on multiple benchmarks and exhibits strong generalization capabilities. The contributions of this work are summarized as below.

- We establish *Gran-29M*, a large-scale pretraining dataset containing 29 million natural and OCR images with comprehensive global annotations and 183 million fine-grained captions.
- We propose a pretraining-adaptation framework that simultaneously enhances the fine-grained feature extraction ability of GranViT with *Bbox2Caption* regression and localized self-distillation

108 using explicit local region supervision and adapts to varying LLMs with stronger capacity for local
109 region localization with *Caption2Bbox* regression.

110 • We demonstrate the robustness and generalization ability of *GranViT* compared with existing vi-
111 sion encoders via comprehensive analysis. *GranViT* achieves state-of-the-art performance on visual
112 grounding and OCR comprehension.
113

114 2 RELATED WORK

116 2.1 MULTIMODAL LARGE LANGUAGE MODELS

117 Multimodal large language models (MLLMs) (Wang et al., 2024; Bai et al., 2025; Dong et al., 2025;
118 Chen et al., 2024c;b; Zhu et al., 2025; Wang et al., 2025; Team et al., 2025b; Li et al., 2025; Lei et al.,
119 2025b) attract wide attention for their potential in image understanding and reasoning. Building
120 on the robust textual understanding and reasoning capabilities of large language models (LLMs)
121 (Touvron et al., 2023; Yang et al., 2025; Guo et al., 2025a; Yang et al., 2025; Lei et al., 2025a),
122 most existing MLLMs augment their functionality with a pretrained vision encoder to enable visual
123 perception. These encoders are usually trained with contrastive learning (Radford et al., 2021; Zhai
124 et al., 2023; Tschannen et al., 2025; Shi et al., 2024) and projectors commonly adopt a two-layer
125 MLP architecture (Li et al., 2024b; Liu et al., 2023c), but pre-trained vision encoders cannot handle
126 high-resolution inputs (Zhai et al., 2023). Early models (Wang et al., 2024; Gu et al., 2024) adopt an
127 image tiling strategy (Chen et al., 2024c; Liu et al., 2023c; Lu et al., 2024; Team et al., 2025b): high-
128 resolution images are divided into patches, from which local features are extracted and aggregated.
129 In comparison, newer MLLMs such as Qwen2.5-VL (Bai et al., 2025), Seed-VL1.5 (Guo et al.,
130 2025b), and Kimi-VL (Team et al., 2025a) train vision encoders from scratch on diverse datasets
131 and support native-resolution input (Bai et al., 2025) to mitigate performance loss from resolution
132 reduction. Beyond architectural improvements, recent MLLMs increasingly focus on post-training
133 strategies (Cheng et al., 2024; Gu et al., 2024; Li et al., 2025). These emphasize large-scale, curated
134 SFT datasets and leverage both SFT and reinforcement learning (Schulman et al., 2017; Shao et al.,
135 2024; Zheng et al., 2024b; 2025; 2024a) to enhance task-specific capability.

136 2.2 VISION FOUNDATION MODELS

137 Vision encoders are a critical component for extracting and representing visual information to sup-
138 port multimodal reasoning in MLLMs. Existing MLLMs usually employ vision encoders pre-trained
139 through contrastive learning, which inherently align visual and textual semantic spaces. Commonly
140 used encoders include CLIP (Radford et al., 2021) using cross-entropy loss (Mao et al., 2023),
141 and SigLIP (Zhai et al., 2023) using sigmoid loss. InternViT (Chen et al., 2024c) combines con-
142 trastive learning with an autoregressive loss and incorporates a text decoder to enhance alignment
143 by decoding visual features into text. SeedViT (Guo et al., 2025b) first undergoes generative self-
144 supervised pretraining (Xie et al., 2022; He et al., 2022) before contrastive learning. AIMv2 (Fini
145 et al., 2025) introduces the first vision encoder trained solely with an autoregressive loss, predict-
146 ing subsequent image patches and text tokens to achieve cross-modal alignment without contrastive
147 learning. SigLIP2 (Tschannen et al., 2025) integrates autoregressive and self-distillation losses in
148 SigLIP to enhance visual representations through multi-objective pretraining. SAILViT (Yin et al.,
149 2025) extends AIMv2 by incorporating alignment with LLMs and multi-stage pretraining with SFT
150 data, facilitating high-dimensional vision-language integration and infusing world knowledge into
151 visual encoding. However, these encoders predominantly emphasize global feature extraction at the
152 cost of fine-grained visual details, and are limited in fine-grained multimodal tasks.

153 3 GRAN-29M: FINE-GRAINED ANNOTATED DATASET

154 In this section, we elaborate on the construction of a large-scale fine-grained *Gran-29M* dataset for
155 pre-training, including data sources, data annotations, filtering criteria, and data reformatting.

156 **Data Source.** We collect diverse large-scale images from public datasets. For natural images, we ex-
157 pand the UMG-41M dataset (Shi et al., 2024) (including CC3M (Sharma et al., 2018), IN21k (Deng
158 et al., 2009), SBU (Ordonez et al., 2011), CC12M (Changpinyo et al., 2021), YFCC15M (Kamath
159 et al., 2021), and VisualGenome (Krishna et al., 2017)) with samples from LAION (Schuhmann
160 et al., 2022) and FLICKR30k (Young et al., 2014). For OCR images, we collect four distinct types
161

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

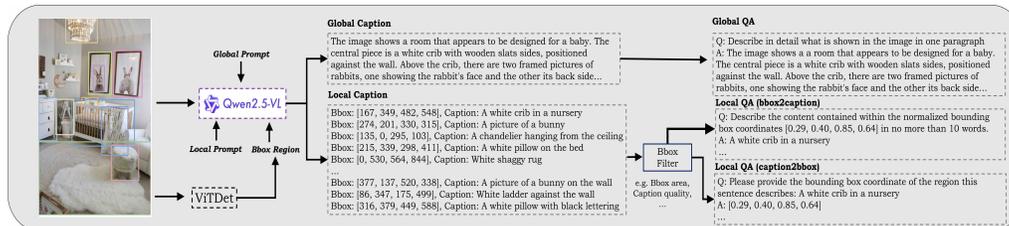


Figure 2: The details of data annotations of *Gran-29M*. We leverage ViTDet (Li et al., 2022) and Qwen2.5-VL-7B (Bai et al., 2025) for bbox and caption generation. Then, we transfer the absolute bbox coordinate to a relative one based on the image resolution and apply rigorous filtering based on image resolution, bbox area, and the number of bboxes per image. Finally, we reformat the global and local captions into QA pairs.

of images from publicly available sources (Li et al., 2024a; 2025): plain text images, chart and table images, receipt images, and rich text images. Refer to Table 6 in the appendix for details.

Data Annotations Workflow. For natural images, we directly utilize bounding box (bbox) coordinates provided by UMG-41M as localized annotations for local regions, and employ Qwen2.5-VL-7B (Bai et al., 2025) to regenerate global and local captions to enhance caption quality. For the LAION dataset (Schuhmann et al., 2022) and FLICKR30k (Young et al., 2014), we utilize ViTDet (Li et al., 2022) to detect bboxes and Qwen2.5-VL-7B (Bai et al., 2025) to generate global and local captions, as shown in Fig. 2. For OCR images, since global descriptions are often vague (e.g., “a page of an academic paper”) and lack details, only local regions are annotated using PaddleOCR (Cui et al., 2025) to provide accurate bboxes and textual contents simultaneously.

Filtering Criteria. To ensure high-quality annotations for both global and local regions, we apply a filtering process based on image resolution and bbox criteria. For local region annotations, we require that the shorter side of each image should be larger than 448 pixels, the aspect ratio of both the entire image and each bbox should be between $\frac{1}{3}$ and 3, the area of each bbox should be greater than 100^2 square pixels, and the number of bboxes per image should be at least one. The filtered results are summarized in Table 7 in the appendix. In total, we obtain 29.51 million images with 183.55 million localized region annotations for large-scale pretraining.

Data Reformatting. To facilitate the training of GranViT, we reorganize existing global and local region captions and reformat them into a standard question-answer pair structure. Using the following question and answer prompts, we rewrite existing data to enhance its suitability for training. For bbox coordinates and corresponding captions, we perform bidirectional annotations through *Bbox2Caption* and *Caption2Bbox* tasks for the vision encoder and LLM pretraining, respectively. Furthermore, we convert the absolute bbox coordinates into relative coordinates based on image resolutions to eliminate the dependence on absolute coordinates.

Global Caption.

Q: Describe in detail what is shown in the image in one paragraph
A: [global captions]

Bbox2Caption.

Q: Describe the content contained within the normalized bounding box coordinates [bbox coordinates] in no more than 10 words.
A: [local captions]

Caption2Bbox.

Q: Please provide the bounding box coordinate of the region this sentence describes: [local captions]
A: [bbox coordinations]

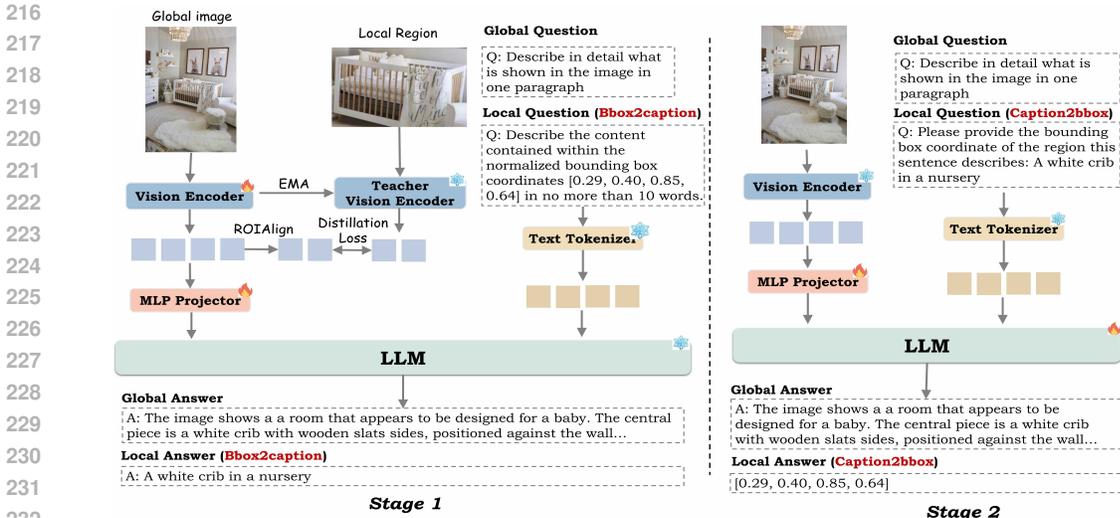


Figure 3: The fine-grained pretraining and transferring paradigm of GranViT. For pretraining, the vision encoder and projector are tuned via the global and *Bbox2Caption* task for fine-grained feature extraction. The teacher vision encoder explicitly supervises the local region of features extracted by the student vision encoder. For vision feature adaptation and transfer, based on the fine-grained vision encoder, we apply LLM tuning to further strengthen the localization capability of the LLM regarding fine-grained visual features via the global and *Caption2Bbox* task.

4 PROPOSED METHOD

4.1 FINE-GRAINED PRETRAINING PARADIGM WITH AUTOREGRESSIVE PERCEPTION

Owing to training solely on images and global captions (Radford et al., 2021; Zhai et al., 2023; Fini et al., 2025), previous vision encoders struggle with fine-grained feature extraction for local regions, while also lacking alignment between visual features and the textual feature of the LLM. To overcome these issues with a unified framework, we leverage the LLM to provide supervision for the fine-grained training of vision encoders. Specifically, we employ the same global image captioning task (Radford et al., 2021; Zhai et al., 2023; Fini et al., 2025) throughout the entire pretraining process to preserve the global perception capability. Furthermore, we enhance the ability of fine-grained feature extraction by cascading the vision encoder with the LLM via the projector during pretraining, and perform large-scale pretraining using both *Bbox2Caption* and *Caption2Bbox* tasks for localized region recognition and grounding, respectively. As depicted in Fig. 3, the proposed framework consists of pretraining and adaptation stages.

- Stage 1: Pretraining that tunes vision encoder and projector with LLM frozen.** We additionally employ the *Bbox2Caption* task for pretraining, which requires the MLLM to generate a localized caption of the object within specified bboxes. The LLM can be viewed as a decoder that converts visual features into texts, where the supervision is directly propagated back to the extracted local features with bboxes, thereby enhancing the fine-grained characteristics of the visual representations. The input prompt to the LLM incorporates the bbox coordinates, thereby facilitating object recognition and localization. It is worth noting that, in this stage, we employ a lightweight LLM (*i.e.*, Qwen2.5-VL-1.5B (Bai et al., 2025)) to compel the vision encoder to extract generic fine-grained features, rather than relying on the powerful reasoning capabilities of large LLMs for output generation.
- Stage 2: Adaptation and Transfer that tune projector and LLM with the vision encoder frozen.** In contrast, we employ the *Caption2Bbox* task in this stage, which requires the MLLM to recognize objects present in the image according to the prompts and output their bbox coordinates. The primary objectives of this stage are to further enhance the localization capability of the LLM based on fine-grained visual features and to ensure the transferability of the vision encoder to other LLMs with comparable or larger size. Since in MLLMs, it is required that the vision encoder provide more fine-grained features, while the LLM should also be capable of utilizing these visual features. On the other hand, the vision

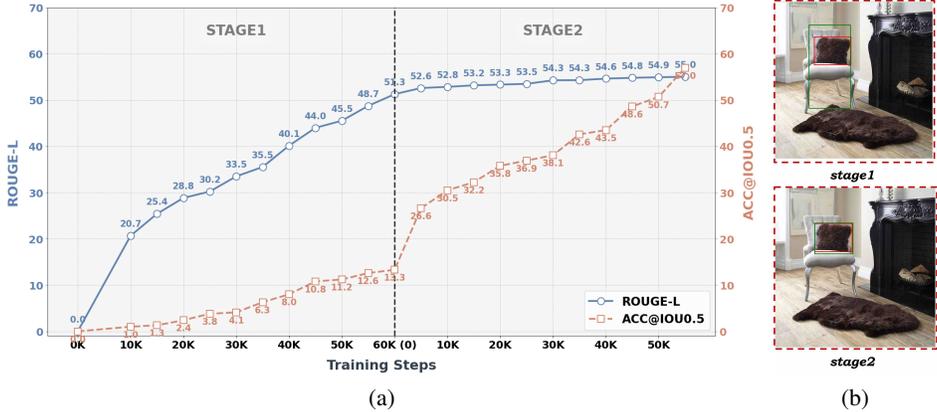


Figure 4: (a) The performance curve of Stage1 and Stage2. We sample 8M *Bbox2Caption* and *Caption2Bbox* samples respectively for pretraining and adaptation and calculate ROUGE-L (Barella & Tortora, 2022) and ACC@IOU0.5 for *Bbox2Caption* and *Caption2Bbox* respectively. In stage 1, the ACC@IOU0.5 of the *Caption2Bbox* task only achieves 13%, while the ROUGE-L of the *Bbox2Caption* task achieves 52%. Conversely, in stage 2, the training of LLM leads to a notable increase in ACC@IOU0.5 for *Caption2Bbox*, while *Bbox2Caption* achieves only a minimal improvement of 3%. (b) Visualization of predicted bbox coordinate of *Caption2Bbox* task in stage 1 and stage 2. Green bboxes indicate predicted regions, while red ones denote ground truth.

encoder can be adapted with different LLMs (*i.e.*, Qwen2.5-VL-3B, Qwen2.5-VL-7B (Bai et al., 2025)), ensuring compatibility and transferability to new architectures.

Across both stages, the autoregressive caption loss is applied to regulate the text output of the LLM for supervising the vision encoder and the LLM, respectively. Given ground truth text T and output text O_{LLM} , the caption loss can be calculated by $L_{caption} = CrossEntropy(O_{LLM}, T)$. Through this pretraining-adaptation pretraining paradigm, the vision encoder gains enhanced fine-grained feature extraction abilities, with its outputs inherently aligned to the semantic space of the LLM. The LLM, in turn, improves its capacity to utilize visual information for accurate localization. These generalized capabilities enhance object recognition and spatial reasoning during downstream SFT, thereby reducing visual understanding errors and mitigating hallucinations.

4.2 EXPLANATION OF FINE-GRAINED PRETRAINING PARADIGM

Pretraining-Adaptation framework with different tasks. To enhance the fine-grained feature extraction capability of the vision encoder, in pretraining, we freeze the LLM and only train the vision encoder and projector in stage 1. Both the *Bbox2Caption* and *Caption2Bbox* tasks are used initially; however, as illustrated in Fig. 4(a), the model shows strong performance in *Bbox2Caption* but limited accuracy in *Caption2Bbox*, primarily due to the frozen LLM, which restricts learning for the more language-reliant *Caption2Bbox*. Thus, we further apply stage 2 to adapt the pre-trained vision encoder to LLMs to utilize fine-grained visual features. Meanwhile, the pretrained vision encoder can be adapted with other LLMs for transferability. Specifically, in stage 2, we keep the vision encoder frozen and tune the projector and LLM using both *Bbox2Caption* and *Caption2Bbox* tasks. Results in Fig. 4(a) demonstrate a notable improvement in ACC@IOU0.5 for *Caption2Bbox*, while *Bbox2Caption* sees only marginal gains. This finding is consistent with the visualization results in Fig. 4(b), confirming the superiority of training *Caption2Bbox* in stage 2. Since further optimizing *Bbox2Caption* in stage 2 yields minimal benefits at twice the computational cost, we finally adopted the pretraining-adaptation paradigm, separately optimizing the two tasks in these two stages.

Freezing the vision encoder in the adaptation stage. We freeze the vision encoder and only tune the LLM for adaptation and transfer in stage 2, enabling the LLM to better use the fine-grained vision feature for the *Caption2Bbox* task. The vision encoder is kept frozen because it has already been sufficiently pretrained in stage 1, and further tuning it would significantly increase adaptation and transfer costs but yield marginal performance gains, as shown in Table 9 in the appendix.

4.3 SELF DISTILLATION FOR FINE-GRAINED PRETRAINING PARADIGM

Although the vision encoder is implicitly supervised through the caption loss $L_{caption}$ from the output text of LLM, it lacks explicit constraints on localized region features. Therefore, we incorporate self-distillation training (Naeem et al., 2024; Maninis et al., 2024; Zhang et al., 2019) into stage 1. As illustrated in Fig. 3, an additional frozen teacher vision encoder is introduced to constrain the feature generated by the student vision encoder, thereby enhancing localized region features.

Specifically, given images x , the student vision encoder extracts the fine-grained features x' . Both prompts and visual features x' are sent to the LLM for generating localized captions. Besides, the image regions x_{crop} are cropped from x according to the bbox coordinates and are sent to the teacher vision encoder for localized feature extraction. The self-distillation loss is calculated between x' and x'_{crop} by $L_{distill} = MSE(x'_{crop}, ROIAlign(x'))$, where x'_{crop} and MSE denotes the extracted features of x_{crop} and mean square error loss, respectively. The weights of the teacher vision encoder are initialized from the student vision encoder and updated by the exponential moving average (EMA) according to the student vision encoder, specifically $\theta_{tea} = \alpha\theta_{tea} + (1-\alpha)\theta_{stu}$. The overall loss can be written as $L = L_{caption} + \lambda L_{distill}$, where λ denotes the weighting coefficient.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETTINGS

Dataset. Three datasets are adopted for distinct purposes of projector pre-training, fine-grained pretraining, and downstream SFT. Following (Liu et al., 2023c), we adopt the BLIP-LAION-CC-SBU-558K dataset (Liu, 2024) for projector pretraining. For fine-grained pretraining, we employ the constructed *Gran-29M* dataset. To balance natural and OCR images, we proportionally sample 50M local region captions from OCR regions and use all the natural image local captions. During stage 1, we use all global samples and approximately 130M *Bbox2Caption* samples. In stage 2, we sample 24M *Caption2Bbox* samples in addition to global samples to reduce the transfer overhead. For SFT, we adopt the Open-LLaVA-NeXT 1M dataset (Chen & Xing, 2024) for downstream adaptation.

Implementation Details. We use the LLaVA-Next framework (Li et al., 2024b) for pre-training and SFT of the vision encoder and LLM. By default, we initialize the vision encoder with SigLIP2 (Tschannen et al., 2025), and adopt Qwen2.5-VL-1.5B (Bai et al., 2025) as the LLM. The projector is implemented with a two-layer MLP. In training, images are resized to 512×512 , padded according to their aspect ratio, and then fed into the vision encoder and LLM for feature extraction and inference. We also implement GranViT with image tiling strategy, as shown in Table 8 in the appendix. We employ AdamW optimizer (Loshchilov & Hutter, 2017) for pretraining and SFT with learning rate 10^{-5} for one epoch. The overall batch size is set to 256 with 128 Ascend 910B NPUs. λ is 1 and α is 0.9 by default in our experiment. Abation study on λ and α is shown in Table 5 in the appendix.

5.2 BENCHMARK EVALUATION

We make extensive evaluations on the well-known OpenCompass benchmark (Contributors, 2023) and additional fine-grained benchmarks. We focus on fine-grained and OCR benchmarks that are divided into four classes: fine-grained (RefCOCO (Yu et al., 2016), RefCOCO+ (Yu et al., 2016), RefCOCOG (Yu et al., 2016), BLINK* (Fu et al., 2024)¹ and MMVP (Tong et al., 2024)), multitmodal VQA (MMBench (Liu et al., 2024a), MMStar (Chen et al., 2024a), HallusionBench (Guan et al., 2024), GQA (Hudson & Manning, 2019) and SEEDBench (Li et al., 2023a)), multimodal reasoning (MMMU (Yue et al., 2024), MathVista MINI (Lu et al., 2023), MMVet (Yu et al., 2023), ScienceQA (Lu et al., 2022) and AI2D (Kembhavi et al., 2016)) and OCR understanding (OCRBench (Liu et al., 2024b), DocVQA (Mathew et al., 2021), ChartQA (Masry et al., 2022), InfoVQA (Mathew et al., 2022) and TextVQA (Singh et al., 2019)). We compare GranViT with diverse vision encoders, *i.e.*, CLIP (Radford et al., 2021), SigLip (Zhai et al., 2023), SigLip2 (Tschannen et al., 2025), AIMv2 (Fini et al., 2025), InternViT (Chen et al., 2024c), and SAILViT (Yin et al., 2025).

¹We calculate the average score of fine-grained evaluation (Counting, Object Localization and Spatial Relation) in BLINK, denoted as BLINK*

Table 1: Performance comparison with low resolution version. The bold font represents the best performance, and the underline represents the second performance.

Capability	Benchmark	CLIP	SigLip	SigLip2	AIMv2	InternViT	SAILViT	GranViT
Fine-Grained	RefCOCO_testA	81.26	69.47	87.78	86.03	85.15	<u>89.65</u>	91.79
	RefCOCO_testB	64.51	56.78	76.90	73.54	71.40	<u>79.82</u>	83.88
	RefCOCO_val	74.71	63.71	83.26	80.28	78.48	<u>85.32</u>	89.13
	RefCOCO+_testA	74.43	63.13	82.92	81.33	77.33	<u>85.01</u>	87.04
	RefCOCO+_testB	51.25	44.65	66.47	62.50	58.58	<u>69.66</u>	73.24
	RefCOCO+_val	65.25	55.84	75.46	72.84	69.18	<u>78.10</u>	81.55
	RefCOCOg_val	68.95	55.02	78.53	76.55	71.62	<u>80.26</u>	83.86
	RefCOCOg_test	68.52	55.57	78.94	75.78	72.71	<u>80.92</u>	83.82
	BLINK*	51.87	50.67	50.35	52.59	<u>52.62</u>	52.54	56.80
	MMVP	63.33	61.66	66.00	65.66	61.00	69.00	<u>66.33</u>
	Average	66.41	57.67	75.61	73.50	70.53	<u>77.95</u>	80.78
VQA	MMBench	61.14	55.95	64.00	65.40	59.44	<u>63.54</u>	62.46
	MMStar	39.46	39.93	39.33	42.53	40.20	43.86	<u>43.73</u>
	HallusionBench	28.83	28.22	31.79	29.39	27.95	<u>31.29</u>	30.34
	GQA	58.80	57.52	60.42	60.14	58.02	<u>60.83</u>	60.95
	SEEDBench	66.89	66.28	69.30	<u>70.20</u>	66.79	69.75	70.36
		Average	51.02	49.58	52.97	<u>53.53</u>	50.48	53.85
Reasoning	MMMU	40.00	38.66	<u>42.00</u>	38.66	44.00	38.66	38.00
	MathVista MINI	38.10	35.50	38.40	37.70	39.50	41.70	<u>40.40</u>
	MMVet	33.53	30.87	38.80	<u>38.34</u>	35.27	35.73	37.29
	ScienceQA	<u>67.87</u>	66.53	66.98	<u>67.57</u>	66.63	72.78	67.42
	AI2D	66.51	65.47	69.26	68.19	66.51	71.21	<u>69.91</u>
		Average	49.20	47.41	<u>51.09</u>	50.09	50.38	52.02
OCR	OCRBench	406	365	515	498	461	590	<u>551</u>
	DocVQA	35.34	34.26	56.32	50.56	44.95	<u>58.75</u>	67.92
	ChartQA	50.84	51.96	64.44	61.36	60.68	<u>63.24</u>	67.96
	InfoVQA	20.67	20.89	24.12	22.60	21.91	<u>24.75</u>	27.19
	TextVQA	46.65	41.71	<u>61.36</u>	56.57	51.74	60.90	61.66
		Average	38.82	37.06	51.55	48.18	45.08	<u>53.33</u>

Performance Comparison. Table 1 provides a performance comparison on various benchmarks. For fine-grained and OCR tasks, GranViT achieves an average top-1 score of 80.78 and 55.97, and surpasses the second best by 2.83 and 2.64, respectively. GranViT is comparable to SAILViT with a marginal difference of only 0.3 in multimodal VQA tasks. For multimodal reasoning tasks, GranViT suffers a slight performance loss of 0.4 compared to SigLIP2. This is because reasoning capability does not rely heavily on fine-grained feature extraction, and GranViT prioritizes fine-grained tasks rather than extensively training reasoning capabilities. Note that reasoning capability can be further improved by applying reasoning VQA data in pretraining.

Transferability. Table 2 reports the performance comparison with larger LLMs (*i.e.*, Qwen2.5-3B, Qwen2.5-7B). Notably, other vision encoders directly employ the larger LLM for SFT, whereas GranViT employs a lightweight LLM (*i.e.*, Qwen2.5-1.5B) for pre-training in stage 1, transfers to the larger LLM in stage 2, and subsequently undergoes SFT. GranViT also demonstrates outstanding performance on fine-grained and OCR tasks, while achieving comparable or even state-of-the-art performance on some VQA tasks (*i.e.*, HallusionBench and SEEDBench).

5.3 SCALING LAWS

We evaluate the scaling capacity of pretraining-adaptation framework in Fig. 5. Specifically, for stage 1, we leverage 8M, 16M and all the 130M regions for *Bbox2Caption* tasks, while we leverage 8M, 16M, 24M and all the 130M regions for *Caption2Bbox* tasks in stage 2. The average score of fine-grained (RefCOCO_testA, RefCOCO+_testA, RefCOCOg_test, BLINK, MMVP) and OCR (OCRBench, DocVQA, ChartQA, InfoVQA, TextVQA) tasks is reported. As the data scale increases, both tasks exhibit significant performance improvements, indicating enhanced fine-grained feature extraction capability of GranViT.

Table 2: Performance comparison for transferring vision encoders to Qwen2.5-3B, Qwen2.5-7B and LLaMA3-8B. The best results are highlighted in bold and the second best underlined. Ref, Ref+ and Refg denote the RefCOCO_testA, RefCOCO+_testA and RefCOCOg_test. MMB, HB, and SB stand for MMBench, HallusionBench, and SEEDBench, and. SQA, OB, DVQA, and IVQA for ScienceQA, OCRBench, DocVQA, and InfoVQA, respectively.

Model	Ref	Ref+	Refg	MMB	HB	SB	MMM	SQA	OB	DVQA	IVQA	Avg
Qwen2.5-3B												
CLIP	86.60	81.73	75.26	65.09	31.67	68.55	<u>39.33</u>	73.12	413	38.70	24.08	56.86
SigLip	87.55	82.83	77.46	69.27	33.12	70.29	36.44	70.64	428	46.27	25.24	58.36
SigLip2	91.03	86.37	83.30	<u>69.34</u>	33.70	71.13	36.00	71.69	529	60.87	29.00	62.30
AIMv2	90.33	87.05	81.83	69.27	31.91	71.84	40.00	<u>74.46</u>	545	56.04	27.34	62.23
SAILViT	<u>91.58</u>	<u>87.82</u>	<u>83.63</u>	69.73	34.08	71.33	36.00	75.75	633	<u>62.53</u>	<u>29.49</u>	64.11
GranViT	93.22	89.32	86.17	67.56	<u>33.77</u>	72.34	38.66	73.24	<u>590</u>	71.09	29.96	64.94
Qwen2.5-7B												
CLIP	90.01	86.25	80.44	70.58	36.39	70.70	46.00	75.26	466	43.02	24.88	60.92
SigLip	90.68	86.27	81.72	<u>74.22</u>	<u>37.94</u>	72.04	46.00	<u>76.99</u>	459	50.03	26.69	62.59
SigLip2	92.06	88.70	85.13	72.21	36.02	72.17	44.00	75.26	540	62.20	29.89	64.69
AIMv2	90.84	87.82	82.39	72.29	37.29	72.11	41.44	72.92	553	56.54	28.59	63.41
SAILViT	92.66	<u>89.50</u>	<u>85.32</u>	74.53	37.76	73.05	<u>44.66</u>	81.11	648	<u>64.13</u>	<u>30.66</u>	67.11
GranViT	92.98	90.46	87.96	73.37	39.37	74.45	<u>44.66</u>	75.85	<u>582</u>	73.14	31.69	67.47
LLaMA3-8B												
CLIP	90.26	86.43	80.59	72.07	37.65	72.18	<u>48.00</u>	74.90	462	44.28	24.14	61.38
SigLip	90.47	86.44	81.44	<u>74.02</u>	35.25	74.55	48.97	76.48	473	50.69	23.37	62.63
SigLip2	92.85	88.59	85.59	72.31	37.99	74.36	45.62	76.04	562	64.72	27.56	65.62
AIMv2	91.30	88.03	83.67	72.49	37.02	73.77	40.77	73.81	529	58.85	24.04	63.33
SAILViT	93.06	<u>89.94</u>	<u>85.89</u>	74.08	<u>39.02</u>	<u>75.35</u>	45.66	79.59	640	68.97	28.08	67.60
GranViT	93.65	91.27	88.23	73.08	39.92	76.37	47.14	<u>78.42</u>	<u>625</u>	77.24	31.44	69.02

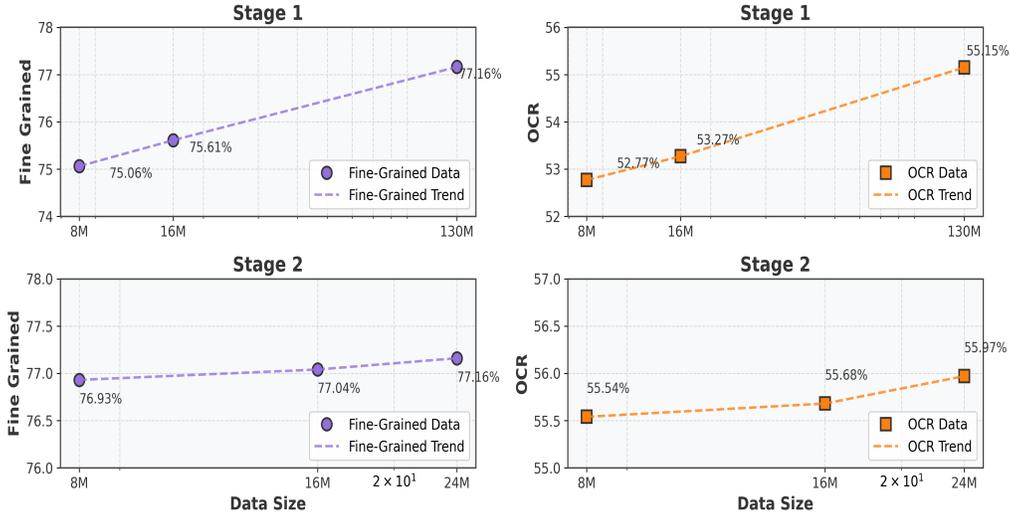


Figure 5: Scaling law of two-stage training.

5.4 ABLATION STUDY

We employ small-scale datasets for ablation studies on the contribution of each module. 8 million global and *Bbox2Caption* QA pairs and 8 million global and *Caption2Bbox* QA pairs are sampled from *Gran-29M* as training data for two stages, respectively. The entire dataset is used for training the projector and conducting supervised fine-tuning (SFT).

Effectiveness of Training Paradigm. In Table 3, the progressive introduction of the two-stage training strategy along with self-distillation results in incremental performance gains on fine-grained and OCR-related tasks. Specifically, with stage 1 pretraining, the MLLM exhibits a substantial improvement in both fine-grained recognition capability and OCR understanding (2.2 and 1.2 gains). Self-

Table 3: Ablation study on each component of the proposed GranViT.

SigLip2	Stage1	Self-Distillation	Stage2	Fine-Grained	VQA	Reasoning	OCR
✓	✗	✗	✗	73.20	52.97	51.09	51.55
✓	✓	✗	✗	75.06	53.64	49.89	52.77
✓	✓	✓	✗	75.55	53.90	50.32	53.02
✓	✓	✓	✓	76.54	53.77	48.99	53.78

Table 4: Performance with different vision encoder initialization for GranViT during pretraining.

Model	Fine-Grained	VQA	Reasoning	OCR
InternViT	69.76	50.48	50.38	45.08
GranViT (InternViT)	75.15	51.78	50.23	50.13
AIMv2	72.28	53.53	50.09	48.18
GranViT (AIMv2)	77.14	55.07	50.59	52.71
SAILViT	75.42	53.85	52.02	54.53
GranViT (SAILViT)	76.79	55.40	51.95	56.61

Table 5: Ablation Study of the coefficient in self-distillation.

λ	α	Fine-Grained	VQA	Reasoning	OCR
1	0.9	75.55	53.90	50.32	53.02
1	0.99	75.25	53.45	50.89	53.30
1	0.999	74.86	53.53	50.13	52.31
0.01	0.9	74.81	53.47	50.13	52.80
0.1	0.9	74.75	53.48	50.58	53.16
0.5	0.9	75.02	53.80	51.65	53.32

distillation training further improves fine-grained and OCR evaluations. With stage 2 adaptation, fine-grained and OCR evaluations yield additional gains of 1.0 and 0.7, respectively.

Different Initialization for Vision Encoder. In Table 4, we compare the performance of vision encoders initialized with different models. All three distinct vision encoders (InternViT-300M (Chen et al., 2024c), AIMv2 (Fini et al., 2025), and SAILViT-Huge (Yin et al., 2025)) exhibit significant performance improvements after pre-training, with the most notable improvements in fine-grained perception (*i.e.*, 5.3 for InternViT, 4.8 for AIMv2, and 1.3 for SAILViT) and OCR understanding (*i.e.*, 5.1 for InternViT, 4.5 for AIMv2, and 2.1 for SAILViT).

Coefficient in Self-Distillation We ablate two parameters in the self-distillation process: λ and α . To efficiently conduct the ablation experiments, during the pre-training stage, we only train stage 1 and then directly proceed to downstream SFT. As shown in Table 5, the evaluation performance of the model on fine-grained tasks gradually improves as λ increases and α decreases. Therefore, we set λ to 1 and α to 0.9 in our experiments.

Visualization. To illustrate the fine-grained feature extraction capability of GranViT, we visualize in Fig. 1(b) the attention maps of different vision encoders. AIMv2 (Fini et al., 2025) and SAILViT (Yin et al., 2025) focus on global features and are severely deficient in local regions. SigLip2 (Tschannen et al., 2025) emphasizes local regions, but exhibits redundant attention to global features. In contrast, GranViT can simultaneously consider local regions and exclude interference from redundant features. This further validates the effectiveness of the proposed pre-training framework.

6 CONCLUSION

This paper proposed GranViT, a novel visual Transformer architecture that integrates fine-grained perception and multimodal alignment for advanced multimodal understanding. GranViT is trained on *Gran-29M*, a newly curated large-scale dataset containing global and region-level descriptive annotations for both natural and OCR images. The region-level bounding box and text annotations enable two dedicated tasks, *i.e.*, *Bbox2Caption* for optimizing the vision encoder to strengthen fine-grained feature extraction and *Caption2Bbox* for adapting vision features to different LLMs with enhanced region localization. Self-distillation loss is further incorporated to explicitly enhance local feature learning. GranViT is potential to serve as a robust foundation MLLM model that offers strong capabilities for complex multimodal reasoning tasks.

REFERENCES

- 540
541
542 abhayzala. abhayzala/AI2D-Caption · Datasets at Hugging Face — huggingface.co. <https://huggingface.co/datasets/abhayzala/AI2D-Caption>, 2024. [Accessed 23-09-
543 2025].
544
- 545 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang,
546 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*,
547 2025.
- 548 Marcello Barbella and Genoveffa Tortora. Rouge metric evaluation for text summarization tech-
549 niques. *Available at SSRN 4120317*, 2022.
550
- 551 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
552 Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on*
553 *computer vision*, pp. 213–229. Springer, 2020.
- 554 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing
555 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the*
556 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 3558–3568, 2021.
557
- 558 Lin Chen and Long Xing. Open-llava-next: An open-source implementation of llava-next se-
559 ries for facilitating the large multi-modal model community. [https://github.com/](https://github.com/xiaochen98/Open-LLaVA-NeXT)
560 [xiaochen98/Open-LLaVA-NeXT](https://github.com/xiaochen98/Open-LLaVA-NeXT), 2024.
- 561 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
562 Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language
563 models? *Advances in Neural Information Processing Systems*, 37:27056–27087, 2024a.
564
- 565 Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shen-
566 glong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source
567 multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*,
568 2024b.
- 569 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
570 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
571 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer*
572 *vision and pattern recognition*, pp. 24185–24198, 2024c.
- 573 Daixuan Cheng, Shaohan Huang, Ziyu Zhu, Xintong Zhang, Wayne Xin Zhao, Zhongzhi Luan,
574 Bo Dai, and Zhenliang Zhang. On domain-adaptive post-training for multimodal large language
575 models. *arXiv preprint arXiv:2411.19930*, 2024.
576
- 577 Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang
578 Lou, and Dongmei Zhang. Hitab: A hierarchical table dataset for question answering and natural
579 language generation. *arXiv preprint arXiv:2108.06712*, 2021.
- 580 OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models,
581 2023.
- 582 Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiaxuan Liu, Xueqing Wang,
583 Zelun Zhang, Changda Zhou, Hongen Liu, et al. Paddleocr 3.0 technical report. *arXiv preprint*
584 *arXiv:2507.05595*, 2025.
585
- 586 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
587 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
588 pp. 248–255. Ieee, 2009.
- 589 Yihao Ding, Siwen Luo, Hyunsuk Chung, and Soyeon Caren Han. Vqa: A new dataset for real-
590 world vqa on pdf documents. In *Joint European Conference on Machine Learning and Knowledge*
591 *Discovery in Databases*, pp. 585–601. Springer, 2023.
592
- 593 Hongyuan Dong, Zijian Kang, Weijie Yin, Xiao Liang, Chao Feng, and Jiao Ran. Scalable vision
language model training via high quality data curation. *arXiv preprint arXiv:2501.05952*, 2025.

- 594 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
595 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
596 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
597 *arXiv:2010.11929*, 2020.
- 598 Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter, Michal Klein, David Haldimann, Sai
599 Aitharaju, Victor G Turrisi da Costa, Louis Béthune, Zhe Gan, et al. Multimodal autoregres-
600 sive pre-training of large vision encoders. In *Proceedings of the Computer Vision and Pattern*
601 *Recognition Conference*, pp. 9641–9654, 2025.
- 602 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A
603 Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but
604 not perceive. In *European Conference on Computer Vision*, pp. 148–166. Springer, 2024.
- 605 Shuhao Gu, Jialing Zhang, Siyuan Zhou, Kevin Yu, Zhaohu Xing, Liangdong Wang, Zhou Cao,
606 Jintao Jia, Zhuoyi Zhang, Yixuan Wang, et al. Infinity-mm: Scaling multimodal performance
607 with large-scale and high-quality instruction data. *arXiv preprint arXiv:2410.18558*, 2024.
- 608 Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang
609 Chen, Furong Huang, Yaser Yacoub, et al. Hallusionbench: an advanced diagnostic suite for
610 entangled language hallucination and visual illusion in large vision-language models. In *Pro-*
611 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14375–
612 14385, 2024.
- 613 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
614 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
615 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025a.
- 616 Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang,
617 Jianyu Jiang, Jiawei Wang, et al. Seed1. 5-vl technical report. *arXiv preprint arXiv:2505.07062*,
618 2025b.
- 619 Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in nat-
620 ural images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
621 pp. 2315–2324, 2016.
- 622 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
623 toencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*
624 *vision and pattern recognition*, pp. 16000–16009, 2022.
- 625 Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawa-
626 har. Icdar2019 competition on scanned receipt ocr and information extraction. In *2019 Interna-*
627 *tional Conference on Document Analysis and Recognition (ICDAR)*, pp. 1516–1520. IEEE, 2019.
- 628 Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning
629 and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer*
630 *vision and pattern recognition*, pp. 6700–6709, 2019.
- 631 Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and
632 Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint*
633 *arXiv:1710.07300*, 2017.
- 634 Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Car-
635 ion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the*
636 *IEEE/CVF international conference on computer vision*, pp. 1780–1790, 2021.
- 637 Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul
638 Hoque, and Shafiq Joty. Chart-to-text: A large-scale benchmark for chart summarization. *arXiv*
639 *preprint arXiv:2203.06486*, 2022.
- 640 Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali
641 Farhadi. A diagram is worth a dozen images. In *European conference on computer vision*, pp.
642 235–251. Springer, 2016.

- 648 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
649 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-*
650 *ings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
- 651
- 652 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie
653 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting lan-
654 guage and vision using crowdsourced dense image annotations. *International journal of computer*
655 *vision*, 123(1):32–73, 2017.
- 656 Jianfeng Kuang, Wei Hua, Dingkan Liang, Mingkun Yang, Deqiang Jiang, Bo Ren, and Xiang
657 Bai. Visual information extraction in the wild: practical dataset and end-to-end solution. In
658 *International Conference on Document Analysis and Recognition*, pp. 36–53. Springer, 2023.
- 659
- 660 Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and bet-
661 ter understanding vision-language models: insights and future directions. *arXiv preprint*
662 *arXiv:2408.12637*, 2024.
- 663 Wenhui Lei, Hanyu Chen, Zitian Zhang, Luyang Luo, Qiong Xiao, Yannian Gu, Peng Gao, Yankai
664 Jiang, Ci Wang, Guangtao Wu, et al. A data-efficient pan-tumor foundation model for oncology
665 ct interpretation. *arXiv preprint arXiv:2502.06171*, 2025a.
- 666
- 667 Wenhui Lei, Wei Xu, Kang Li, Xiaofan Zhang, and Shaoting Zhang. Medlsam: Localize and seg-
668 ment anything model for 3d ct images. *Medical Image Analysis*, 99:103370, 2025b.
- 669
- 670 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan
671 Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint*
672 *arXiv:2408.03326*, 2024a.
- 673 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Bench-
674 marking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*,
675 2023a.
- 676
- 677 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
678 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*
679 *preprint arXiv:2407.07895*, 2024b.
- 680 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
681 pre-training with frozen image encoders and large language models. In *International conference*
682 *on machine learning*, pp. 19730–19742. PMLR, 2023b.
- 683
- 684 Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multi-
685 modal arxiv: A dataset for improving scientific comprehension of large vision-language models.
686 *arXiv preprint arXiv:2403.00231*, 2024c.
- 687 Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhoujun Li, and Ming Zhou. Docbank:
688 A benchmark dataset for document layout analysis. *arXiv preprint arXiv:2006.01038*, 2020.
- 689
- 690 Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer back-
691 bones for object detection. In *European conference on computer vision*, pp. 280–296. Springer,
692 2022.
- 693
- 694 Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang,
695 Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from
696 scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025.
- 697 Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large
698 multi-modal model with robust instruction tuning. *CoRR*, 2023a.
- 699
- 700 Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Ya-
701 coob, and Dong Yu. Mmc: Advancing multimodal chart understanding with large-scale instruc-
tion tuning. *arXiv preprint arXiv:2311.10774*, 2023b.

- 702 Haotian Liu. liuhaotian/LLaVA-Pretrain · Datasets at Hugging Face — huggingface.co. [https://](https://huggingface.co/datasets/liuhaotian/LLaVA-Pretrain)
703 huggingface.co/datasets/liuhaotian/LLaVA-Pretrain, 2024. [Accessed 23-
704 09-2025].
- 705 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
706 *in neural information processing systems*, 36:34892–34916, 2023c.
- 707
708 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan,
709 Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around
710 player? In *European conference on computer vision*, pp. 216–233. Springer, 2024a.
- 711 Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin,
712 Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large
713 multimodal models. *Science China Information Sciences*, 67(12):220102, 2024b.
- 714 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
715 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
716 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 717 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
718 *arXiv:1711.05101*, 2017.
- 719 Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren,
720 Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding.
721 *arXiv preprint arXiv:2403.05525*, 2024.
- 722 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
723 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
724 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521,
725 2022.
- 726 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-
727 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of
728 foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- 729 Kevis-Kokitsi Maninis, Kaifeng Chen, Soham Ghosh, Arjun Karpur, Koert Chen, Ye Xia, Bingyi
730 Cao, Daniel Salz, Guangxing Han, Jan Dlabal, et al. Tips: Text-image pretraining with spatial
731 awareness. *arXiv preprint arXiv:2410.16512*, 2024.
- 732 Anqi Mao, Mehryar Mohri, and Yutao Zhong. Cross-entropy loss functions: Theoretical analysis
733 and applications. In *International conference on Machine learning*, pp. 23803–23828. pmlr, 2023.
- 734 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
735 mark for question answering about charts with visual and logical reasoning. *arXiv preprint*
736 *arXiv:2203.10244*, 2022.
- 737 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document
738 images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*,
739 pp. 2200–2209, 2021.
- 740 Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar.
741 Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer*
742 *Vision*, pp. 1697–1706, 2022.
- 743 Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual
744 question answering by reading text in images. In *2019 international conference on document*
745 *analysis and recognition (ICDAR)*, pp. 947–952. IEEE, 2019.
- 746 Muhammad Ferjad Naeem, Yongqin Xian, Xiaohua Zhai, Lukas Hoyer, Luc Van Gool, and Fed-
747 erico Tombari. Silc: Improving vision language pretraining with self-distillation. In *European*
748 *Conference on Computer Vision*, pp. 38–55. Springer, 2024.
- 749 Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million
750 captioned photographs. *Advances in neural information processing systems*, 24, 2011.

- 756 Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwal-
757 suk Lee. Cord: a consolidated receipt dataset for post-ocr parsing. In *Workshop on Document*
758 *Intelligence at NeurIPS 2019*, 2019.
- 759 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
760 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
761 models from natural language supervision. In *International conference on machine learning*, pp.
762 8748–8763. PmLR, 2021.
- 763
764 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham
765 Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images
766 and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- 767
768 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
769 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
770 open large-scale dataset for training next generation image-text models. *Advances in neural in-*
771 *formation processing systems*, 35:25278–25294, 2022.
- 772 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
773 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 774
775 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
776 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
777 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 778 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,
779 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th*
780 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
781 2556–2565, 2018.
- 782
783 Bowen Shi, Peisen Zhao, Zichen Wang, Yuhang Zhang, Yaoming Wang, Jin Li, Wenrui Dai, Junni
784 Zou, Hongkai Xiong, Qi Tian, et al. Umg-clip: A unified multi-granularity vision generalist for
785 open-world understanding. In *European Conference on Computer Vision*, pp. 259–277. Springer,
786 2024.
- 787 Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for
788 image captioning with reading comprehension. In *European conference on computer vision*, pp.
789 742–758. Springer, 2020.
- 790 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,
791 and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF*
792 *conference on computer vision and pattern recognition*, pp. 8317–8326, 2019.
- 793
794 Brandon Smock, Rohith Pesala, and Robin Abraham. Pubtables-1m: Towards comprehensive ta-
795 ble extraction from unstructured documents. In *Proceedings of the IEEE/CVF Conference on*
796 *Computer Vision and Pattern Recognition*, pp. 4634–4642, 2022.
- 797
798 Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen,
799 Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. Kimi-vl technical report. *arXiv preprint*
800 *arXiv:2504.07491*, 2025a.
- 801 Kwai Keye Team, Biao Yang, Bin Wen, Changyi Liu, Chenglong Chu, Chengru Song, Chongling
802 Rao, Chuan Yi, Da Li, Dunju Zang, et al. Kwai keye-vl technical report. *arXiv preprint*
803 *arXiv:2507.01949*, 2025b.
- 804
805 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
806 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF*
807 *Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.
- 808 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
809 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

- 810 Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdul-
811 mohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2:
812 Multilingual vision-language encoders with improved semantic understanding, localization, and
813 dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- 814 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
815 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
816 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- 817 Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu,
818 Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal
819 models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025.
- 820 wendlerc. wendlerc/RenderedText · Datasets at Hugging Face — huggingface.co. <https://huggingface.co/datasets/wendlerc/RenderedText>, 2024. [Accessed 23-09-
821 2025].
- 822 Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.
823 Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF*
824 *conference on computer vision and pattern recognition*, pp. 9653–9663, 2022.
- 825 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
826 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
827 *arXiv:2505.09388*, 2025.
- 828 Weijie Yin, Dingkang Yang, Hongyuan Dong, Zijian Kang, Jiacong Wang, Xiao Liang, Chao Feng,
829 and Jiao Ran. Sailvit: Towards robust and generalizable visual backbones for mllms via gradual
830 feature refinement. *arXiv preprint arXiv:2507.01643*, 2025.
- 831 Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual
832 denotations: New similarity metrics for semantic inference over event descriptions. *Transactions*
833 *of the association for computational linguistics*, 2:67–78, 2014.
- 834 Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context
835 in referring expressions. In *European conference on computer vision*, pp. 69–85. Springer, 2016.
- 836 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,
837 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv*
838 *preprint arXiv:2308.02490*, 2023.
- 839 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
840 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multi-
841 modal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF*
842 *Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- 843 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
844 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,
845 pp. 11975–11986, 2023.
- 846 Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your
847 own teacher: Improve the performance of convolutional neural networks via self distillation. In
848 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3713–3722, 2019.
- 849 Yilun Zhao, Chen Zhao, Linyong Nan, Zhenting Qi, Wenlin Zhang, Xiangru Tang, Boyu Mi, and
850 Dragomir Radev. Robut: A systematic study of table qa robustness against human-annotated
851 adversarial perturbations. *arXiv preprint arXiv:2306.14321*, 2023.
- 852 Guanghao Zheng, Yuchen Liu, Wenrui Dai, Chenglin Li, Junni Zou, and Hongkai Xiong. Mc-
853 dit: Contextual enhancement via clean-to-clean reconstruction for masked diffusion models. *Ad-*
854 *vances in Neural Information Processing Systems*, 37:97353–97377, 2024a.
- 855 Tianyi Zheng, Peng-Tao Jiang, Ben Wan, Hao Zhang, Jinwei Chen, Jia Wang, and Bo Li. Beta-tuned
856 timestep diffusion model. In *European Conference on Computer Vision*, pp. 114–130. Springer,
857 2024b.

864 Tianyi Zheng, Jiayang Zou, Peng-Tao Jiang, Hao Zhang, Jinwei Chen, Jia Wang, and Bo Li. Bidi-
865 rectional beta-tuned diffusion model. *IEEE Transactions on Pattern Analysis and Machine Intel-*
866 *ligence*, 2025.

867 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen
868 Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for
869 open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.

872 A APPENDIX

874 A.1 DECLARATION ON THE USE OF LARGE LANGUAGE MODELS

875
876 In preparing this paper, we used large language models (LLMs) solely to improve the clarity of the
877 writing. The core motivation, method design, and experimental setup are conceived and developed
878 entirely by ourselves, without the involvement of LLMs. In addition, during the construction of the
879 Gran-29M dataset, we employ Qwen2.5-7B for recaptioning image descriptions to improve textual
880 quality. All final writing, editing, and formatting of the manuscript are carefully reviewed and
881 completed by us.

882 A.2 ETHICS STATEMENT

883
884 We hereby confirm that the data collection, model development, and experimental methodologies
885 presented in this work adhere to the ICLR Code of Ethics. The Gran-29M dataset is constructed
886 exclusively from publicly available sources and does not contain personally identifiable information
887 or offensive content. The proposed GranViT model is designed for general-purpose visual-language
888 understanding, and we have conducted thorough analyses to identify and mitigate potential biases in
889 both training data and model outputs. All experiments are performed in accordance with responsible
890 research practices, and the model will be released for research use only to prevent potential misuse.
891 We acknowledge our responsibility to uphold ethical standards in all aspects of this research.

892 A.3 REPRODUCIBILITY STATEMENT

893
894 We have taken several steps to ensure the reproducibility of this work. All datasets used in our
895 experiments are publicly available, with their sources and processing methods detailed in Section
896 4. Model architectures and initializations are based on publicly released visual encoders and large
897 language models, as described in Section 5. All training hyperparameters, including optimizer set-
898 tings and learning rate schedules, are explicitly provided in Section 5.2. Code implementations and
899 configuration files will be made publicly available upon acceptance to further facilitate replication
900 of our results.

901 A.4 DETAILS ABOUT *Gran-29M*

902
903 We systematically document the data sources of both natural and OCR images utilized in the
904 *Gran-29M* dataset, as shown in Table 6 and Table 7. For natural images, CC3M (Sharma et al.,
905 2018), IN21k (Deng et al., 2009), SBU (Ordonez et al., 2011), CC12M (Changpinyo et al., 2021),
906 YFCC15M (Kamath et al., 2021), VisualGenome (Krishna et al., 2017), together with LAION
907 (Schuhmann et al., 2022) and FLICKR30k (Young et al., 2014), are contained. For OCR im-
908 ages, there are 30 datasets contained in total for diversity: Arxiv, InfoVQA (Mathew et al., 2022),
909 LRV-Instruction (Liu et al., 2023a), OCRVQA (Mishra et al., 2019), PDFVQA (Ding et al., 2023),
910 POIE (Kuang et al., 2023), SROIE (Huang et al., 2019), PubTables.en (Smock et al., 2022), Ren-
911 deredText (wendlerc, 2024), MMC-Instruction (Liu et al., 2023b), AI2D_gpt4v (abhayzala, 2024),
912 AI2D_internvl (Chen et al., 2024c), ArxivQA (Li et al., 2024c), Chart2Text (Kantharaj et al., 2022),
913 Diagram_Image_To_Text, Robut_SQA (Zhao et al., 2023), Robut_WikiSQL (Zhao et al., 2023),
914 Docx_en, AI2D_original (Kembhavi et al., 2016), FigureQA (Kahou et al., 2017), Hitab (Cheng
915 et al., 2021), Robut_wtq (Zhao et al., 2023), TextCaps (Sidorov et al., 2020), TextOCR, Uber_Text,
916 CORD (Park et al., 2019), ChartQA (Masry et al., 2022), DocBank (Li et al., 2020), SynthText
917 (Gupta et al., 2016) and Docmatrix (Laurençon et al., 2024). [Fig 8](#) and [Fig 9](#) visualize some data
samples in *Gran-29M*.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

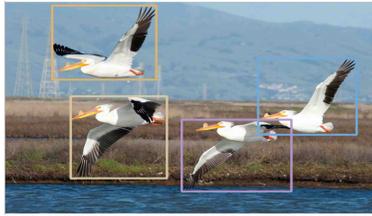
Table 6: Detailed data sources of datasets used in *Gran-29M*.

Data Type	Data Source	<i>#images</i>	<i>#regions</i>
Natural	CC3M	565521	2342622
	IN21k	614367	1628363
	LAION	17194230	54356988
	SBU	21479	52259
	CC12M	4909682	21714139
	FLICKR30k	1269	4351
	YFCC15M	655400	1884172
	VisualGenome	2150	14825
OCR	Arxiv	2655630	22574019
	InfoVQA	4343	21129
	LRV-Instruction	8304	22152
	OCRVQA	86	238
	PDFVQA	12253	98783
	POIE	898	12881
	SROIE	994	25586
	PubTables_en	121330	522516
	RenderedText	7031	20003
	MMC-Instruction	58583	218288
	AI2D_gpt4v	1958	27919
	AI2D_intervl	11981	110438
	ArxivQA	52517	1532997
	Chart2Text	22051	698874
	Diagram_Image_To_Text	154	2244
	Robut_SQA	5714	740443
	Robut_WikiSQL	38935	7092396
	Docx_en	429182	3720371
	AI2D_original	2364	27348
	FigureQA	96000	1316177
	Hitab	2495	392662
	Robut_wtq	38241	5677381
	TextCaps	20548	186867
	TextOCR	23511	215809
	Uber_Text	118042	571779
	CORD	955	3482
	ChartQA	14650	45871
	DocBank	25482	172933
SynthText	756552	3563689	
Docmatrix	1019766	51945528	

Table 7: Data sources of natural and OCR images in *Gran-29M*. *#images* and *#regions* denote the number of images and annotated bounding boxes after filtering, respectively.

Data Type	Data Source	<i>#images</i>	<i>#regions</i>
Natural	UMG-41M	6.7M	27.63M
	LAION	17.19M	54.35M
	FLICKR30k	1269	4351
OCR	Text Images	1.9M	42.57M
	Chart, Table	325K	2.8M
	Invoice Receipt	2847	41K
	Rich Text Images	3.3M	56M
TOTAL		29.51M	183.55M

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



Global Caption

The image captures a moment of flight featuring four birds. These birds, with their predominantly white plumage accentuated by black and grey markings, are mid-flight above what appears to be a body of water. Their wings are spread wide, showcasing the impressive span that these birds have. The beak details suggest they may be pelicans, known for their large size and distinctive bill shape. In the background, there's a landscape with some structures that resemble towers or transmission lines, suggesting this might be near a human settlement or infrastructure.

Local Region

[546, 141, 759, 334], Caption: The bird in the lead.
[381, 296, 618, 477], Caption: A white bird with an orange beak.
[141, 237, 349, 444], Caption: A bird flying with its wings up.
[101, 8, 326, 195], Caption: A bird flying above three other birds



Global Caption

The image captures a charmingly decorated children's bedroom. Dominating is a white crib, adorned with a bedsheet and pillowcase featuring playful cartoon characters in vibrant colors. The text overlay on the image indicates that this bedding set includes 7 pieces: a bumper sheet, pillowcase, and a duvet cover (without filling). To the right of the crib, a plush pink rabbit toy sits comfortably, adding to the room's child-friendly ambiance. On the opposite side, a red cup is visible, perhaps indicating for an adult or older sibling who might be taking care of the young child. The overall setting suggests careful consideration of color coordination and character themes, creating a welcoming environment for a child.

Local Region

[558, 268, 679, 558], Caption: A nightstand with a white top and red bottom.
[287, 131, 445, 285], Caption: A pink pillow on the bed.
[575, 125, 679, 304], Caption: A pink stuffed animal with a bow on its head.
[582, 312, 679, 389], Caption: The drawer is white.

7 pieces: bumper sheet + pillowcase + duvet cover (without duvet filling)

Figure 6: Visualization of *Gran-29M*.

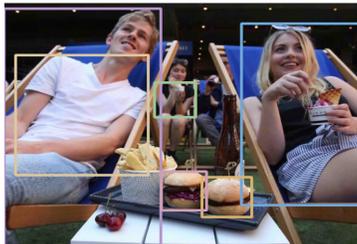


Global Caption

The image captures a white Chevrolet SUV parked in front of a building. The sign on the wall behind it says "Hablamos Español Mal Comprendemos" which is Spanish for "Let's speak Spanish Well Understand Us". This indicates that the business may be bilingual or catering to Spanish-speaking customers.

Local Region

[106, 91, 553, 399], Caption: A white chevrolet truck.
[488, 195, 588, 231], Caption: A black pickup truck.
[230, 271, 307, 401], Caption: A black tire on a white truck.
[0, 1, 285, 273], Caption: A white building with a sign that says 'sund'.
[496, 102, 640, 225], Caption: A white house in the background.
[106, 74, 249, 161], Caption: Green and red sign.
[362, 132, 511, 194], Caption: A white truck parked in.



Global Caption

The image captures a casual and relaxed outdoor gathering. Two individuals are seated in deck chairs, with the man on the left and the woman on the right. The man is wearing a white t-shirt and jeans, while the woman is dressed in a black top and striped shorts. In front of them, there's a tray holding two burgers, one with purple toppings that may be pickles or onion relish, and another with red condiment that could be ketchup or tomato sauce. Accompanying the burgers, there's a basket filled with breadsticks and a single-use bottle of what appears to be beer. The background shows other more people sitting, suggesting a social event.

Local Region

[415, 39, 629, 390], Caption: A woman in a black shirt.
[0, 10, 274, 470], Caption: A man in a white shirt.
[20, 109, 252, 321], Caption: The shirt is white.
[346, 336, 433, 413], Caption: A small sandwich with purple filling.
[263, 153, 335, 220], Caption: A person sitting in the background.
[275, 325, 354, 400], Caption: The sandwich has purple sauce

Figure 7: Visualization of *Gran-29M*.

A.5 HIGH RESOLUTION PERFORMANCE OF GRANViT

Additionally, we provide the evaluation results of GranViT with an image tiling strategy in the pretraining. According to image tiling (Chen et al., 2024c; Gu et al., 2024; Wang et al., 2024; Lu et al., 2024; Team et al., 2025b), images are firstly converted into $N \times 512 \times 512$ local patches and one global patch. All patches are simultaneously fed into the vision encoder for feature extraction. With a patch size of 16, we obtain $N \times 1024$ visual patch features. We use pixel shuffle (Gu et al., 2024; Wang et al., 2024) to compress these visual features to $N \times 256$ patches. These visual features are then recombined based on positions and fed into the projector and LLM for understanding. The evaluation results are reported in Table 8.

A.6 THE REASON FOR FREEZING THE VISION ENCODER IN STAGE 2

In Table 9, we compare the performance gap between the vision encoder that is frozen and tunable in stage 2. Tuning the vision encoder in stage 2 does not achieve significant improvement, while leading to more training cost, since the vision encoder is trained well in stage 1 for fine-grained

Table 8: Performance comparison with image tiling. The bold font represents the best performance, and the underline represents the second performance.

Capability	Benchmark	CLIP	SigLip	SigLip2	AIMv2	SAILViT	GranViT
Fine-Grained	RefCOCO_testA	82.03	82.58	84.95	84.58	<u>87.92</u>	90.71
	RefCOCO_testB	66.53	67.47	74.91	70.02	<u>76.85</u>	82.04
	RefCOCO_val	76.01	77.25	80.11	77.97	<u>83.12</u>	87.21
	RefCOCO+_testA	75.16	77.68	80.24	79.23	<u>82.48</u>	85.15
	RefCOCO+_testB	32.00	55.67	63.81	59.97	<u>67.13</u>	70.52
	RefCOCO+_val	65.74	68.32	72.26	70.71	<u>75.59</u>	79.05
	RefCOCOg_val	70.73	72.74	75.23	74.91	<u>78.92</u>	81.98
	RefCOCOg_test	70.81	72.18	74.75	74.48	<u>77.98</u>	81.63
	BLINK*	51.40	52.84	<u>53.49</u>	56.25	52.33	52.42
	MMVP	61.33	64.00	65.33	<u>67.33</u>	68.00	65.66
	Average		67.31	69.07	72.51	71.55	<u>75.03</u>
VQA	MMBench	60.44	64.39	<u>64.62</u>	65.01	64.00	63.77
	MMStar	39.06	41.26	<u>42.73</u>	41.80	45.60	40.80
	HallusionBench	30.18	29.89	30.08	26.95	<u>30.56</u>	35.93
	GQA	58.99	59.88	59.95	60.04	<u>60.72</u>	61.32
	SEEDBench	67.35	68.75	<u>69.60</u>	69.52	70.07	69.54
	Average		51.20	52.83	53.40	52.66	<u>54.19</u>
Reasoning	MMMU	38.11	35.44	41.33	<u>40.66</u>	40.44	35.66
	MathVista MINI	36.80	35.80	37.40	<u>38.90</u>	38.50	40.00
	MMVet	35.27	32.47	40.59	40.59	<u>40.13</u>	38.30
	ScienceQA	66.18	68.66	68.36	<u>66.98</u>	74.71	66.28
	AI2D	67.13	68.65	69.62	69.52	71.85	<u>69.88</u>
	Average		48.70	48.20	<u>51.46</u>	51.33	53.13
OCR	OCRBench	414	450	545	522	<u>551</u>	583
	DocVQA	53.16	58.36	68.74	64.58	<u>71.81</u>	72.81
	ChartQA	60.12	62.60	65.04	65.48	<u>67.36</u>	71.96
	InfoVQA	24.17	27.81	31.71	31.38	<u>33.47</u>	33.59
	TextVQA	56.74	60.95	67.78	66.33	69.47	<u>69.40</u>
	Average		47.12	50.94	57.55	55.99	<u>59.44</u>

Table 9: Performance comparison of whether the vision encoder is frozen in stage 2.

Vision Encoder State	FLOPs	MACs	Fine-Grained	VQA	Reasoning	OCR
Frozen	3.24T	1.62T	77.24	54.83	51.34	54.02
Tunable	4.25T	2.09T	77.17	54.83	52.48	54.06

feature extraction. Therefore, to reduce the training complexity, we freeze the vision encoder in stage 2.

A.7 THE DIFFERENCE BETWEEN GRANViT AND SAILViT

SAILViT (Yin et al., 2025) addresses the problem of insufficient visual-language alignment through a three-stage pre-training strategy that co-optimizes the vision encoder, projector, and LLM to achieve better alignment. GranViT differs from SAILViT in two key aspects. First, GranViT leverages both *Bbox2Caption* and *Caption2Bbox* tasks to strengthen fine-grained feature extraction and local region localization in the vision encoder and LLM, while SAILViT relies solely on global

Table 10: Performance comparison when the vision encoder is frozen during SFT training.

Vision Encoders	Fine-Grained	VQA	Reasoning	OCR
SigLip2	70.51	53.36	51.41	49.16
AIMv2	57.31	52.24	48.75	46.90
SAILViT	71.90	54.16	50.93	52.79
GranViT	75.16	53.07	49.12	54.81

Table 11: Continue training performance of GranViT. Ref, Ref+ and Refg denotes the Ref-COCO_testA, RefCOCO+_testA and RefCOCOg_test respectively. MMB, HB, and SB denote the MMBench, HallusionBench, and SEEDBench. SQA, OB, DVQA, and IVQA denote ScienceQA, OCRBench, DocVQA, and InfoVQA, respectively.

Model	Ref	Ref+	Refg	MMB	HB	SB	MMMU	SQA	OB	DVQA	IVQA	Avg
SigLip2	88.43	83.26	79.54	62.61	32.54	69.77	38.66	71.93	584	60.56	26.44	61.10
GranViT	91.26	86.67	84.32	61.37	32.44	70.48	36.77	70.59	623	68.16	28.61	62.99

question-answering in the pretraining. Second, while SAILViT injects world knowledge into the vision encoder using large-scale SFT data to improve task-specific performance, GranViT focuses on enhancing the generic representation ability of the vision encoder. We contend that improved generic representations facilitate better adaptation to downstream tasks. Notably, the two strategies are complementary: after learning stronger generic features, task-specific adaptation following the paradigm of SAILViT can further enhance the performance of MLLM on specialized applications.

A.8 CONTINUE PRETRAINING

We augment our two-stage training paradigm via continuous pretraining with SFT data, following stage 3 in SAILViT (Yin et al., 2025). We leverage LLaVA-One-Vison (Li et al., 2024a) dataset for pretraining and utilize Open-LLaVA-NeXT 1M dataset (Chen & Xing, 2024) for SFT. GranViT is initialized from SigLip2 and pretrained with *Gran-29M* dataset (8M samples for both stages) first. Then, we apply continuous pretraining with SFT data to both models. As shown in Table 11, GranViT outperforms SigLip2 on fine-grained and OCR evaluation significantly. This experiment demonstrates that our method is compatible with the pretraining approach of SAILViT and that SFT data can be subsequently incorporated after our pretraining paradigm to further performance improvement.

A.9 FROZEN VISION ENCODER IN SFT

To isolate the training gains of the vision encoder during the SFT stage, we compared the performance of different vision encoders with a frozen MLLM. As shown in Table 10, GranViT achieved the best performance in both fine-grained perception and OCR evaluations, outperforming SailViT by an average of 3.2 and 2.1, respectively. This fully demonstrates the significantly stronger fine-grained perception capability of GranViT.

A.10 ATTENTION MAP VISUALIZATION

Similar to Fig. 1(b), we provide additional visualizations of attention maps in complex multi-object and OCR text scenarios in Fig. 8 and Fig. 9. Furthermore, to enhance the clarity of the attention maps, we filter out pixels with attention values below a threshold of 0.3, which helps eliminate diffuse activations and accentuate the model’s primary focus areas. The results demonstrate that GranViT tends to concentrate more on local regions, while SailViT, SigLIP2, and AIMv2 exhibit a stronger focus on global areas.

A.11 *Bbox2Caption* VISUALIZATION

In Fig. 10, we provide several visualizations for the *Bbox2Caption* task and compare the answers between GranViT and SigLip2. The red rectangle box visualizes the bbox coordinates mentioned in the questions. Unlike SigLIP2, which is limited to global image-level captioning, GranViT supports fine-grained description generation for objects inside specified bboxes.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

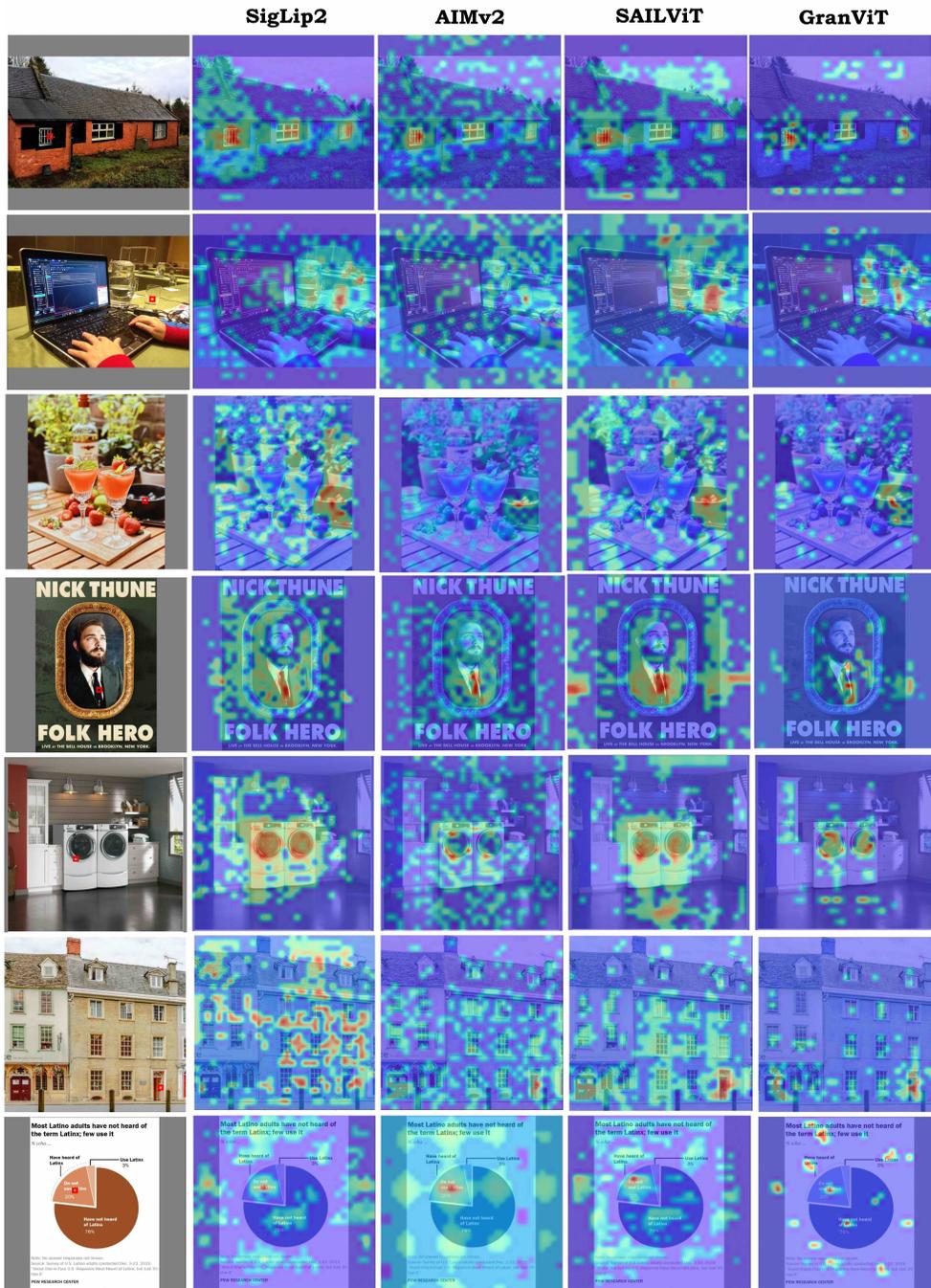


Figure 8: Attention map visualization. The red rectangle box denotes the query region for attention maps.

A.12 *Caption2Bbox* VISUALIZATION

Fig. 11 and Fig. 12 visualizes the bbox prediction results of the *Caption2Bbox* task in two-stage training. Green bboxes indicate predicted regions, while red ones denote ground truth. It is evident that training the *Caption2Bbox* task in Stage 1 yields inferior performance compared to training it in Stage 2.



1232 Figure 9: Attention map visualization. The red rectangle box denotes the query region for attention
1233 maps.

1235 A.13 BAD CASE VISUALIZATION

1236
1237 we visualize some failure cases in Fig. 13 in the revised manuscript. GranViT still exhibits sub-
1238 optimal performance in three challenging scenarios (from top to bottom): 1) when the bounding
1239 boxes are excessively small or text is overly dense, 2) when objects severely overlap, and 3) when
1240 precise localization is hindered by its relative coordinates. To address these limitations, we plan to
1241 explore several directions in future work. First, we will investigate the use of absolute coordinates
(similar to Qwen3-VL) to improve spatial precision. Second, we will develop multi-scale pretrain-

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



Q: Describe the content contained within the normalized bounding box coordinates [0.366, 0.627, 0.575, 0.873] in no more than 10 words.

SigLip2: The 2019 KAWAII Mule is a utility vehicle with an engine displacement of 450cc.

SAILViT : A green off-road vehicle

GranViT : The tire is black

Ground Truth: The tire is black



Q: Describe the content contained within the normalized bounding box coordinates [0.19, 0.756, 0.417, 0.888] in no more than 10 words.

SigLip2: A kitchen with a black vent hood above the stove and pots on the counter

SAILViT : This kitchen scene suggests a space ready for cooking, featuring a prominent black vent hood over the stove. Several pots sit on the adjacent counter, adding a touch of everyday life and practicality to the space.

GranViT : A white pot on the stove

Ground Truth: A white bowl on a stove



Q: Describe the content contained within the normalized bounding box coordinates [0.138, 0.208, 0.326, 0.38] in no more than 10 words.

SigLip2: A table with a marble top and pink glasses, forks, knife, napkins, and macarons on it.

SAILViT : An elegantly set table with a sleek marble top is prepared for a gathering. The setting includes pink-tinted drinking glasses, polished silverware, neatly folded napkins, and a delightful selection of macarons, awaiting guests.

GranViT : A glass of water

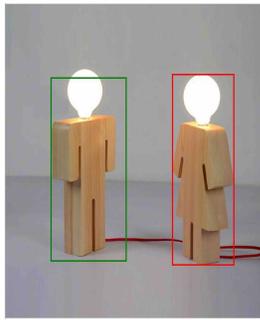
Ground Truth: A glass of water with a white label

Figure 10: *Bbox2Caption* visualization. The red rectangle box visualizes the bbox coordinates mentioned in the questions.

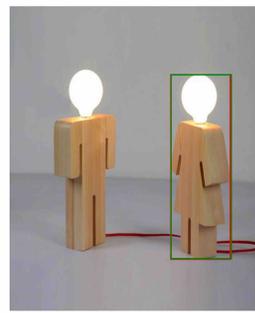
ing strategies to enhance performance on extremely small objects. Third, we plan to incorporate

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Q: Please provide the bounding box coordinate of the region this sentence describes: ['a wooden man with a female partner']



Train caption2bbox in stage1



Train caption2bbox in stage2

Q: Please provide the bounding box coordinate of the region this sentence describes: ['a wicker basket with a pillow in']



Train caption2bbox in stage1



Train caption2bbox in stage2

Q: Please provide the bounding box coordinate of the region this sentence describes: Mean age(yrs)'

	Group I	Group II	Statistic	P value	NS
Number of patients	24	24			
Male	15	11	$X^2 = 1.343$	0.247	*
Female	9	13			
Mean age(yrs)	52.31 (range:22-83)	54.55 (range:25-78)	$t = 0.448$	0.657	*
Operation side			$X^2 = 2.116$	0.146	*
Left	11	16			
Right	13	8			
Aetiology of indications			$X^2 = 1.532$	0.655	*
Femoral head necrosis	16	18			
Femoral neck fracture	3	4			
Primary osteoarthritis	5	2			

Train caption2bbox in stage1

	Group I	Group II	Statistic	P value	NS
Number of patients	24	24			
Male	15	11	$X^2 = 1.343$	0.247	*
Female	9	13			
Mean age(yrs)	52.31 (range:22-83)	54.55 (range:25-78)	$t = 0.448$	0.657	*
Operation side			$X^2 = 2.116$	0.146	*
Left	11	16			
Right	13	8			
Aetiology of indications			$X^2 = 1.532$	0.655	*
Femoral head necrosis	16	18			
Femoral neck fracture	3	4			
Primary osteoarthritis	5	2			

Train caption2bbox in stage2

Figure 11: *Caption2Bbox* visualization. Green bboxes indicate predicted regions, while red ones denote ground truth.

advanced data augmentation techniques specifically designed for dense and overlapping scenarios. These improvements will help build a more robust and accurate visual grounding system.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
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

Q: Please provide the bounding box coordinate of the region this sentence describes: ['flowers on the table']



Train caption2bbox in stage1



Train caption2bbox in stage2

Q: Please provide the bounding box coordinate of the region this sentence describes: ['a picture on the wall']



Train caption2bbox in stage1



Train caption2bbox in stage2

Q: Please provide the bounding box coordinate of the region this sentence describes: ['the shirt is gray']



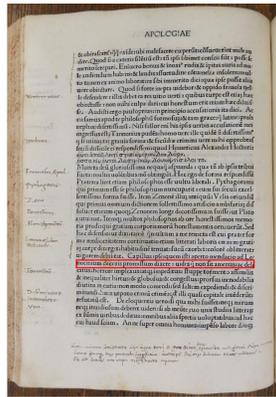
Train caption2bbox in stage1



Train caption2bbox in stage2

Figure 12: *Caption2Bbox* visualization. Green boxes indicate predicted regions, while red ones denote ground truth.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457



Q: Describe the content contained within the normalized bounding box coordinates [0.334, 0.651, 0.762, 0.673] in no more than 10 words.
GranViT: &inequaliturhuis globosacconegtusdelfilpepanduldeiclini
Ground Truth: nocinium decoris promffum dixcre:uidesg non fit amenusac deli



Q: Describe the content contained within the normalized bounding box coordinates [0.288, 0.171, 0.545, 0.509] in no more than 10 words.
GranViT: a man in a suit holding a microphone
Ground Truth: a sign in spanish



Q: Describe the content contained within the normalized bounding box coordinates [0.41, 0.363, 0.551, 0.467] in no more than 10 words.
GranViT: the nose of a man
Ground Truth: the eye of a man

Figure 13: Bad case visualization of *bbox2caption* tasks. The red rectangle box visualizes the *bbox* coordinates mentioned in the questions. We visualize three cases (from top to bottom): 1) when the bounding boxes are excessively small or text is overly dense, 2) when objects severely overlap, and 3) when precise localization is hindered by its relative coordinates.