Non-Natural Image Understanding with Advancing Frequency-based Vision Encoders

Supplementary Material

001 002	Overview In this supplementary material, we present following content.	the
003		
004	1. Additional Examples	1
005	2. Evaluation Details	1
006	3. Training Data	1
007	4. Baselines for the Caption task	1
800	4.1. Closed-Source MLLMs	1
009	4.2. Open-Source MLLMs	2
010	5. Baselines for the Q&A task	2
011	5.1. For GeoQA	2
012	5.2. For FunctionQA	3
013	5.3. For ChartQA	3
014	6. Ablation of Vision Encoder in Stage 3	3
015	7. Limitation and Future Work	5

1. Additional Examples

In Figure 1 and Figure 2, we present additional qualitative
examples illustrating the performance of our model on both
captioning and question-answering tasks. These figures underscore the versatility of our model across multiple types
of non-natural image understanding tasks.

022 2. Evaluation Details

Traditional metrics for evaluating the quality of descrip-023 tions, such as BLEU [17] and CIDEr [20], are limited by 024 their reliance on the linguistic style of the reference descrip-025 026 tions. For instance, even for the GPT-4V, if not fine-tuned on a corresponding training set, may generate descriptions 027 028 that are correct but receive low scores. Conversely, a model fine-tuned on a training set with a similar linguistic style 029 030 can easily achieve high scores. To overcome this limitation, we are inspired by [14], and we consider using GPT-4 to 031 evaluate the correctness and level of detail in the generated 032 descriptions. The prompts used for geometric caption eval-033 uation are shown in Figure 3 and Figure 4(similar prompts 034 for chart and function), which allows for a more objective 035 036 assessment of the results produced by different models.

<u> </u>		N T 1
Category	Statistic	Number
Total	- Total number	190K
	- Average length (words)	61.91
	- Average length (characters)	350.67
	- Vocabulary size	21127
Geometry	- Total number	60K
	- Average length (words)	44.23
	- Average length (characters)	230.28
	- Vocabulary size	2410
Chart	- Total number	30K
	- Average length (words)	109.40
	- Average length (characters)	726.73
	- Vocabulary size	20598
Function	- Total number	100K
	- Average length (words)	58.49
	- Average length (characters)	321.11
	- Vocabulary size	145

Table 1. Statistics of Caption Training Data.

3. Training Data

In Tables 1 and Tables 2, we present the training data used 038 for alignment and instruction fine-tuning across three types 039 of non-natural images: geometry, charts, and functions. 040 The geometry data is sourced from Geo170K [4], the func-041 tion data from Mavis-caption [22], and the chart data from 042 Chartllama [6]. Notably, for the chart data, we extended the 043 textual responses to provide more detailed analyses of the 044 visual modality's impact on the results, rather than directly 045 answering with a single number or word. All of our training 046 data will be publicly available. 047

4. Baselines for the Caption task

To evaluate the performance of our model, we compare it
against both closed-source and open-source multimodal049large language models (MLLMs).050vided into two categories:051

4.1. Closed-Source MLLMs

These closed-Source models represent the state-of-the-art054in multimodal understanding and captioning tasks. They055include: Qwen-VL-Plus [2], Gemini-1.0-Pro [19], Qwen-056VL-Max [2], and GPT-4V [16]. These models serve as the057performance upper bound, given their advanced design and058access to extensive training resources.059

048

053

CVPR 2025 Submission #10821. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 1. More presentation of caption tasks.

Table 2. Statistics of Instruction Training Data.

Category	Statistic	Number
Total	- Number of unique questions	247K
	- Average question length	44.35
	- Average answer length	61.62
Geometry	- Number of unique questions	117K
	- Average question length	69.88
	- Average answer length	70.19
Chart	- Number of unique questions	100K
	- Average question length	12.51
	- Average answer length	60.51
Function	- Number of unique questions	30K
	- Average question length	50.82
	- Average answer length	31.45

4.2. Open-Source MLLMs

These models are publicly available, making them more ac-061 cessible for research and development. All three data types 062 include LLaVA [13] and its variants (LLaVA-adapter [5], 063 LLaVA-next [8]), known for their adaptability to multi-064 modal inputs, with visual encoders fine-tuned on open 065 datasets. For these models, we directly use their pre-trained 066 weights. In addition, we introduce 3 baselines for each data 067 type, including: 068

Geo-Caption: For geometric image caption we selected
3 state-of-the-art baselines respectively: G-LLaVA [4],
MAVIS [22], and EAGLE [9].

Function Caption: For function image caption, to the
best of our knowledge, only MAVIS has attempted, so we
use the function image caption dataset provided by MAVIS
and train according to G-llava and EAGLE as the corre-

sponding baseline model.

- G-LLaVA: First train the projection linear layer and then both the projection linear layer and the language model are trainable.
 077
 078
 079
- **MAVIS**: First Train the visual encoder and then freeze the CLIP-Math in both the alignment phase and instruction fine-tuning phase, and train the projection layer along with the LoRA-based LLM.
- **EAGLE**: First train the visual encoder and projection layer and then use Lora to fine-tune the visual encoder and train the projection layer and LLM with full parameters

Chart Caption: For the chart image caption, we compare 2 chart understanding models, ChartLlama [6] and Unichart [15]. We additionally use the same chart dataset to fine-tune MAVIS, G-LLAVA, and EAGLE. The complete experimental results are shown in the **Table 3**.

5. Baselines for the Q&A task

5.1. For GeoQA

In our experimental setup, we compare our method against 095 three categories of baselines: heuristics-based baselines, 096 conventional models, and multi-modal large language mod-097 els (MLLMs). The heuristics-based baselines include Ran-098 dom Chance and Frequent, which provide simple reference 099 points for evaluating performance. Conventional models 100 such as Geoformer [3] and UniMath [10] represent tra-101 ditional task-specific solutions designed for mathematical 102 and geometric reasoning. The MLLMs include state-of-103 the-art open-source models like LLAVA-1.5 [13] (7B and 104 13B), Math-LLAVA [18], G-LLAVA [4] (7B and 13B), 105 MAVIS [22], and EAGLE [9], which are based on ad-106

076

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094



Figure 2. More presentation of question-answer tasks.

Table 3. The complete experimental results of Chart Caption

Model	Correctness	Detail
Closed-source MLLMs		
Qwen-VL-Plus [2]	1.75	2.15
Gemini-1.0-Pro [19]	2.00	2.20
Qwen-VL-Max [2]	2.35	2.30
GPT-4V [16]	1.88	2.33
Open-source MLLMs		
LLAVA [13]	2.05	2.45
LLAMA-Adapter [5]	2.19	2.65
LLAVA-NeXT [8]	2.26	2.72
G-LLAVA [4]	2.25	2.79
MAVIS [22]	2.31	2.78
EAGLE [9]	2.33	2.75
Ours	2.35	2.84

vanced architectures such as Vicuna and LLAMA-2. These
baselines ensure a comprehensive evaluation of our method
against a range of existing approaches.

110 5.2. For FunctionQA

In our experimental setup, we compare EDGE against three 111 categories of baselines: heuristics-based baselines, closed-112 source MLLMs, and open-source MLLMs. The heuristics-113 based baseline includes Random Chance, which serves as 114 a simple reference point. The closed-source MLLMs in-115 clude advanced proprietary models such as CoT GPT-4 [1], 116 117 PoT GPT-4 [1], Multimodal Bard [19], and GPT-4V [16], which represent the state-of-the-art in multimodal reason-118 ing tasks. The open-source MLLMs consist of models 119 like LLAVA [21], LLAMA-Adapter [5], LLAVA-NeXT [8], 120 SPHINX-MoE [11], and MAVIS [22], which leverage ac-121 cessible architectures such as Vicuna-1.5 and MAmmoTH-122 123 2-7B for multimodal tasks. These baselines provide a comTable 4. Ablation experiments of the visual encoder with different adjustment strategies in stage 3.

Vision Encoder in Stage 3	Full Fine-tune	Freeze	LoRA
Performance on GeoQA	65.7	67.3	68.2

prehensive comparison, highlighting the effectiveness of
our model, which achieves competitive performance against124both open-source and closed-source approaches.125

5.3. For ChartQA

In our experiments on the ChartQA benchmark, we com-128 pare our method against both conventional models and 129 multi-modal large language models (MLLMs). The conven-130 tional models include Pix2struct [7], Matcha [12], Chart-131 T5 [23], and Unichart [15], which are specialized meth-132 ods designed for processing and reasoning with chart 133 data. The multi-modal large language models consist of 134 SPHINX [11], Qwen [2], and Chartllama [6]. 135

6. Ablation of Vision Encoder in Stage 3

The visual encoder is fully fine-tuned in stages 1 and 2. In 137 stage 3, visual encoder is frozen and LoRA is used for fur-138 ther adjustments. This approach is based on the fact that, as 139 the LLM backbone is fine-tuned on the QA dataset and ac-140 quires knowledge from step-by-step rationales, the LoRA-141 based visual encoder progressively shifts focus toward key 142 geometric cues critical to the resolution process, rather than 143 the entire image. This synergy enhances a deeper under-144 standing of geometric features. We include an ablation 145 study of the adjustment strategy, as shown in the Table 146 above, which confirms that LoRA tuning yields the best per-147 formance. The projection layer is fine-tuned in stages 2 and 148 3, while the LLM is adjusted only in stage 3. 149

127

136

// Evaluation of the correctness of the geometry

{

"role": "system",

"content":

"You are an intelligent chatbot designed for evaluating the accuracy of generative outputs for geometry-based image questionanswer pairs."

"Your task is to compare the predicted answer with the correct answer and determine its accuracy in terms of geometric shapes and properties. Here's how you can accomplish the task:"

"_____"

"##INSTRUCTIONS: "

"- Check if the predicted answer includes all the correct geometric shapes and properties present in the correct answer.\n"

"- Evaluate whether the predicted answer accurately describes the geometric elements, including their relationships and characteristics.\n"

"- Consider synonyms or paraphrases as valid matches.\n"

"- Provide a single evaluation score that reflects the accuracy of the prediction in terms of geometric correctness."

},

{

"role": "user",

"content":

"Please evaluate the following geometry-based image question-answer pair:\n\n"

f"Question: {question}\n"

f"Correct Answer: {answer}\n"

f"Predicted Answer: pred n n"

"Provide your evaluation only as an accuracy score where the accuracy score is an integer value between 0 and 5, with 5 indicating the highest level of accuracy."

"Please generate the response in the form of a Python dictionary string with keys 'score', where its value is the accuracy score in INTEGER, not STRING."

"DO NOT PROVIDE ANY OTHER OUTPUT TEXT OR EXPLANATION. Only provide the Python dictionary string."

"For example, your response should look like this: {'score': 4}."

}

Figure 3. Evaluation of the correctness of the geometry.

//Evaluation of geometric details

"role": "system", "content":

"You are an intelligent chatbot designed for evaluating the detail orientation of generative outputs for geometry-based image question-answer pairs." "Your task is to compare the predicted answer with the correct answer and determine its level of detail, considering both completeness and specificity. Here's how you can accomplish the task:"

"##INSTRUCTIONS: "

"- Check if the predicted answer covers all major points from the image. The response should not leave out any key aspects related to the geometric elements.\n"

"- Evaluate whether the predicted answer includes specific details rather than just generic points. It should provide comprehensive information that is tied to specific elements of the image.\n"

"- Consider synonyms or paraphrases as valid matches.\n"

"- Provide a single evaluation score that reflects the level of detail orientation of the prediction, considering both completeness and specificity."

}, {

"role": "user",

"content":

"Please evaluate the following geometry-based image question-answer pair:\n\n"

f"Question: {question}\n"

f"Correct Answer: {answer}\n"

f"Predicted Answer: {pred} $\n\n$ "

"Provide your evaluation only as a detail orientation score where the detail orientation score is an integer value between 0 and 5, with 5 indicating the highest level of detail orientation."

"Please generate the response in the form of a Python dictionary string with keys 'score', where its value is the detail orientation score in INTEGER, not STRING."

"DO NOT PROVIDE ANY OTHER OUTPUT TEXT OR EXPLANATION. Only provide the Python dictionary string."

"For example, your response should look like this: {'score': 4}."

}

Figure 4. Evaluation of geometric details.

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225 226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255 256

257

150 7. Limitation and Future Work

Despite the significant progress achieved by FM-ViT and 151 152 the edge model in understanding non-natural images, there are still some limitations that warrant further investigation. 153 154 First, while the frequency modulation technique enhances the ability to capture high-frequency information, it remains 155 156 insufficient in extracting fine-grained details, such as subtle protrusions on a line or isolated points. Second, this study 157 primarily focuses on improving the visual encoding capabil-158 159 ities for non-natural images, without delving into the impact of larger and more diverse datasets, leaving this aspect for 160 future exploration. 161

In future work, we plan to further enhance the model's 162 163 ability to extract fine-grained information from non-natural images. Additionally, we aim to generate synthetic datasets 164 with greater diversity and realism through advanced data 165 augmentation methods to further improve the model's un-166 derstanding and generalization capabilities for broader ap-167 plications, such as education, accessibility tools, and data-168 169 driven policymaking.

170 References

- In Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al.
 Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 3
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 1(2):3, 2023. 1, 3
- [3] Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. *arXiv preprint arXiv:2212.02746*, 2022. 2
- [4] Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal large language model. *arXiv preprint arXiv:2312.11370*, 2023. 1, 2, 3
- [5] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie
 Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023. 2, 3
- [6] Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang,
 Gang Yu, Bin Fu, and Hanwang Zhang. Chartllama: A multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*, 2023. 1, 2, 3
- [7] Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu,
 Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova.
 Pix2struct: Screenshot parsing as pretraining for visual lan-

guage understanding. In *International Conference on Machine Learning*, pages 18893–18912. PMLR, 2023. 3

- [8] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild, 2024. 2, 3
- [9] Zhihao Li, Yao Du, Yang Liu, Yan Zhang, Yufang Liu, Mengdi Zhang, and Xunliang Cai. Eagle: Elevating geometric reasoning through llm-empowered visual instruction tuning. arXiv preprint arXiv:2408.11397, 2024. 2, 3
- [10] Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. Unimath: A foundational and multimodal mathematical reasoner. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7126–7133, 2023. 2
- [11] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. arXiv preprint arXiv:2311.07575, 2023. 3
- [12] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. arXiv preprint arXiv:2212.09662, 2022. 3
- [13] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 2, 3
- [14] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. arXiv preprint arXiv:2306.05424, 2023. 1
- [15] Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty. Unichart: A universal visionlanguage pretrained model for chart comprehension and reasoning. arXiv preprint arXiv:2305.14761, 2023. 2, 3
- [16] OpenAI. Openai: Gpt-4v(ision) system card. https:// openai.com/research/gpt-4v-system-card, 2023. 1, 3
- [17] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 1
- [18] Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. Mathllava: Bootstrapping mathematical reasoning for multimodal large language models. arXiv preprint arXiv:2406.17294, 2024. 2
- [19] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1, 3
- [20] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evalua 259

260

261

tion. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 4566–4575, 2015. 1

- [21] Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, Maosong Sun, and Gao Huang. Llava-uhd: an lmm perceiving any aspect ratio and high-resolution images. *arXiv preprint arXiv:2403.11703*, 2024. 3
- [22] Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Yichi Zhang, Ziyu
 Guo, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Bin Wei,
 Shanghang Zhang, et al. Mavis: Mathematical visual instruction tuning. *arXiv preprint arXiv:2407.08739*, 2024. 1,
 2, 3
- [23] Mingyang Zhou, Yi R Fung, Long Chen, Christopher
 Thomas, Heng Ji, and Shih-Fu Chang. Enhanced chart
 understanding in vision and language task via crossmodal pre-training on plot table pairs. *arXiv preprint arXiv:2305.18641*, 2023. 3