EUCLID: SUPERCHARGING MULTIMODAL LLMS WITH SYNTHETIC HIGH-FIDELITY VISUAL DESCRIPTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) have made rapid progress in recent years, yet continue to struggle with *low-level visual perception* (LLVP) particularly the ability to accurately describe the geometric details of an image. In this paper, we first demonstrate this limitation by introducing *Geoperception*, a benchmark designed to evaluate an MLLM's ability to accurately transcribe 2D geometric information from an image. We then conduct a comprehensive empirical study to explore strategies for improving LLVP performance through the use of synthetic high-fidelity visual description data. Our findings highlight the benefits of certain model architectures and training techniques, including the use of CNN-based visual encoders and multi-stage training with a data curriculum. Notably, we find that a data curriculum enables models to learn challenging geometry understanding tasks which they fail to learn from scratch. Lastly, we develop *Euclid*, a family of models specifically optimized for strong low-level geometric perception. Although trained on synthetic multimodal data, Euclid shows strong generalization ability on novel real-world geometry shapes. For instance, Euclid outperforms the best closed-source model in our benchmark by up to 58.56% on certain Geoperception benchmark tasks and 10.65% on average across all tasks.

028 1 INTRODUCTION

029

004

010

011

012

013

014

015

016

017

018

019

021

025

026 027

Multimodal large language models (MLLMs) have rapidly progressed in recent years, demonstrating remarkable potential in understanding and reasoning about the visual world through the powerful capabilities of large language models (LLMs) (Liu et al., 2024c;a; Achiam et al., 2023; Team et al., 2023; Hu et al., 2023; Tong et al., 2024a; Wang et al., 2024a). These models have showcased strong performance in tasks such as visual question answering (VQA) (Goyal et al., 2017), image captioning (Lin et al., 2014), and multimodal reasoning (Liu et al., 2023) – for example, LLaVA-NeXT-34B (Liu et al., 2024b) achieves an impressive 83.7% accuracy on the VQAv2 benchmark (Goyal et al., 2017), a comprehensive natural image VQA benchmark.

While MLLMs achieve impressive results on tasks like VQA, their performance relies on high-level semantic extraction (Tong et al., 2024b); in contrast, they often fall short on *low-level visual perception* (LLVP) — i.e., the ability to accurately describe the geometric details of an image, such as 040 the points, lines, angles, and spatial relationships among its constituent objects. This limitation be-041 comes especially apparent in tasks requiring precise visual descriptions, such as mathematical visual 042 problem solving (Zhang et al., 2024a; Lu et al., 2023), scientific visual understanding (Yue et al., 043 2024; Fu et al., 2024a), abstract visual reasoning (Jiang et al., 2024; Ahrabian et al., 2024), and even 044 simple visual comprehension (Rahmanzadehgervi et al., 2024; Wang et al., 2024b). Furthermore, LLVP is also vital in real-world applications, including spatial understanding for robotics, medical image analysis for accurate diagnosis, GUI agents, quality control in manufacturing to detect 046 subtle defects, autonomous driving systems that rely on exact object localization, and augmented 047 reality applications that demand precise overlay of virtual objects onto the real world. In this paper, 048 we study the challenges of LLVP abilities in MLLMs, take steps to understand the root cause, and improve their performance. Our study focuses specifically on 2D geometry, a domain where components and relationships are rigorously defined, and both images and textual descriptions can be 051 generated at scale synthetically. 052

053 We start by developing *Geoperception*, a 2D geometry multimodal benchmark that focuses exclusively on basic LLVP questions without incorporating higher-order reasoning. Our findings reveal



Figure 1: Four examples from our *Geoperception* dataset. The questions are sourced from the Geometry-3K corpus (Lu et al., 2021), which compiles problems from two widely-used high school textbooks. We perform filtering, validation, and generate question-and-answer text for each image.

that current MLLMs consistently struggle on this benchmark. This limitation raises an important re-066 search question: what is causing the difficulties exhibited by contemporary MLLMs in LLVP tasks? 067 We hypothesize the main factor to be the insufficient availability of high-fidelity visual descrip-068 tion datasets. Furthermore, in the absence of sufficient data, it is challenging to identify and adopt 069 architectural choices and training strategies that could improve the effectiveness and efficiency of LLVP-specific training. To this end, we develop a synthetic dataset engine for large-scale generation 071 of 2D geometry images paired with high-fidelity textual descriptions. This synthetic dataset enables us to conduct a comprehensive empirical study with controlled experiments to explore strategies for 073 improving MLLMs' performance on LLVP tasks within the 2D geometry domain. Our key insights 074 include: 1. Scaling LLM size does not benefit LLVP learning. 2. CNN visual encoders are more 075 suitable in LLVP learning than ViT architectures. 3. Tuning vision encoders does not offer a strong advantage. 4. Curriculum learning significantly improves a model's performance, es-076 pecially in understanding complex geometric shapes. With these lessons learned, we then train 077 a family of models—using a carefully designed curriculum of synthetic data—that are specifically optimized for strong LLVP, which we call Euclid. Upon evaluation, we show that our models excel 079 on real-world low-level geometric perception tasks.

081

063

064

065

- 082
- 083 084

087

090

GEOPERCEPTION: A BENCHMARK FOR GEOMETRIC LLVP 2

Although the shortcomings of MLLMs in LLVP are commonly recognized, there is no comprehen-085 sive benchmark that focuses purely on this task. Our goal is to construct a benchmark focusing solely on the LLVP ability of MLLMs, which is also representative enough of real-world applications. As a fundamental and broadly representative LLVP ability in many applications, we select geometry understanding as our domain of dataset construction. Following basic geometry defini-880 tions and axioms, we define seven basic geometric LLVP tasks in our benchmark. We include the full details of benchmark tasks and construction in Appendix B.

Current MLLMs struggle to perceive low-level geometry annotations and relationships. We 092 evaluate seven leading MLLMs, both open source and closed source. Their performances are shown 093 in Table 1. Despite the simplicity of Geoperception for humans, it remains a considerable challenge 094 for even the most advanced MLLMs. Notably, all models fall short of achieving 30% accuracy on 095 the *PointLiesOnLine* task and do not outperform the text-only GPT-40 mini model in *AngleClassi*-096 fication task. Closed source models generally outperform open source ones, with Gemini-1.5-pro attaining the highest average score of 56.98%, followed by gemini-1.5-flash at 54.76%. Among open 098 source models, Pixtral-12B achieves the best performance with an overall score of 41.95%. Noted 099 Cambrian-1 (Tong et al., 2024a), which is reported to be trained on Geo-170K (Gao et al., 2023), a geometry multimodal instruction tuning dataset built on the logical annotation of Geometry-3K, 100 the same source with Geoperception, still faces challenges in our Geoperception task, despite being 101 trained on the dataset having the same images and augmented text annotations. 102

103 104

105 106

3 EMPIRICAL STUDY ON MLLM DESIGN SPACE

We hypothesize the insufficient availability of high-fidelity visual description datasets to be a main 107 factor involved in this shortcoming. This may be because LLVP is intuitive for humans, and the Table 1: Performance (average evaluation score) of different models on Geoperception. POL: PointLiesOnLine, POC: PointLiesOnCircle, ALC: AngleClassification, LHC: LineComparison, PEP: Perpendicular, PRA: Parallel, EOL: Equals. As the Random Baseline method, we use GPT-40-mini, given the same textual instruction but without an image. The best model for each task is **bolded**.

	Log	gical	Num	erical	А	nnotatio	ns	
Model	POL	POC	ALC	LHC	PEP	PRA	EQL	Overall
Random Baseline	1.35	2.63	59.92	51.36	0.23	0.00	0.02	16.50
		Open So	ource					
Molmo-7B-D (Deitke et al., 2024)	11.96	35.73	56.77	16.79	1.06	0.00	0.81	17.59
Llama-3.2-11B (Dubey et al., 2024)	16.22	37.12	59.46	52.08	8.38	22.41	49.86	35.08
Qwen2-VL-7B (Wang et al., 2024a)	21.89	41.60	46.60	63.27	26.41	30.19	54.37	40.62
Cambrian-1-8B (Tong et al., 2024a)	15.14	28.68	58.05	61.48	22.96	30.74	31.04	35.44
Pixtral-12B (AI, 2023)	24.63	53.21	47.33	51.43	21.96	36.64	58.41	41.95
		Closed S	ource					
GPT-4o-mini (Achiam et al., 2023)	9.80	61.19	48.84	69.51	9.80	4.25	44.74	35.45
GPT-40 (Achiam et al., 2023)	16.43	71.49	55.63	74.39	24.80	60.30	44.69	49.68
Claude 3.5 Sonnet (Anthropic, 2024)	25.44	68.34	42.95	70.73	21.41	63.92	66.34	51.30
Gemini-1.5-Flash (Team et al., 2023)	29.30	67.75	49.89	76.69	29.98	63.44	66.28	54.76
Gemini-1.5-Pro (Team et al., 2023)	24.42	69.80	57.96	79.05	38.81	76.65	52.15	56.98



Figure 2: LLM size experiments. Training loss and testing accuracy curve comparing three choices of LLM size with a fixed visual encoder and multimodal connector. Training losses are windowsmoothed using a window size of 10 for better visibility.

datasets used to train MLLMs are primarily annotated by humans¹. When describing visual in-formation, humans often overlook LLVP details, assuming such details to be straightforward and self-evident. Furthermore, in the absence of sufficient datasets, it is challenging to identify and adopt architectural choices and training strategies that could improve the effectiveness and efficiency of LLVP-specific training. To this end, we develop a synthetic dataset generation engine to programmatically produce geometry shapes for training our MLLM. We include the full detail of our dataset generation engine in Appendix F. With the sufficient training dataset, we are able to conduct fully controlled experiments on different aspects of MLLMs focusing specifically on their LLVP abilities on geometry domain.

Exploration of MLLM design space. We follow the current most popular design of MLLMs (Liu et al., 2024c): a visual encoder, a LLM and an MLP connector in between. Although there exists many ready-to-use MLLMs off the shelf, they differ significantly in many aspects, such as archi-tectural choices (e.g., vision encoder, LLM, and multimodal connector), training datasets, and hy-perparameters. Consequently, due to the high cost in re-pretraining the visual encoders or LLMs, we opt to train our own MLLM with existing vision encoders and LLMs. This is similar with some recent empirical exploration in MLLM design space (Tong et al., 2024a; McKinzie et al., 2024), but with our large-scale controllable 2D geometry dataset. Specifically, we start with a typical setting of MLLMs: CLIP-ViT-L/14 (Radford et al., 2021) as the visual encode, a two layer MLP as multi-modal connector and the latest Qwen-2.5 series (Team, 2024b) as LLM. During training, we actively tune the MLP and LLM, while keeping visual encoder frozen.

¹This includes both explicit annotation in multimodal instruction-tuning datasets and implicit annotation in large-scale internet collections of text-image pairs



Figure 3: Vision encoder experiments. Training loss and testing accuracy (on a 1500 instances holdout test set) curve comparing eight visual encoders, with a fixed multimodal encoder and LLM.For a fair comparison, all visual encoder transcribe an image into 256 visual tokens. Training losses are window-smoothed using a window size of 10 for better visibility.

174

175

Lesson 1: Scaling LLM size does not benefit LLVP learning. Reflected in most MLLM releases (Liu et al., 2024a; Tong et al., 2024a; Wang et al., 2024a), scaling up the LLM often results in improved MLLM performance when trained on the same dataset. However, since their training datasets are complex mixtures from multiple sources, it is unclear whether this improvement stems from enhanced language-space reasoning or better visual perception. In contrast, our training dataset is designed with simplicity in the language space (following specific templates) and focuses on LLVP abilities. This setup allows us to better isolate and analyze the source of performance gains, providing clearer insights into the ability to learn better LLVP.

186 We use three variants of Qwen-2.5 (Team, 2024b): 0.5B, 1.5B, and 3B, while keeping other com-187 ponents in the MLLMs consistent and training them on the same dataset. The results are shown 188 in Fig. 2. First, we observe a sharp decrease in loss at the start of training, which corresponds to 189 the LLM adapting to answer templates, indicating no significant difficulty across different model sizes. The subsequent loss decrease after the plateau marks the beginning of learning LLVP, as also 190 evidenced by the testing accuracy. For LineComparison, Qwen-2.5-1.5B performs the best, while 191 Qwen-2.5-3B learns most slowly. For PointLiesOnLine, Qwen-2.5-1.5B and Qwen-2.5-3B perform 192 nearly identically. while Qwen-2.5-0.5B learns relatively slower but eventually reaches a similar fi-193 nal performance as the other models. In summary, we do not observe a clear trend that larger LLMs 194 learn LLVP tasks faster or better². Based on these findings, we will use Qwen-2.5-1.5B for further 195 exploration. 196

197

Lesson 2: CNN visual encoders are more suitable in LLVP learning than ViT architectures. We then study the choice of visual encoder architectures, including two families of architectures: 199 Vision Transformer (ViT) (Dosovitskiy, 2020) and ConvNeXT (Liu et al., 2022); as well as two 200 visual representation learning objectives: language-supervised learning (Radford et al., 2021) and 201 self-supervised learning (Oquab et al., 2023). We summarize the visual encoders in our experiment 202 in Table 4, for all vision encoders in our table, control the number of visual tokens to 256. The re-203 sult is shown in Fig. 3. We find that ConvNeXt-XXLarge and ConvNeXt-Large consistently learns 204 the fastest among all of the visual encoders. Notably, ConvNeXT-Large shows superior learning 205 performance with the vision transformers which are 3-5 times larger. We hypothesize that CNN 206 architecture extract visual features globally, effectively preserving low-level visual features. In con-207 trast, ViT architectures split images into discrete patches, making it more challenging to retain the original low-level visual information. Self-supervised learning (SSL) visual encoders, DINO-v2, 208 struggles to learn the geometry concept; we hypothesis this is due to the weak vision-language rep-209 resentation in these models. Surprisingly, although the SigLIP-family is widely-recognized as a 210 better visual encoder (Tong et al., 2024a; Li et al., 2024a), we find that their performance in learning 211 basic visual geometry attributes is limited. 212

 ²The BLINK (Fu et al., 2024c) benchmark shares similar observations. For example, among 14 tasks,
 LLaVA-1.5-13B outperforms its 7B variant in only 4 tasks. Additionally, LLaVA-one-vision (Li et al., 2024a)'s 0.5B variant outperforms its 7B variant by 3.9% and underperforms its 72B variant by only 3.3%.

229 230

231

232

233

234

235

236

237 238

239

240

241 242

243

244 245

246

247

248 249

250

251

252



Figure 4: Tuning/freezing vision encoder experiments. Testing accuracy (on a 1500 instances holdout test set) curve comparing freezing versus tuning the visual encoder during training.

Lesson 3: Tuning vision encoder does not offer strong advantage. By intuition, actively tuning the visual encoder can help it learn better visual representations which is helpful for LLVP. We aim to empirically justify the effect of tuning versus freezing the visual encoder. In Fig. 4, we show the testing accuracy curves of tuning and freezing visual encoders. Surprisingly, we find that compared with using a frozen encoder, tuning the visual encoder does not help the model learn LLVP faster or better. This suggests that current visual encoders seem to be able to preserve adequate information for LLVP, and it's sufficient to train LLMs to make better use of the visual features. In what follows, we will freeze the encoder for simplicity.



Figure 6: Curriculum learning experiments. Test accuracy on difficulty level hard of three training strategies: purely training on difficulty level hard, mixed training of difficulty levels easy/medium/hard, and curriculum training.

253 Lesson 4: Curriculum learning unleashes full po-254 tential. Finally, we study training data composi-In our preliminary experiment Fig. 19, we tion. observe that the model fails to converge on diffi-256 culty level 3 of PointLiesOnLine and difficulty level 257 2 and 3 of LineComparison. However, when using 258 mixed training set of all three difficulty levels, the 259 model achieves convergence, despite using the same 260 amount of data for each difficulty levels. We hy-261 pothesize that including easier levels aids the model 262 in learning more complex levels. To test this hypothesis, we report the test accuracy for three dif-264 ficulty levels separately during the mixed training of



Figure 5: Separate testing accuracy curves on difficulty levels easy, medium, and hard, shown over the course of training on a mixture of all difficulty levels.

ConvNeXt-XXLarge, in Fig. 5, on both tasks. We notice that the testing accuracy for easier tasks
 increase earlier and more quickly than difficulty tasks. In *PointLiesOnLine* tasks, we notice an appar ent plateau for hard level tasks until the model has trained on approximately 20K samples. During
 this period, the testing accuracy for easy and medium continue to increase. This suggests that learn ing easier shapes can significantly help the model tackle more challenging shapes, comparing with
 directly learning the challenging ones, aligning with the principles of curriculum learning.

270 While mixed training enables effective spontaneous curriculum learning, we investigate whether a 271 structured curriculum can further enhance model efficiency on challenging shapes. To this end, we 272 monitor the model's performance and dynamically increase the difficulty level of training data (i.e., 273 the curriculum stage) based on this performance. Specifically, the model starts by training on the 274 easy level data. and is evaluated when it finishes a training round, using testing accuracy from the current level of data. Upon evaluation, if the model achieves an accuracy exceeding a predefined 275 threshold θ , the framework advances the level to the next. Formally, the update rule for advancing 276 stages is given by: if accuracy $_{s} > \theta$ $c \leftarrow c+1$. The model is trained on a total of M \Rightarrow 277 rounds and K steps within each round. To avoid forgetting, we apply data smoothing at each stage. 278 Specifically, we smooth our dataset distribution over all stages using an exponential attenuation 279 function: ratio_s = exp $(-\alpha \cdot | \text{stage}_s - c |)$, where α denotes the attenuation rate. This ensures that 280 stages proximal to the current stage receive higher sampling probabilities. 281

We refer to this as our curriculum training strategy. Specifically, the accuracy threshold for advancing training stage θ is set to 0.99. We train all the models for M = 30 rounds, each round with K = 50 steps. The results are shown in Fig. 6. Firstly, we find that all of the models fail to converge when trained purely on hard level for *PointLiesOnLine* task. In contrast, the mixed training strategy shown by the red curve, consistently reaches faster convergence on hard level. Curriculum training strategy, shown by the purple curve, proves more efficient than mixed training.

287 288 289

4 EUCLID: A FAMILY OF MLLMS FOR GEOMETRIC VISUAL PERCEPTION

We take all of the lessons we learned in the previous sections and train Euclid, a family of MLLMs specifically designed for strong geometric LLVP. We use the same strategy as the curriculum training in Section 3, but scale our training to all tasks in Geoperception. For each task, we create N stages of training dataset shapes with progressively increasing geometric complexity.

Table 2: Performance comparison between Euclid and the best leading open source and closed source MLLMs on the seven tasks. Note that Euclid is *not* trained on any of the in-distribution data from the benchmark tasks below. The best model for each task is **bolded**.

	Log	gical	Num	erical	А	nnotatio	ns	
Model	POL	POC	ALC	LHC	PEP	PRA	EQL	Average
Random Baseline	0.43	2.63	59.92	51.36	0.25	0.00	0.02	16.37
Pixtral-12B (AI, 2023)	24.63	53.21	47.33	51.43	21.96	36.64	58.41	41.95
Gemini-1.5-Pro (Team et al., 2023)	24.42	69.80	57.96	79.05	38.81	76.65	52.15	56.98
Euclid-ConvNeXt-Large Euclid-ConvNeXt-XXLarge	80.54 82.98	57.76 61.45	86.37 90.56	88.24 90.82	42.23 46.96	64.94 70.52	34.45 31.94	64.93 67.89

Evaluation results. The results are shown in Table 2. Overall, although only trained on very 304 simple synthetic geometry shapes, and using only a 1.5B language model, Euclid significantly out-305 performs current leading MLLMs in most of the tasks, showing strong generalization abilities on 306 real-world geometry LLVP. Notably, in the *PointLiesOnLine* task, which is particularly challenging 307 for existing MLLMs, Euclid achieves up to 82.98% accuracy, more than three times the performance 308 of Gemini-1.5-Pro. On all both numerical tasks, LineComparison and AngleClassification, Euclid 309 keeps superior performance. However, on three annotation tasks, Euclid's performance is limited. 310 We hypothesis this is due to the limited setting of our annotation types and styles, making the model 311 hard to generalize to diverse human geometry annotations.

312 313

5 CONCLUSION

314 In this work, we highlight the importance of accurate *low-level visual perception* (LLVP) in MLLMs. 315 To this end, we first introduce Geoperception, a large-scale multimodal benchmark focused exclu-316 sively on geometry-domain LLVP. We find that even top models such as Gemini-1.5-Pro struggle 317 significantly, although it is straightforward for humans. We then conduct an empirical study to 318 explore the design space of MLLM training and architectures using the dataset generated by a geo-319 metric high-fidelity synthetic-data engine that we develop. Our key insides include that CNN-based 320 visual encoders outperform ViT in our tasks; and employing a curriculum-based training approach 321 yields much more model potential than direct task training. Based on insights from this study, we 322 develop Euclid, a model trained purely on synthetic data generalizes effectively to real-world geometric LLVP tasks, surpassing the leading MLLMs by a substantial margin. 323

324 REFERENCES 325

341

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-326 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical 327 report. arXiv preprint arXiv:2303.08774, 2023. 328
- Kian Ahrabian, Zhivar Sourati, Kexuan Sun, Jiarui Zhang, Yifan Jiang, Fred Morstatter, and Jay 330 Pujara. The curious case of nonverbal abstract reasoning with multi-modal large language models. 331 *arXiv preprint arXiv:2401.12117*, 2024.
- 332 Mistral AI. Pixtral 12b. https://mistral.ai/news/pixtral-12b/, 2023. Accessed: 333 2024-09-27. 334
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel 335 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language 336 model for few-shot learning. Advances in neural information processing systems, 35:23716-337 23736, 2022. 338
- 339 Anthropic. The claude 3 model family: Opus, Sonnet, Haiku, March 2024. URL https:// 340 www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/ Model_Card_Claude_3.pdf.
- Roman Bachmann, Oğuzhan Fatih Kar, David Mizrahi, Ali Garjani, Mingfei Gao, David Griffiths, 343 Jiaming Hu, Afshin Dehghan, and Amir Zamir. 4m-21: An any-to-any vision model for tens of 344 tasks and modalities. arXiv preprint arXiv:2406.09406, 2024. 345
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, 346 Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. 347 Paligemma: A versatile 3b vlm for transfer. arXiv preprint arXiv:2407.07726, 2024. 348
- 349 Markus J Buehler. Cephalo: Multi-modal vision-language models for bio-inspired materials analysis 350 and design. Advanced Functional Materials, pp. 2409531, 2024. 351
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 352 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In Proceedings 353 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14455–14465, 354 2024. 355
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Moham-356 madreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin 357 Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne 359 Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron 360 Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, 361 Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, 362 Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for 364 state-of-the-art multimodal models, 2024. URL https://arxiv.org/abs/2409.17146. 365
- Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. 366 arXiv preprint arXiv:2010.11929, 2020. 367
- 368 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 369 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 370 arXiv preprint arXiv:2407.21783, 2024.
- 371 Deqing Fu, Ruohao Guo, Ghazal Khalighinejad, Ollie Liu, Bhuwan Dhingra, Dani Yogatama, 372 Robin Jia, and Willie Neiswanger. Isobench: Benchmarking multimodal foundation models 373 on isomorphic representations. In First Conference on Language Modeling, 2024a. URL 374 https://openreview.net/forum?id=KZd1EErRJ1. 375
- Deqing Fu, Tong Xiao, Rui Wang, Wang Zhu, Pengchuan Zhang, Guan Pang, Robin Jia, and 376 Lawrence Chen. Tldr: Token-level detective reward model for large vision language models. 377 arXiv preprint arXiv:2410.04734, 2024b.

378	Xingyu Fu Yushi Hu Bangzheng Li Yu Feng Haoyu Wang Xudong Lin Dan Roth Noah A
379	Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but
380	not perceive. arXiv preprint arXiv:2404.12390, 2024c.
381	Jiahui Gao, Renije Pi, Jipeng Zhang, Jiacheng Ye, Waniun Zhong, Yufei Wang, Langing Hong
383	Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal

- 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.
 We L G and This We are the formula for the Data of the Multi-modal large language model. *arXiv preprint arXiv:2312.11370*, 2023.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings* of the IEEE conference on computer vision and pattern recognition, pp. 6904–6913, 2017.
- Jinyi Hu, Yuan Yao, Chongyi Wang, Shan Wang, Yinxu Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu,
 Yue Zhao, Haoye Zhang, et al. Large multilingual models pivot zero-shot multimodal learning
 across languages. *arXiv preprint arXiv:2308.12038*, 2023.
- Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. *arXiv preprint arXiv:2406.09403*, 2024.
 - Yifan Jiang, Jiarui Zhang, Kexuan Sun, Zhivar Sourati, Kian Ahrabian, Kaixin Ma, Filip Ilievski, and Jay Pujara. Marvel: Multidimensional abstraction and reasoning through visual evaluation and learning. *arXiv preprint arXiv:2404.13591*, 2024.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. Ge omverse: A systematic evaluation of large models for geometric reasoning. *arXiv preprint arXiv:2312.12241*, 2023.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed- ings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- Xianghao Kong, Ollie Liu, Han Li, Dani Yogatama, and Greg Ver Steeg. Interpretable diffusion via
 information decomposition. *arXiv preprint arXiv:2310.07972*, 2023.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024a.
- 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*, 2024b.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Zhiqiu Lin, Xinyue Chen, Deepak Pathak, Pengchuan Zhang, and Deva Ramanan. Revisiting the
 role of language priors in vision-language models. *arXiv preprint arXiv:2306.01879*, 2023.
- 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*, pp. 26296–26306, 2024a.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
- 430

392

396

397

398 399

406

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024c. 432 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, 433 Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around 434 player? arXiv preprint arXiv:2307.06281, 2023. 435 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 436 A convnet for the 2020s. In Proceedings of the IEEE/CVF conference on computer vision and 437 pattern recognition, pp. 11976–11986, 2022. 438 439 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. 440 Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. 441 arXiv preprint arXiv:2105.04165, 2021. 442 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-443 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of 444 foundation models in visual contexts. arXiv preprint arXiv:2310.02255, 2023. 445 446 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights 447 from multimodal llm pre-training. arXiv preprint arXiv:2403.09611, 2024. 448 449 David Mizrahi, Roman Bachmann, Oguzhan Kar, Teresa Yeo, Mingfei Gao, Afshin Dehghan, and 450 Amir Zamir. 4m: Massively multimodal masked modeling. Advances in Neural Information 451 Processing Systems, 36, 2024. 452 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, 453 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning 454 robust visual features without supervision. arXiv preprint arXiv:2304.07193, 2023. 455 456 Aitor Ormazabal, Che Zheng, Cyprien de Masson d'Autume, Dani Yogatama, Deyu Fu, Donovan 457 Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, Isaac Ong, et al. Reka core, flash, and edge: A 458 series of powerful multimodal language models. arXiv preprint arXiv:2404.12387, 2024. 459 Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. Multi-460 math: Bridging visual and mathematical reasoning for large language models. arXiv preprint 461 arXiv:2409.00147, 2024. 462 463 Cyril Picard, Kristen M Edwards, Anna C Doris, Brandon Man, Giorgio Giannone, Md Ferdous 464 Alam, and Faez Ahmed. From concept to manufacturing: Evaluating vision-language models for 465 engineering design. arXiv preprint arXiv:2311.12668, 2023. 466 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 467 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 468 models from natural language supervision. In International conference on machine learning, pp. 469 8748-8763. PMLR, 2021. 470 471 Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind. arXiv preprint arXiv:2407.06581, 2024. 472 473 Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, 474 real-time object detection. In Proceedings of the IEEE conference on computer vision and pattern 475 recognition, pp. 779-788, 2016. 476 477 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can 478 teach themselves to use tools. Advances in Neural Information Processing Systems, 36, 2024. 479 480 Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, De-An Huang, Hongxu 481 Yin, Karan Sapra, Yaser Yacoob, Humphrey Shi, et al. Eagle: Exploring the design space for 482 multimodal llms with mixture of encoders. arXiv preprint arXiv:2408.15998, 2024a. 483 Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy 484 Ka-Wei Lee. Math-llava: Bootstrapping mathematical reasoning for multimodal large language 485 models. arXiv preprint arXiv:2406.17294, 2024b.

486 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for 487 reasoning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 488 11888-11898, 2023. 489 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. arXiv preprint 490 arXiv:2405.09818, 2024a. 491 492 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, 493 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly 494 capable multimodal models. arXiv preprint arXiv:2312.11805, 2023. 495 Qwen Team. Qwen2.5: A party of foundation models, September 2024b. URL https: 496 //qwenlm.github.io/blog/qwen2.5/. 497 498 Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha 499 Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, 500 vision-centric exploration of multimodal llms. arXiv preprint arXiv:2406.16860, 2024a. 501 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide 502 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pp. 9568–9578, 2024b. 504 505 Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry 506 without human demonstrations. Nature, 625(7995):476-482, 2024. 507 Kirill Vishniakov, Zhiqiang Shen, and Zhuang Liu. Convnet vs transformer, supervised vs clip: 508 Beyond imagenet accuracy. arXiv preprint arXiv:2311.09215, 2023. 509 510 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, 511 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the 512 world at any resolution. arXiv preprint arXiv:2409.12191, 2024a. 513 Zhenhailong Wang, Joy Hsu, Xingyao Wang, Kuan-Hao Huang, Manling Li, Jiajun Wu, and Heng 514 Ji. Text-based reasoning about vector graphics. arXiv preprint arXiv:2404.06479, 2024b. 515 516 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, 517 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-518 agent conversation framework. arXiv preprint arXiv:2308.08155, 2023. 519 Sifan Wu, Amir Khasahmadi, Mor Katz, Pradeep Kumar Jayaraman, Yewen Pu, Karl Willis, and 520 Bang Liu. Cadvlm: Bridging language and vision in the generation of parametric cad sketches. 521 arXiv preprint arXiv:2409.17457, 2024. 522 523 Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth 524 anything: Unleashing the power of large-scale unlabeled data. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10371–10381, 2024. 525 526 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, 527 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multi-528 modal understanding and reasoning benchmark for expert agi. In Proceedings of the IEEE/CVF 529 Conference on Computer Vision and Pattern Recognition, pp. 9556–9567, 2024. 530 531 Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? arXiv preprint 532 arXiv:2210.01936, 2022. 534 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language 535 image pre-training. In Proceedings of the IEEE/CVF International Conference on Computer 536 Vision, pp. 11975–11986, 2023. 537 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, 538 Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the

diagrams in visual math problems? arXiv preprint arXiv:2403.14624, 2024a.

540 541 542	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. <i>arXiv preprint arXiv:2407.08739</i> , 2024b.
543	
544	wenwen Zhuang, Xin Huang, Xiantao Zhang, and Jin Zeng. Math-puma: Progressive upward
545	multimodal alignment to enhance mathematical reasoning. arXiv preprint arXiv:2408.08640,
546	2024.
547	
548	
549	
550	
551	
552	
553	
554	
555	
556	
557	
558	
559	
560	
561	
562	
563	
564	
565	
566	
567	
568	
569	
570	
571	
572	
573	
574	
575	
576	
577	
578	
579	
580	
581	
582	
583	
584	
585	
586	
587	
588	
589	
590	
591	
592	
593	

APPENDIX

596 597 598

594 595

599 600

601

613

A EXTENDED BACKGROUND AND RELATED WORK

602 Vision-Language MLLMs. While recent iterations of LLMs feature a standardized model architecture and pretraining recipe, MLLMs still often differ in design choices for infusing visual inputs. 603 One popular design is to align *continuous* visual features with the embedding space of a backbone 604 LLM (Liu et al., 2024a;b; Dubey et al., 2024; McKinzie et al., 2024; Tong et al., 2024a; Beyer 605 et al., 2024; AI, 2023; Wang et al., 2024a); another approach involves tokenizing visual inputs to 606 be trained jointly with language tokens (Team et al., 2023; Team, 2024a). These modules are often 607 infused with a decoder-only LLM, but others have explored encoder-decoder architectures to inte-608 grate a more varied collection of modalities (Alayrac et al., 2022; Mizrahi et al., 2024; Ormazabal 609 et al., 2024; Bachmann et al., 2024). Our study focuses on decoder MLLMs with a continuous vi-610 sual encoder, and we carry out an empirical study to explore the effect of synthetic dataset mixture, 611 training recipe, and encoder design (Liu et al., 2022; Radford et al., 2021; Zhai et al., 2023; Oquab 612 et al., 2023).

614 Geometry-Oriented MLLMs. At the core of these choices is the hardness in designing a module adept in general visual reasoning (McKinzie et al., 2024; Tong et al., 2024a). In this work, we 615 explore the optimal design of MLLMs specialized in low-level visual perception, a crucial aspect for 616 (among other applications) multimodal mathematical understanding (Lu et al., 2023; Zhang et al., 617 2024a). This paper supplements prior efforts in improving mathematical reasoning (Gao et al., 2023; 618 Zhang et al., 2024b; Zhuang et al., 2024; Li et al., 2024b; Peng et al., 2024; Shi et al., 2024b) with a 619 detailed study on the effect of dataset mixture, curriculum, and visual encoder, to reach a recipe that 620 elicits strong performance on geometric tasks (Kazemi et al., 2023) that require low-level perception. 621

622 **Evaluating LLVP.** Many benchmarks (Rahmanzadehgervi et al., 2024) have reported that 623 frontier-class MLLMs struggle with visual perception tasks, which are prerequisites for applica-624 tions that emphasize low-level geometric perception (Chen et al., 2024; Fu et al., 2024c), including 625 mathematical (Yue et al., 2024; Lu et al., 2023; Zhang et al., 2024a; Jiang et al., 2024) and spatial 626 reasoning (Chen et al., 2024; Fu et al., 2024b). These findings collectively identify that MLLMs ex-627 hibit a language prior (Lin et al., 2023)—a preference of textual inputs over visual inputs—leading 628 to a performance gap between modalities (Wang et al., 2024b; Zhang et al., 2024a; Fu et al., 2024a). Meanwhile, there lacks a high-quality benchmark that evaluates low-level geometric perception in 629 MLLMs, and the Geoperception benchmark represents a first effort to narrow this gap. This type of 630 efforts have led to significant improvements in certain capabilities of MLLMs, such as composition-631 ality of objects (Yuksekgonul et al., 2022; Kong et al., 2023). 632

633 **Improving LLVP.** Many prior works study *data-driven* approaches to improve low-level percep-634 tion skills. For example, Gao et al. (2023); Li et al. (2024b); Zhuang et al. (2024) employ a stan-635 dardized supervised finetuning recipe, and optionally adjust the training data mixture. This type of 636 training data is often synthesized from text-only math problems (Lu et al., 2021; Trinh et al., 2024) 637 or via rule-based systems (Kazemi et al., 2023). In parallel, Vishniakov et al. (2023); Shi et al. 638 (2024a); Tong et al. (2024b) have explored the design space of visual encoders for general-purpose 639 vision-language reasoning. We identify best practices over the union of these design spaces, and 640 then train small MLLMs with strong performance in low-level perception tasks.

Lastly, several works (Schick et al., 2024; Surís et al., 2023; Hu et al., 2024) have opted to augment
an MLLM with external APIs that process low-level features with specialized vision modules, such
as object detection (Redmon et al., 2016), segmentation (Kirillov et al., 2023), and depth estimation (Yang et al., 2024). While these agentic frameworks (Wu et al., 2023) present a promising
alternative that directly addresses the shortcomings of visual encoders, they are limited by their
scalability to novel use cases, and may be insufficient for precise tool routing that requires low-level
perception as a primer (Picard et al., 2023; Wu et al., 2024; Buehler, 2024).

648 B GEOPERCEPTION BENCHMARK DETAILS

650 **Benchmark Tasks.** Although there are many complex geometry problems in textbooks and real-651 world, the basic constitutions of it is relatively simple: only five axioms can already underpin all 652 further geometric shapes and reasoning steps, introduced by Euclid over two thousand years ago. 653 These axioms involve establishing and extending lines using points (Axioms 1 and 2), constructing circles from a point and a radius (Axiom 3), and defining perpendicularity (Axiom 4) and 654 parallelism (Axiom 5). Additionally, Euclid provided common notions regarding the properties 655 of equality. Accordingly, we define seven tasks in our Geoperception dataset: PointLiesOnLine 656 (POL), PointLiesOnCircle (POC), Parallel (PRA), Perpendicular (PEP), Equal (EQL), AngleClas-657 sification (ALC) and LineComparison (LHC). In geometric diagrams, perpendicularity, parallelism, 658 and equality are often indicated by annotation symbols. Thus, we classify *Parallel*, *Perpendicular*, 659 and Equal as annotated geometry understanding. Meanwhile, PointLiesOnLine, PointLiesOnCir-660 cle, AngleClassification, and LineComparison fall under primitive geometry shape understanding, 661 which includes both logical (PointLiesOnLine, PointLiesOnCircle) and numerical (AngleClassifica-662 tion, LineComparison) tasks. 663

Benchmark Construction. Thanks to the precise annotated logical forms for geometric diagrams
 from Geometry-3K (Lu et al., 2021), we are able to build a large scale benchmark based on it,
 focusing just on geometric LLVP without requiring any further reasoning. Four examples from
 Geoperception are illustrated in Fig. 1.

Data Filtering. Geoperception is sourced from the Geometry-3K (Lu et al., 2021) corpus, which offers precise logical forms for geometric diagrams, compiled from popular high-school textbooks. However, certain points in these logical forms are absent in the corresponding diagrams. To resolve this, we use GPT-40-mini MLLM to confirm the presence of all points listed in the logical forms. This process filters the 3,002 diagrams to retain 1,584, where at least one logical form fully represents its points in the diagram. A random inspection of 100 annotations reveals only two errors, indicating high annotation accuracy.

676 Converting Logical Forms Into Questions. We convert logical forms into question-and-answer 677 pairs for each of the seven tasks in Geoperception. In the Equals task, for example, we directly convert the logical form (e.g., Equals (LengthOf (Line (Q, T)), 86)) into a question-678 679 answer pair (e.g., Q: What is the length of line QT as annotated? A: 86). For PointLiesOnLine, two points on the line are chosen to form the question, with the remain-680 ing points on the line as the answer. Similarly, for PointLiesOnCircle, we ask which points lie 681 on the circle, using its center as the basis for the question. For Parallel and Perpendicular, 682 we represent each line by two points and query which other lines are parallel or perpendicular to it. 683 In AngleClassification, we ensure the queried angle is in the range of $[10, 80] \cup [100, 170]$ 684 degrees to avoid ambiguity. For LineComparison, we ensure that the shorter line is less than 685 70% of the length of the longer line. Since multiple equivalent questions can be generated for a sin-686 gle logical form (e.g., a line containing five points generates ${}^{5}P_{2}$ equivalent questions), we randomly 687 select one to avoid redundancy. Table 3 summarizes the question statistics for each task, as well as 688 the number of images involved. Extended examples from Geoperception are illustrated in Fig. 7. 689

Statistics. In Table 3, we provide more details on the Geoperception benchmark, such as the number of logic forms present before and after filtering, the number of questions, and the number of images. AngleClassification and LineComparison are directly derived from points coordinates without filtering.

694

668

Evaluation Details. During evaluation, we carefully craft the evaluation prompts for each question
 type to ensure the clarity, the full prompts can be found in Appendix E. We use greedy sampling
 during evaluation to get deterministic results. Additionally, GPT-40-mini without image input is
 used for generating the random baseline, employing the same textual instructions.

The open source models include Molmo-7B-D (Deitke et al., 2024), Cambrian-1-8B (Tong et al., 2024a), Qwen2-VL-7B (Wang et al., 2024a), Llama-3.2-11B (Dubey et al., 2024), and Pixtral-12B (AI, 2023). The closed-source models include GPT-4o-mini (Achiam et al., 2023), GPT-4o (Achiam et al., 2023), Claude-3.5-Sonnet (Anthropic, 2024), Gemini-1.5-flash (Team et al., 2023),

702 703	Predicate	# LF Before Filter	# LF After Filter	# Q	# I
704	PointLiesOnLine	6988	2567	1901	924
705	PointLiesOnCircle	1966	1240	359	322
706	Parallel	222	123	106	101
700	Perpendicular	1111	680	1266	456
707	Equals	6434	4123	4436	1202
708	AngleClassification	-	-	2193	1389
709	LineComparison	-	-	1394	1394
710					

Table 3: Statistics of the five predicates in our Geoperception dataset. Including number of logic forms before filter, after filter and the number of questions and images.

and Gemini-1.5-pro (Team et al., 2023). Additionally, GPT-4o-mini without image input is used for generating the random baseline, employing the same textual instructions. To prevent stretching, all images are padded to square dimensions before being fed into the models. During evaluation of a given question by an MLLM, let G denote the ground truth set of answers, and let P denote the predicted set of answers; then the evaluation score is defined as

Evaluation score =
$$\begin{cases} \frac{|P|}{|G|} & \text{if } P \subseteq G, \\ 0 & \text{otherwise.} \end{cases}$$
(1)



⁸¹⁰ C FUTURE DIRECTIONS.

Our work examines the potential of using synthetic multimodal data to improve MLLM performance in low-level geometric perception tasks. However, there are still directions that remain under-explored:

- 1. Automatic curriculum learning. Incorporating a more diverse dataset, including varied geometric shapes and different domain dataset, introduces challenges in defining the learning order. Rule based definition and manual curation may become impractical, necessitating automated strategies like hard negative sampling to organize the curriculum based on training loss or testing accuracy. This approach could streamline the process, reduce human effort, provide more suitable and efficient curriculum learning orders.
 - 2. Using a more-diverse training dataset. Currently, the text portion of our synthetic multimodal training data uses a restricted set of templates, and the model trained on such templates could fail to generalize to other question types; it could therefore be beneficial to increase the diversity of our training images as well as the instruction-following formats.
 - 3. Generalizing to other task domains. In this work, our study is focused on data from 2D geometry, as it provides a focused test bed of fundamental tasks. We believe the lessons we learn from this domain can be effectively generalized to a broader set of downstream domains that benefit from high-quality LLVP.

D EXPERIMENT DETAILS IN EMPIRICAL STUDY AND EUCLID TRAINING.

Experimental setting for empirical study. We use *PointLiesOnLine* and *LineComparison* as the test bed tasks for the exploration. For each task, we carefully create three levels with incremental difficulties. We name them as difficulty level easy, medium and hard. Based on the insight from our preliminary experiments, to increase the difficulty levels, for *PointLiesOnLine*, we increase the complexity of geometry shapes as is shown in Fig. 16, for *LineComparison*, we increase the total number of letters in letter pool while mixing geometry shapes. To report stable results, we run the training for three times and report the best run among them (i.e., having the lowest overall training loss or testing accuracy).

Table 4 summarize the visual encoder we use in our empirical study.

Table 4: Summary of Visual Encoders

Model	Params	Objective
ConvNeXt Large@512	200M	CLIP
ConvNeXt XXLarge@512	847M	CLIP
ViT-g/14@224	1.01B	CLIP
ViT-H/14@224	632M	CLIP
ViT-L/14@224	303M	CLIP
SigLIP@224 (ViT)	428M	CLIP-like
DINOv2 Giant@224 (ViT)	1.14B	Self-Sup
DINOv2 Large@224 (ViT)	304M	Self-Sup

Experimental setting for Euclid training. For models, we select the best visual encoder archi-tecture we found in our investigation, ConvNeXt, including ConvNeXt-Large@512 and ConvNeXt-XXLarge@512, and keep the same multimodal connector (2 layers MLP) and LLM (Qwen2.5-1.5Binstruct). The accuracy threshold for advancing training stage θ is set to 0.99. All models are trained on N = 3 stages with manually curated geometry shapes and M = 50 rounds with K = 500 steps in each round, and the batch size is 64 for each training step. The total training dataset volume for both of the models is 1.6M.

Euclid error analysis. We take a deep look into Euclid's prediction on *Geoperception* and find that its performance is hindered when diagrams are heav-ily annotated. An example is shown in Fig. 8, where a line is annotated by "x", preventing the model from choosing the correct point. We hypothesize that in-corporating training data with more diverse annota-tion types and geometry shapes could help the model with such scenarios.



Figure 8: An error case where Euclid fails to predict the correct point on a line, potentially distracted by the annotation "x".

E PROMPTS FOR THE GEOPERCEPTION DATASET EVALUATION

PROMPT TEMPLATE FOR THE POINTLIESONLINE TASK

```
Answer me directly just with the all points lie on the line
mentioned in the question (do not include the point mentioned in
the question).
Answer template:
  (If only one point) The other point is: "your point".
Or
  (if multiple points) The other points are: "your points".
For example:
  The other point is: A
Or
  The other points are: A, B, C
```

Figure 9: TEMPLATE FOR THE POINTLIESONLINE TASKS

PROMPT TEMPLATE FOR THE POINTLIESONCIRCLE TASK

```
Answer me directly just with the all points lie on the circle
mentioned in the question.
Answer template:
(If only one point) The point is: "your point".
Or
(If multiple points) The points are: "your points".
For example:
The point is: A
Or:
The points are: A, B, C
```

Figure 10: TEMPLATE FOR THE POINTLIESONCIRCLE TASKS

PROMPT TEMPLATE FOR THE PARALLEL TASK

```
Answer me directly just with the all lines which are parallel
to the line mentioned in the question (do not include the line
mentioned in the question).
Answer template:
  (If only one line) The line is: "your line".
Or
   (If multiple lines) The lines are: "your lines".
For example:
   The line is: BC
Or:
   The lines are: BC, DE
```

Figure 11: TEMPLATE FOR THE PARALLEL TASKS

```
972
973
974
975
976
976
977
978
978
```

```
PROMPT TEMPLATE FOR THE PERPENDICULAR TASK
Answer me directly just with the all lines which are perpendicular
to the line mentioned in the question (do not include the line
mentioned in the question).
Answer template:
  (If only one line) The line is: "your line".
Or
  (If multiple lines) The lines are: "your lines".
For example:
  The line is: BC
Or:
  The lines are: BC, DE
```

Figure 12: TEMPLATE FOR THE PERPENDICULAR TASKS

PROMPT TEMPLATE FOR THE EQUALS TASK

```
Answer me directly just with the annotations presented on the
image.
Answer template:
The annotation is: "your annotation".
For example:
The annotation is: 2x+4
Or:
The annotations is: 90
```

Figure 13: TEMPLATE FOR THE EQUALS TASKS

PROMPT TEMPLATE FOR THE ANGLE CLASSIFICATION TASK

```
Answer me directly just with the classification of the angle
mentioned in the question.
Answer template:
The angle is: "your angle".
For example:
The angle is: acute
Or:
The angle is: obtuse
```

Figure 14: TEMPLATE FOR THE ANGLE CLASSIFICATION TASKS

```
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
```

PROMPT TEMPLATE FOR THE LINECOMPARISON TASK

```
Answer me directly just with the longer line mentioned in the
question.
Answer template:
The longer line is: "your line".
For example:
The longer line is: BC
Or:
The longer line is: DE
Figure 15: TEMPLATE FOR THE LINECOMPARISON TASKS
```

1026 F DETAILS FOR TRAINING DATA ENGINE

In this section, we provide all geometry shapes we use for Euclid training, including the pseudocode for generating text describing the geometry shapes and diagram examples.

1032 Our geometry dataset generation engine is built on 1033 AlphaGeometry (Trinh et al., 2024). Given an in-1034 put formal language describing a geometry shape, 1035 the geometry engine is able to render infinite ac-1036 tual geometry images (exampled in Fig. 16 with full 1037 metadata (e.g. point coordinates) then create dif-1038 ferent types of accurate textual annotations to train MLLMs with corresponding images. 1039 1040

FOR TRAINING TEXTUAL DATASET SYNTHESIS

1041 F.1 PSEUDOCODE

1042

```
A B C = triangle A B C;

D = midpoint B C

D = midpoint B C

E = midpoint A B;

E = midpoint A C;

D = midpoint A C;

E = midpoint A C;

F = intersection, IA D B E

Geometry Dataset Generation Engine
```

Figure 16: Three geometry logical shapes, of increasing complexity, used in our empirical study. Our geometry image generation engine is able to produce infinite visual instances for each of these logical shapes. All letters are randomly sampled from the alphabet and reassigned to each of the points before drawing.

1043 1044 Algorithm 1 Data Synthesis for the POINTLIESONLINE Task 1045 1: Input: data_info, points_set 1046 2: Output: data 1047 3: for points_set ∈ data_info do 1048 4: for (A, B) \in permutations (points_set, 2) do 1049 5: all_rest_points \leftarrow [p for p in points_set if p not in [A, B]] 1050 for rest_points ∈ permutations (all_rest_points) do 6: 1051 7: verb_agreement \leftarrow 'is' if len(rest_points) == 1 else 1052 'are' 1053 8: rest_points \leftarrow [f"{p}" for p in rest_points] 1054 9: rest_points ← sorted(rest_points) 1055 question \leftarrow 'What is the point lying on line ' + A + B + 10: 1056 1 ?1 1057 answer \leftarrow 'The point lying on line ' + A + B + ' ' + 11: verb_agreement + ' ' + ', '.join(rest_points) 1058 $qt \leftarrow ''.join(rest_points)$ 1059 12: data \leftarrow {'question': question, 'answer': answer, 'gt': 13: gt} 1061 14: end for 1062 15: end for 1063 16: end for 1064 Algorithm 2 Data Synthesis for the POINTLIESONCIRCLE Task 1067 1: Input: data_info 1068 2: Output: data 1069 3: point_set ← random.choice(list(data_info.items())) 1070 4: center_point ← point_set[0] 1071 5: target_points ← point_set[1] 1072 6: target_points ← sorted(target_points) 7: question \leftarrow 'What are the point lying on circle ' + center_point + '?' 1074 1075 8: answer \leftarrow 'The point lying on circle ' + center_point + ' are ' + ', '.join(target_points) 1076 9: gt <- ''.join(target_points) 1077 10: data \leftarrow {'question': question, 'answer': answer, 'gt': gt} 1078 1079

```
1080
       Algorithm 3 Data Synthesis for the ANGLECLASSIFICATION Task
1081
       1: Input: data_info
1082
       2: Output: data
1083
       3: angle \leftarrow data_info
1084
       4: angle_options \leftarrow [f'{angle[1][0]}{angle[1][1]}{angle[1][2]}',
1085
          f'{angle[1][2]}{angle[1][1]}{angle[1][0]}']
1086
       5: angle_letter ← random.choice(angle_options)
1087
       6: angle_class ← 'acute' if angle[0] < 90 else 'obtuse'
1088
       7: question \leftarrow 'Is angle ' + angle_letter + ' acute or obtuse?'
       8: answer <- 'Angle ' + angle_letter + ' is ' + angle_class
1089
       9: gt \leftarrow angle_class
1090
       10: data \leftarrow {'question': question, 'answer': answer, 'qt':
                                                                              gt}
1091
1092
1093
       Algorithm 4 Data Synthesis for the LINECOMPARISON Task
1094
       1: Input: data_info
1095
       2: Output: data
1096
       3: names \leftarrow [data_info[0][1], data_info[1][1]]
1097
       4: lengths ← [data_info[0][0], data_info[1][0]]
1098
       5: if lengths[0] > lengths[1] then
1099
             longer_name, shorter_name ← names[0], names[1]
       6:
1100
       7: else
1101
       8:
            longer_name, shorter_name ← names[1], names[0]
1102
       9: end if
1103
       10: data \leftarrow [
1104
            { 'question': 'Which line is longer, ' + longer_name + ' or '
       11:
1105
          + shorter_name + '?',
              'answer': 'The longer line is ' + longer_name,
       12:
1106
              'gt': longer_name },
       13:
1107
             { 'question': 'Which line is longer, ' + shorter_name + ' or
       14:
1108
          ' + longer_name + '?',
1109
       15:
              'answer': 'The longer line is ' + longer_name,
1110
       16:
              'gt': longer_name }
1111
       17: ]
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
```

	STUIN 5 Data Synthesis for the PARALLEL Task
1:	Input: data_info
2: 3.	Oulpul: data
3. 4:	for line points E points set do
5:	for (A, B) \in permutations (line_points, 2) do
6:	all_rest_lines ← [p for p in points_set if p !=
	line_points]
7:	gts \leftarrow [''.join(
8:	f`{p}' for line in all_rest_lines for p in line)
9:]
10:	rest_point_pairs ← []
11:	IOF rest_line < all_rest_lines do
12:	$C, D \leftarrow \text{random.sample(rest_rine, 2)}$
13. 14·	end for
15:	all possible answer \leftarrow '. '. join(
16:	[f'{C}{D}' for C, D in rest_point_pairs]
17:)
18:	verb_agreement \leftarrow `is' if len(rest_point_pairs) == 1 el:
	`are'
19:	question \leftarrow `What is the line parallel to line ' + A +
20	·?/
20:	answer \leftarrow (
21.222	A + B + verb agreement + all possible answer
23:)
24:	$qt \leftarrow ', '.join(qts)$
25:	data \leftarrow {
26:	`question': question, `answer': answer, `task': t
	'gt': gt
27:	}
28:	end for
29:	end for
27: 28: 29:	} end for end for

1188 1189 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200 1201 1202 1203 Algorithm 6 Data Synthesis for the PERPENDICULAR Task 1204 1: **Input:** data_info 1205 2: Output: data 1206 3: source_lines, target_lines ← data_info 1207 4: all_possible_answer \leftarrow [] 1208 5: gts \leftarrow target_lines > Randomly choose two points from each target line 1209 6: for target_line \in target_lines do 1210 C, D \leftarrow random.sample(target_line, 2) 7: 1211 8: all_possible_answer.append(f`{C}{D}') 1212 9: end for 1213 10: verb_agreement \leftarrow 'is' if len(all_possible_answer) == 1 else 1214 'are' 1215 11: for (A, B) \in permutations (source_line, 2) do 1216 12: question \leftarrow 'What is the line perpendicular to line ' + A + B + `?' 1217 13: answer \leftarrow (1218 'According to the diagram, the line perpendicular to ' + 14: 1219 A + B + verb_agreement + `, '.join(all_possible_answer 15: 1220 16:) 1221 gt \leftarrow ', '.join(gts) 17: 1222 data \leftarrow { 18: 1223 'question': question, ', 'answer': answer, '`ask': task, 19: 1224 `gt': gt 1225 20: } 1226 21: end for 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239 1240 1241

```
1243
1244
1245
1246
      Algorithm 7 Data Synthesis for the EQUAL Task
1247
      1: Input: data_info
1248
      2: Output: data
1249
      3: statement, content ← data_info.split(`;')
1250
      4: if statement == 'angles_value' then
1251
      5:
           angle_letter, angle_measure ← content.split('=')
1252
           angle_letter ← random.choice([angle_letter,
      6:
        angle_letter[::-1]])
1253
          question \leftarrow `What is the measure of angle ' + angle_letter +
      7:
1254
         ' as annotated?'
1255
          answer \leftarrow 'Angle ' + angle_letter + ' is annotated as ' +
      8:
1256
        angle_measure
1257
      9:
           gt \leftarrow angle_measure
1258
      10: else if statement == `segments_value' then
1259
           segment_letter, segment_length \leftarrow content.split(`=')
      11:
1260
      12:
           segment_letter 
we random.choice([segment_letter,
1261
         segment_letter[::-1]])
1262
           question \leftarrow 'What is the length of line ' + segment_letter +
      13:
         ' as annotated?'
1263
      14:
          answer \leftarrow 'Line ' + segment_letter + ' is annotated as ' +
1264
        segment_length
1265
      15:
           gt \leftarrow segment\_length
1266
      16: else if statement == `angles' then
1267
           angle1, angle2 ← content.split('=')
      17:
1268
           18:
1269
      19:
           1270
      20:
           query_angle 
    random.choice([angle1, angle2])
1271
      21:
           answer_angle ← angle2 if query_angle == angle1 else angle1
1272
      22.
           question \leftarrow 'What is the angle in the diagram that is equal
1273
        to angle ' + query_angle
           answer \leftarrow 'Angle ' + query_angle + ' is equal to angle ' +
1274
      23:
         answer_angle
1275
      24:
           gt \leftarrow answer_angle
1276
      25: else if statement == `segments' then
1277
           26:
1278
           27:
1279
      28:
           segment2 ← random.choice([segment2, segment2[::-1]])
1280
      29:
           1281
      30:
           answer_segment 

segment2 if query_segment == segment1 else
1282
         segment1
1283
      31:
           question \leftarrow `What is the segment in the diagram that is equal
1284
         to segment ' + query_segment
           answer \leftarrow 'Segment ' + query_segment + ' is equal to segment
1285
      32:
         ' + answer_segment
1286
      33:
           gt ← answer_segment
1287
      34: end if
1288
      35: data \leftarrow
                {
1289
                 'question': question, 'answer': answer, 'task': task,
      36:
1290
         `gt′:
               gt
1291
      37: }
1292
1293
1294
1295
```

1296 F.2 GEOMETRY SHAPES USED FOR EUCLID TRAINING

GEOMETRY SHAPE GENERATION CODE

1300 PointLiesOnLine = triangle A B C; D = midpoint B C 1301 (stage 1) A B C (stage 1) A B C = triangle A B C; D = midpoint B C; O = circle O A B C (stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C 1302 (Stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C; O = circle O A B C (stage 3) A B C = triangle A B C; D = midpoint B C; E = midpoint A C; F = intersection.ll A D B E (stage 3) A B C = triangle A B C; D = midpoint B C; E = midpoint A C; F = intersection.ll A D B E; O 1303 1304 circle O A B C 1305 PointLiesOnCircle (stage 1) A B = segment A B; C = on_circle C A B (stage 1) A B = segment A B; C = on_circle C A B; D = on_circle D A B 1306 (stage 1) A B = segment A B; C = on.circle C A B; D = on.circle D A B; E = on.circle E A B (stage 1) A B = segment A B; C = on.circle C A B; D = on.circle D A B; E = on.circle E A B; F = 1307 on_circle F A B 1308 (stage 1) A B = segment A B; C = on_circle C A B; D = on_circle D A B; E = on_circle E A B; F = on_circle F A B; G = on_circle G A B 1309 (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = on_circle E A B 1310 1311 (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = on_circle E A B; F = 1312 on_circle F A B (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = on_circle E A B; F = 1313 on_circle F A B; G = on_circle G A B (stage 2) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = on_circle E A B; F = 1314 on_circle F A B; G = on_circle G A B; H = on_circle H A B (stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = midpoint A C (stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = midpoint A C; F = on_circle 1315 1316 FAB (stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = midpoint A C; F = on_circle 1317 F A B; G = on_circle G A B (stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = midpoint A C; F = on_circle 1318 F A B; G = on_circle G A B; H = on_circle H A B
(stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = midpoint A C; F = on_circle 1319 F A B; G = on_circle G A B; H = on_circle H A B; I = on_circle I A B 1320 (stage 3) A B = segment A B; C = on_circle C A B; D = midpoint A B; E = on_circle E A B; F = on_circle F A B; G = on_circle G A B; H = on_circle H A B; I = midpoint B C 1321 (stage 3) A B = segment A B; C = on.circle C A B; D = midpoint A B; E = midpoint B C (stage 3) A B = segment A B; C = on.circle C A B; D = midpoint A B; E = lc.tangent E C A (stage 3) A B = segment A B; C = on.circle C A B; D = midpoint A B; E = on.circle E A B; F = 1322 1323 on_circle F A B; G = on_circle G A B; H = lc_tangent H C A AngleClassification 1324 (stage 1) A B C = triangle A B C (stage 3) A B C = triangle A B C; D = midpoint B C (stage 3) A B C = triangle A B C; D = midpoint B C; E = midpoint A C; F = intersection.ll F A D B E 1325 1326 LengthComparison (stage 1) A B C = triangle A B C (stage 2) A B C = triangle A B C; D = midpoint B C (stage 3) A B C = triangle A B C; D = midpoint A B; E = midpoint A C 1327 1328 Parallel (stage 1) A B C = triangle A B C; D = midpoint A B; E = midpoint A C 1329 (stage 1) A B C = triangle A B C; D = midpoint A B; E = midpoint A C (stage 1) A B C = triangle A B C; D = midpoint A B; E = midpoint A C (stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C 1330 1331 (stage 3) A B C = triangle A B C; D = midpoint A B; E = midpoint A C; F = midpoint B C Perpendicular 1332 (stage 1) A B C = triangle A B C; D = foot A B C (stage 1) A B C = r.triangle A B C (stage 1) A B C = r.triangle A B C (stage 1) A B = segment A B; C = eq.triangle C A B; D = eq.triangle D A B; E = on_circle E A B 1333 1334 (stage 2) A B C = triangle A B C; D = foot A B C; E = foot C A B (stage 2) A B C = r.triangle A B C; D = foot A B C (stage 2) A B C = triangle A B C; D = foot A B C (stage 2) A B C = triangle A B C; O = circle A B C; D = foot O A B; E = foot O C A 1335 (stage 3) A B C D = rectangle A B C D; E = intersection.ll A C B D (stage 3) A B C = triangle A B C; O = incenter A B C; D = foot O A C; E = foot O B C; F = foot O A 1336 1337 (stage 3) A B C = r.triangle A B C; D = foot A B C; E = foot D A B (stage 3) A B C = triangle A B C; D = foot A B C; E = foot C A B; F = foot B A C 1338 Equal 1339 (stage 1) A B C = triangle A B C; D = midpoint C B (stage 1) A B C = triangle A B C; D = midpoint C B; O = circle O A B C 1340 (stage 1) A B C = triangle A B C; D = angle_bisector B A C, online D C B (stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C (stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C; O = circle O A B C 1341 1342 (stage 2) A B C = triangle A B C; D = midpoint A B; E = midpoint A C 1343 (stage 3) A B C = triangle A B C; O = circle A B C; D = on circle D O C, angle bisector C A B 1344 Figure 17: GEOMETRY SHAPE GENERATION CODE FOR EUCLID TRAINING 1345 1346 1347

1348

1298

1299



Figure 18: Examples of the geometry diagrams used to train Euclid, the diagrams are generated by our dataset engine.



Figure 19: Result of our preliminary experiments, we use a standard setting of MLLMs: an OpenAI-CLIP@224 as visual encoders (Radford et al., 2021), two-layer MLP as multimodal connector and Qwen-2.5-1.5B as language model. We find that the model can reach convergence in some of the easy tasks, while struggle to learn hard tasks. We also find mixed training is better than separate training, given the same amount of training data in each difficulty level.





Figure 20: The complete result of the effect of LLM size. The finding is similar with Fig. 2.



Figure 21: The complete result of the effect of different visual encoders. The finding is similar with Fig. 3.



Figure 22: The complete result of the effect of tuning visual encoders. The finding is similar with Fig. 4.