

EgoReasoner: Metric Spatial Reasoning in VLMs for Embodied Agents

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Embodied agents require metric spatial self-localization*
 002 *to act effectively, yet vision-language models (VLMs)—*
 003 *increasingly used as foundation models for such agents—*
 004 *consistently struggle with metric camera geometry. We*
 005 *present **EgoReasoner**, a two-stage framework that trains*
 006 *VLMs to estimate camera pose trajectories via geomet-*
 007 *ric chain-of-thought reasoning, combining supervised fine-*
 008 *tuning with GRPO reinforcement learning on geometry-*
 009 *grounded rewards. Trained on 684,455 frames from*
 010 *RealEstate10K, our model achieves 3.2× and 1.7× higher*
 011 *structured output parse rates over a base 8B VLM and Gem-*
 012 *ini 3.1 Pro, and surpasses Gemini in translation accuracy*
 013 *by over 6×. These results demonstrate that structured train-*
 014 *ing can instill metric geometric reasoning in VLMs, advanc-*
 015 *ing them toward spatially-aware foundation models for em-*
 016 *bodied agents.*

017 1. Introduction

018 Embodied agents must localize themselves spatially to act
 019 purposefully—a prerequisite for long-horizon interaction,
 020 navigation, and object manipulation in the physical world.
 021 Vision-language models (VLMs) have become a dominant
 022 backbone for embodied agent perception [1, 10], yet they
 023 consistently fail at metric camera geometry: existing VLMs
 024 cannot reliably infer rotation magnitudes, translation scales,
 025 or trajectory consistency from image sequences. Traditional
 026 pose estimation pipelines [11, 13, 14] fill this role in isola-
 027 tion, but remain specialized modules disconnected from the
 028 semantic reasoning layers central to modern agent systems.
 029 Closing this gap—enabling a single VLM to jointly reason
 030 about scene semantics *and* its own metric motion—is a key
 031 step toward foundation models for spatially-aware embod-
 032 ied agents.

033 We present **EgoReasoner**, to our knowledge the first
 034 framework to train a VLM to estimate metric camera pose
 035 trajectories as an explicit reasoning task. EgoReasoner
 036 employs a two-stage pipeline: (1) supervised fine-tuning
 037 (SFT) on structured five-step geometric chain-of-thought

chains that decompose pose estimation into epipolar reason- 038
 ing, depth estimation, 3D triangulation, and relative 039
 pose integration; followed by (2) GRPO [12] reinforce- 040
 ment learning with geometry-grounded reward signals to 041
 further refine metric accuracy without degrading output for- 042
 mat reliability. Trained on 684,455 frames from two non- 043
 overlapping RealEstate10K splits, the pipeline yields consis- 044
 tent, parseable trajectory outputs and strong pose accu- 045
 racy across all evaluation metrics. 046

On a 200-sequence holdout, EgoReasoner (SFT+RL) 047
 achieves 3.2× higher structured output parse reliability 048
 than an unmodified 8B VLM and 1.7× higher than Gem- 049
 ini 3.1 Pro [5], while surpassing Gemini in translation accu- 050
 racy by over 6×. These results demonstrate that metric ge- 051
 ometric reasoning is learnable within VLMs through struc- 052
 tured training, with direct implications for building founda- 053
 tion models capable of unified spatial and semantic percep- 054
 tion for embodied agents. 055

Contributions: 056

- **Novel direction.** To our knowledge, EgoReasoner is the 057
 first work to formulate metric camera pose trajectory es- 058
 timation as an explicit VLM reasoning task, bridging 059
 language-grounded perception and metric spatial self- 060
 localization for embodied agents. 061
- **Method.** A two-stage SFT+RL training pipeline com- 062
 bining five-step geometric chain-of-thought supervision 063
 with GRPO geometry-grounded rewards, enabling reli- 064
 able structured output and progressive metric accuracy 065
 improvement. 066
- **Results.** EgoReasoner surpasses Gemini 3.1 Pro in trans- 067
 lation accuracy by over 6× (RTE median 0.028 vs. 0.190) 068
 and achieves 3.2× and 1.7× higher output parse reliabil- 069
 ity than a base 8B VLM and Gemini 3.1 Pro respectively, 070
 with RRE median 1.21° on the full 200-sequence evalua- 071
 tion set. 072

2. Related Work 073

Camera Pose Estimation and 3D Reconstruction. 074
 Structure-from-Motion systems such as COLMAP [11] re- 075
 cover camera poses and sparse 3D structure from fea- 076

077	ture correspondences via bundle adjustment. Recent feed-	125
078	forward methods including DUS _t 3R [14], MAS _t 3R [9],	126
079	VGGT [13], and π^3 [15] estimate camera geometry directly	127
080	in a single forward pass, providing strong learned priors for	128
081	multi-view reconstruction. However, these methods treat	129
082	pose estimation as direct geometric regression without inter-	130
083	pretable reasoning. EgoReasoner instead studies metric	131
084	camera trajectory estimation as a language-mediated geo-	132
085	metric reasoning problem within a VLM.	133
086	Vision-Language Models for Spatial Reasoning. Large	134
087	VLMs such as LLaVA [10], InternVL [4], and Qwen-	135
088	VL [1] have shown strong capabilities in visual question	136
089	answering, captioning, and general visual reasoning. Re-	137
090	cent works extend VLMs toward 3D spatial understanding,	138
091	including SpatialVLM [3] for spatial relation reasoning and	139
092	3D-LLM [7] for 3D scene understanding. EgoReasoner tar-	140
093	gets the less-explored setting of metric camera pose trajec-	141
094	tory inference through explicit step-by-step geometric rea-	142
095	soning.	143
096	Reasoning via Supervised and Reinforcement Learn-	144
097	ing. Chain-of-thought prompting [16] and RL with veri-	145
098	fiable rewards (DeepSeek-R1 [6]) show that structured su-	146
099	pervision and ground-truth signals can elicit complex rea-	147
100	soning in LLMs beyond what fine-tuning alone achieves.	148
101	EgoReasoner adapts this to geometric reasoning: ground-	149
102	truth camera poses serve as a natural verifiable reward, en-	150
103	abling VLMs to improve both final trajectory accuracy and	151
104	the consistency of intermediate geometric reasoning steps.	152
105	3. Method	153
106	3.1. Problem Formulation	154
107	Given a sequence of N consecutive RGB images	155
108	$\{I_1, I_2, \dots, I_N\}$ captured by a moving camera, EgoRe-	156
109	asoner aims to predict the camera pose trajectory	157
110	$\{\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_N\}$, where each $\mathbf{T}_i \in \text{SE}(3)$ represents	158
111	the camera-to-world transformation at frame i . Following	159
112	standard practice, poses are represented as rotation matri-	160
113	ces $\mathbf{R}_i \in \text{SO}(3)$ and translation vectors $\mathbf{t}_i \in \mathbb{R}^3$.	161
114	Unlike conventional regression-based approaches,	162
115	EgoReasoner formulates this as a <i>reasoning task</i> : the VLM	163
116	is asked to produce a structured natural language chain of	164
117	thought that explains <i>how</i> it arrives at the pose estimates,	165
118	before outputting the final numerical trajectory. This	166
119	formulation enables the model to leverage its pre-trained	167
120	visual and language understanding to decompose the	168
121	problem into interpretable sub-steps.	169
122	3.2. Geometric Reasoning Chain Construction	170
123	A key design principle is that <i>all mathematical steps are</i>	171
124	<i>computed by a deterministic program</i> , while the VLM sup-	172
	plies logical inter-step reasoning insights and linguistically	173
	diverse natural language descriptions—populating a varied	174
	prompt pool and preventing the training data from degener-	175
	ating into formulaic concatenations of numerical outputs.	
	Synthetic correspondence generation. Rather than rely-	
	ing on a feature matcher, we generate pixel correspondences	
	purely by geometric projection using GT poses and depth	
	maps from Depth Anything V3 [17] (applied jointly to 5	
	frames). For each sequence, 5 reference points $\{(u_j, v_j)\}$	
	are sampled in Frame 0. Each point is back-projected to 3D	
	as $\mathbf{X}_j = d_j \cdot \mathbf{K}^{-1}[u_j, v_j, 1]^\top$, transformed to Frame i via	
	the GT relative pose: $\mathbf{X}'_j = \mathbf{R}\mathbf{X}_j + \mathbf{t}$, and re-projected to	
	get the corresponding pixel $[u'_j, v'_j] = \mathbf{K}(\mathbf{X}'_j/X'_{j,z})$ along	
	with its depth. This yields exact 2D–3D correspondences	
	across all frames with no feature matching required.	
	Chain-of-thought assembly. Each training sample con-	
	sists of a user query (5 RGB frames + \mathbf{K} + instruction)	
	paired with a five-step assistant response: (1) point se-	
	lection and matching description— <i>VLM-generated</i> natural	
	language; (2) depth estimation at selected points— <i>VLM-</i>	
	<i>generated</i> ; (3) 3D back-projection— <i>program-filled</i> exact	
	arithmetic; (4) Procrustes pose estimation via SVD on 3D	
	correspondences to solve (\mathbf{R}, \mathbf{t}) — <i>program-filled</i> ; (5) tra-	
	jectory summary— <i>VLM-generated</i> . The final numerical	
	pose sequence is delimited by regex-parseable tokens, en-	
	abling automatic evaluation of both intermediate step qual-	
	ity and end-to-end trajectory accuracy.	
	3.3. Supervised Fine-Tuning	
	EgoReasoner is initialized from Qwen3-VL-8B-	
	Instruct [2], a state-of-the-art open-source VLM capable of	
	multi-image input. The model is fine-tuned using LoRA [8]	
	applied to all attention projection layers ($\mathbf{q}, \mathbf{k}, \mathbf{v}, \mathbf{o}$) and	
	MLP projection layers ($\text{gate}, \text{up}, \text{down}$), with rank	
	$r = 16$, scaling factor $\alpha = 32$, and dropout rate 0.05.	
	The model is trained on the geometric reasoning chain	
	dataset using standard next-token prediction loss. Special	
	formatting tokens demarcate the chain-of-thought segment	
	from the final numerical output.	
	Training is performed for 1 epoch over 62,834	
	RealEstate10K sequences (from 30,702 source videos), us-	
	ing a cosine learning rate of 1×10^{-4} with 10% warmup	
	and an effective batch size of 24 across 3 GPUs. Training	
	completes in 7,062 steps (72.8 hours), reaching a final eval	
	loss of 0.2238.	
	3.4. Reinforcement Learning with Geometry Re-	
	wards	
	To further improve numerical accuracy beyond SFT, we ap-	
	ply GRPO with geometry-grounded reward signals, build-	
	ing on the SFT-initialized model.	
	The reward function is $r = r_{\text{fmt}} + r_{\text{val}} + r_{\Delta R} + r_{\Delta t}$,	
	where r_{fmt} (weight 0.2) rewards well-formed output struc-	

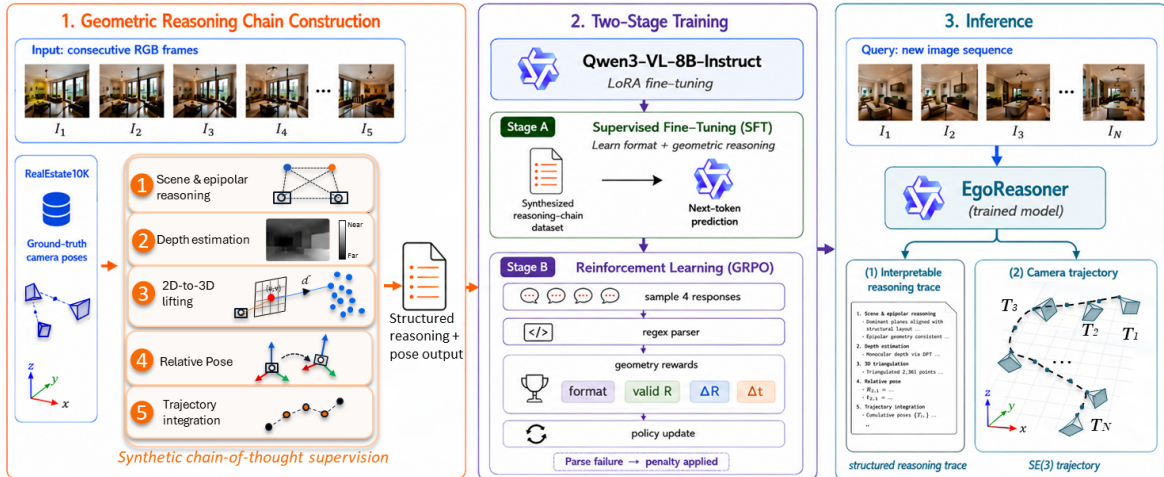


Figure 1. **EgoReasoner pipeline overview.** Given a sequence of consecutive RGB frames, the model produces a five-step geometric chain-of-thought—covering epipolar reasoning, depth estimation, 3D triangulation, relative pose estimation, and trajectory integration—before outputting the final numerical pose trajectory. The base VLM is first fine-tuned via SFT on synthesized geometric reasoning chains, then further refined by GRPO reinforcement learning with geometry-grounded reward signals.

176 ture, r_{val} (weight 0.1) rewards a valid rotation matrix, $r_{\Delta R}$
 177 (weight 0.35) rewards rotation accuracy, and $r_{\Delta t}$ (weight
 178 0.35) rewards translation accuracy. A strong parse fail-
 179 ure penalty of -1.0 is applied when the output cannot be
 180 parsed, ensuring that the RL stage does not degrade the
 181 structured reasoning format established by SFT.

182 The RL training loop proceeds as follows: given 5 con-
 183 secutive frames and intrinsic matrix \mathbf{K} , the model generates
 184 $n=4$ candidate responses per prompt; a regex-based parser
 185 extracts \mathbf{R} and \mathbf{t} ; the reward function scores each candidate;
 186 and GRPO updates the policy from relative reward differ-
 187 ences. Training uses a learning rate of 1×10^{-6} , KL coeffi-
 188 cient 0.01, temperature 0.7, and runs on 492 sequences for
 189 1 epoch.

190 4. Experimental Evaluation

191 4.1. Dataset

192 EgoReasoner is trained on two non-overlapping splits of the
 193 **RealEstate10K** dataset, a large-scale collection of indoor
 194 video sequences with ground-truth camera poses. The **SFT**
 195 **split** comprises 62,834 five-frame sequences from 30,702
 196 source videos; the SFT model is trained on the training por-
 197 tion of this split, and 60 held-out sequences are reserved
 198 for step-by-step chain-of-thought evaluation. The **RL split**
 199 comprises 74,027 usable five-frame sequences from 13,532
 200 source videos, with *zero* video overlap with the SFT split;
 201 492 sequences are used for RL training, and a **200-sequence**
 202 holdout is reserved for the main comparative evaluation on
 203 which all four methods are assessed. Together the two splits
 204 cover 684,455 frames in total.

Metric	Base Qwen3	Gemini 3.1 Pro [†]	Ours (SFT)	Ours (SFT+RL)
Parse Rate (%)	31.0	59.0	100.0	100.0
ATE median*	N/A	0.0113	0.0086	0.0082
RRE median (°) ↓	4.87	2.72	1.35	1.21
RRE mean (°) ↓	6.32	2.91	1.79	1.71
RTE median ↓	0.2607	0.1895	0.0308	0.0283
TDE median (°) ↓	80.33	18.76	16.45	15.47
TDE mean (°) ↓	69.63	37.33	25.25	23.69
TDE $\leq 15^\circ$ (%) ↑	8.1	43.2	45.9	48.6

Table 1. **Camera pose estimation on the 200-seq RealEstate10K holdout.** Bold: best per metric. [†]Gemini prompted with geometric hints. *Base Qwen3 ATE unavailable: model outputs only 1 image-pair pose (2 frames) vs. the required 5-frame trajectory, causing frame-count mismatch.

205 4.2. Baselines and Metrics

206 We compare against two baselines: **Base Qwen3-VL-8B-**
 207 **Instruct** [2] (the unmodified base VLM, no geometric fine-
 208 tuning) and **Gemini 3.1 Pro** [5] prompted with explicit geo-
 209 metric reasoning hints. Pose accuracy is evaluated with
 210 four metrics: Relative Rotation Error (**RRE**, degrees), Rel-
 211 ative Translation Error (**RTE**, normalized magnitude), Ab-
 212 solute Trajectory Error (**ATE**, RMSE after rigid alignment),
 213 and Translation Direction Error (**TDE**, degrees). For CoT
 214 evaluation we additionally report **parse rate**, epipolar er-
 215 ror, depth relative error, and 3D triangulation error at each
 216 reasoning step.

217 4.3. Comparative Evaluation

218 Three findings stand out from Table 1, where pose metrics
 219 are computed on the commonly parsed subset. First, the

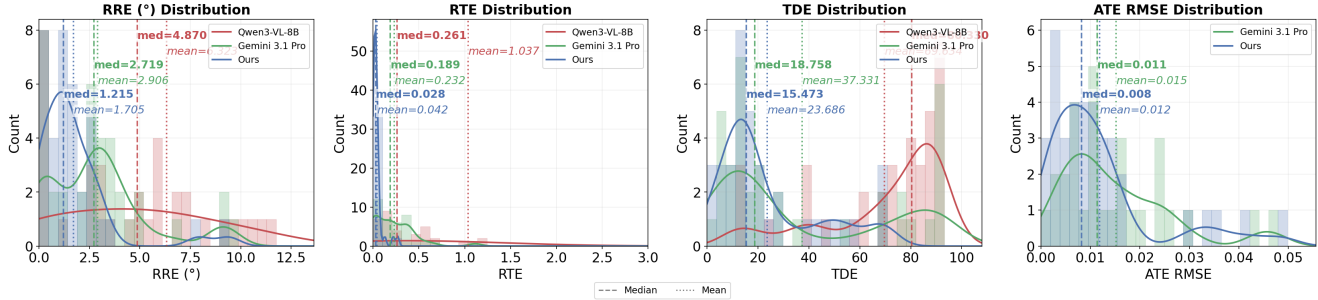


Figure 2. **Error distributions on the commonly parsed evaluation subset (200-sequence holdout).** EgoReasoner (SFT+RL) achieves consistently tighter distributions across all four metrics. Base Qwen3 ATE is absent due to frame-count mismatch (Table 1).

Step	Metric	Base Qwen3	Gemini 3.1 Pro	Ours
S1	Parse (%)	85.0	83.3	68.3
Epipolar	Epipolar (px) ↓	6.34	2.30	2.22
S2	Parse (%)	76.7	5.0	96.7
Depth	Depth Err ↓	3.261	1.779	0.760
S3	Parse (%)	70.0	1.7	81.7
3D Tri.	3D Err ↓	7.430	59.436	1.044
S4	Parse (%)	51.7	0.0	86.7
Rel. Pose	RTE mean ↓	0.7616	N/A	0.0741
S5	Parse (%)	76.7	3.3	60.0
Final	RRE med. (°) ↓	3.32	3.99	2.21
Traj.	RTE med. ↓	0.1013	0.1344	0.0596

Table 2. **Step-by-step CoT evaluation** on the 60-seq SFT holdout. Each step is evaluated independently. Bold: best per metric.

base model’s TDE median of 80.33° reveals near-random translation direction estimation, confirming that metric geometric reasoning does not emerge naturally in VLMs without geometric supervision—despite a reasonable parse rate. Second, Gemini 3.1 Pro, despite being a much larger model prompted with explicit geometric hints, parses only 59% of sequences and still yields poor translation estimation (RTE median 0.190), demonstrating that prompt engineering alone cannot substitute for geometric training. Third, EgoReasoner (SFT) achieves full parse reliability and surpasses Gemini by $6.1\times$ in RTE, confirming that geometric chain-of-thought supervision instills both structured output behavior and substantially improved metric reasoning; error distributions are shown in Figure 2.

4.4. Chain-of-Thought Step-by-Step Evaluation

To assess whether EgoReasoner learns genuine intermediate geometric reasoning rather than learning to produce superficially formatted outputs, we evaluate each step of the reasoning chain independently on the 60-sequence SFT evaluation set. Table 2 reports parse rates and accuracy metrics at each of the five reasoning steps.

The step-by-step results reveal a key structural finding: Gemini’s parse rate collapses to near zero at Steps 2–4, indicating its occasional successful final outputs in Table 1

are heuristic guesses rather than the product of structured geometric reasoning. In contrast, EgoReasoner maintains high parse rates across all intermediate steps and achieves substantial accuracy gains at each stage. The largest improvements occur at depth estimation ($4.3\times$ over the base model) and 3D triangulation ($7.1\times$), which are the geometric foundations that the final trajectory estimate depends on—explaining why intermediate CoT supervision translates directly into better end-to-end pose accuracy.

4.5. Effect of Reinforcement Learning

The GRPO RL stage yields consistent gains across all metrics (Table 1), with translation-related metrics improving more than rotation—RTE improves by 8.1% versus 10.4% for RRE—consistent with the reward function’s equal weighting of rotation and translation accuracy and RL’s capacity to exploit reward gradients that SFT supervision alone cannot provide. The distribution plots in Figure 2 corroborate this: the SFT+RL curves shift most noticeably in RTE and TDE tails, indicating RL corrects the harder-to-learn translation direction errors. Importantly, parse reliability is fully preserved across all evaluation sequences, confirming that the parse failure penalty effectively prevents the RL stage from sacrificing output format for score.

5. Conclusion

We presented **EgoReasoner**, an SFT+RL framework that trains VLMs to estimate camera pose trajectories via a five-step geometric chain-of-thought reasoning process. On a 200-sequence holdout set, EgoReasoner achieves a median RRE of 1.21° , a median RTE of 0.028, and $TDE \leq 15^\circ$ for 48.6% of samples—surpassing Gemini 3.1 Pro in translation accuracy by over $6\times$ and achieving $3.2\times$ higher parse reliability than an unmodified 8B VLM—demonstrating that metric geometric reasoning can be effectively instilled in VLMs through structured training. This capability is foundational for embodied agents in long-horizon interaction tasks, where a single VLM backbone must jointly handle semantic scene understanding and precise metric spatial self-localization.

282

References

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

- [1] 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*, 2023. 1, 2
- [2] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025. 2, 3
- [3] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13537–13547, 2024. 2
- [4] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2024. 2
- [5] Google DeepMind. Gemini 3.1 pro, 2025. Accessed: 2026-05-08. 1, 3
- [6] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 2
- [7] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. In *Advances in Neural Information Processing Systems*, 2023. 2
- [8] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022. 2
- [9] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 2
- [10] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, 2023. 1, 2
- [11] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1
- [12] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 1
- [13] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggg: Visual geometry grounded transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5294–5306, 2025. 1, 2
- [14] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vi-

- sion made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20697–20709, 2024. 1, 2
- [15] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang, Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua Shen, and Tong He. π^3 : Permutation-equivariant visual geometry learning. *arXiv preprint arXiv:2507.13347*, 2025. 2
- [16] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022. 2
- [17] Lihe Yang et al. Depth anything v3. *arXiv preprint*, 2025. 2

339

340

341

342

343

344

345

346

347

348

349

350

351

352