

# Improving Thinking Process in Visual Grounding via Free Thinking Rewards

Anonymous CVPR submission

## Abstract

001 We study egocentric visual intention grounding, where an  
002 assistant must infer and localize the object implied by a  
003 first-person view and an intention sentence without explicit  
004 object naming. Existing approaches either use a two-stage  
005 reasoning-then-grounding pipeline or apply reinforcement  
006 learning (RL) to train think-then-answer VLMs, but they op-  
007 timize only IoU-based rewards, which often leads to reward  
008 hacking, improving box accuracy while neglecting reason-  
009 ing quality. We introduce a label-free thinking-process  
010 reward that needs no human Chain-of-Thought labels or  
011 teacher models; it scores each sampled reasoning trace  
012 by how much it increases the likelihood of the correct an-  
013 swer, thus favoring reasoning that truly supports predic-  
014 tion. We also propose a data filtering strategy that selects  
015 informative easy-to-medium samples for RL using rollout  
016 error rate and reward variance. Together, these form a  
017 general recipe for process-aware RL finetuning of vision-  
018 language assistants for egocentric intention grounding. Our  
019 method achieves new state-of-the-art results on EgoInten-  
020 tion, boosting Precision@0.5 by +3.2 (3B) and +2.1 (7B)  
021 over strong Qwen2.5-VL baselines, and generalizes zero-  
022 shot to RefEgo-Int with +10.2 (3B) and +7.5 (7B).

## 023 1. Introduction

024 Recent reinforcement learning (RL) algorithms [8, 10,  
025 20, 26, 41] such as Group Relative Policy Optimiza-  
026 tion (GRPO [8]) enable think-then-answer vision-language  
027 models (VLMs), where the model first produces an ex-  
028 plicit reasoning trace and then commits to a final answer.  
029 For AI assistance applications [27, 32, 37], this capabil-  
030 ity allows VLMs to serve as a bridge from ambiguous hu-  
031 man requests to downstream embodied agents’ with action-  
032 able steps. Given a natural-language request, the system  
033 must first infer what the user actually needs, then determine  
034 which object in the scene should be acted on and where it  
035 is, and finally decompose the goal into actionable substeps  
036 for interacting with the environment. For example, when  
037 shown a cluttered kitchen and told “I need something to



Figure 1. An example of a correct bounding box with incorrect reasoning. VLM trained with IoU reward only (colored in red) mistakenly identifies the food packaging bag in the user’s hand as a handbag. The bounding box in the final answer correctly refers to the carrier paper box, but due to the flawed reasoning, the box is mislabeled as a handbag. The model with our free thinking reward corrects this mismatch (colored in green).

wipe the counter;” a useful system should infer the towel 038  
as the intended object, localize it in the scene, and only 039  
then plan how to pick it up. In this paper, we focus on this 040  
crucial upstream stage, which is visual intention grounding 041  
task [23, 27, 31]: given an intention sentence and a scene, 042

043 the model reasons about the intended target and localizes it  
044 as a foundation for subsequent action.

045 In visual grounding and related perception tasks [4, 13,  
046 14, 17], RL is typically plugged in as a post-training stage  
047 that optimizes a verifiable reward, most often derived from  
048 the overlap (Intersection over Union, IoU) between the pre-  
049 dicted and ground-truth bounding boxes. This setup is ef-  
050 fective because it directly aligns training with the evaluation  
051 metric, but the policy is only told whether the final box is  
052 good, not whether the intermediate reasoning was helpful  
053 or misleading. Recent think-then-answer approaches [4] try  
054 to address this by supervising the reasoning process itself,  
055 for example, with human-written rationales [25] or chains  
056 of thought (CoT) [35] distilled from stronger teacher mod-  
057 els [1, 8]; however, such labels are expensive to collect at  
058 scale, and can overfit models to narrow reasoning templates.

059 Our experiments further show that naïvely combining  
060 GRPO with IoU rewards only can lead to reward hack-  
061 ing [26], where the model improves box accuracy while  
062 its internal reasoning drifts away from semantically faith-  
063 ful explanations, as shown in Figure 1. Therefore, these  
064 observations motivate our goal: to use RL to shape both  
065 the final outcome and the intermediate reasoning without  
066 any extra reasoning labels, yielding a label-free, process-  
067 aware training signal. For each input, we generate multiple  
068 rollouts and assess how strongly each rollout supports the  
069 ground-truth answer by inspecting the model’s token-level  
070 likelihoods conditioned on that rollout (detailed in Sec. 3.2).  
071 We then aggregate these signals to reward thoughts that  
072 raise probability on the ground truth answer tokens. This  
073 process-aware signal complements IoU reward, mitigates  
074 reward hacking that overfits to final boxes.

075 Beyond making the model’s thinking more reliable, we  
076 also ask how to use RL under realistic compute budgets.  
077 Because each training step requires multiple rollouts per  
078 sample, naïvely running RL on the full supervised finetun-  
079 ing (SFT) corpus would multiply computation cost by the  
080 number of rollouts. We therefore adopt a static data filter-  
081 ing strategy that prunes the RL set before training, using  
082 the current model’s rollout statistics to estimate sample dif-  
083 ficulty *e.g.*, by rollout error rate and informativeness level  
084 *e.g.*, by the variance of rollout rewards across attempts. In  
085 contrast to prior data-selection practices [30, 39, 40] in re-  
086 asoning tasks [21, 22, 24] that favor hard unsolved cases [33],  
087 we find that focusing RL on informative samples of easy-  
088 to-medium difficulty stabilizes optimization and improves  
089 generalization. In practice, we perform RL on only a small  
090 fraction (about 5%) of the supervised fine-tuning samples  
091 while still achieving clear gains over the SFT baseline.

092 Bringing these two components together, we obtain a  
093 general training framework for visual intention grounding.  
094 We evaluate this framework on benchmarks such as EgoIn-  
095 tention [27] and RefEgo-Int [14] using mainstream open-

source backbones Qwen2.5-VL-3B/7B-Instruct [3]. Our  
method yields consistent gains in both reasoning accuracy  
and grounding accuracy: the former measures how often  
the model correctly predicts the intended target object cat-  
egory, while the latter is reported as Precision@0.5 based  
on IoU between predicted and ground-truth boxes, improv-  
ing by 1.2 and 3.2 points respectively. More generally, the  
proposed label-free thinking reward and data filtering are  
generic and can be applied to other vision-language tasks  
that emphasize high-quality thinking processes.

Our contributions can be summarized as follow:

1. We develop a label-free reward that encourages high-  
quality intermediate reasoning without requiring human  
annotations of thought traces or distillation from larger  
models.
2. We propose a data filtering strategy that jointly consid-  
ers sample difficulty and informativeness, and show that  
informative easy–medium samples benefits visual inten-  
tion grounding, in contrast to prior reasoning tasks that  
favor unsolved hard examples.
3. We provide a training recipe for reliable and effective vi-  
sual intention grounding, achieving consistent gains over  
Qwen2.5-VL-3B/7B baselines on the EgoIntention and  
RefEgo-Int benchmarks.

## 2. Related Work

### 2.1. RL for Vision Language Models

Recent industrial foundation models such as DeepSeek-  
R1 [8], Kimi k1.5 [29], and OpenAI’s o1 system [10] have  
popularized GRPO-style reinforcement learning (RL) as a  
scalable post-training recipe, establishing the think-then-  
answer paradigm and triggering a wave of RL research on  
vision–language models (VLMs). Most subsequent acade-  
mic efforts apply RL to tasks that naturally require long-  
chain reasoning, including mathematical and scientific QA,  
visual reasoning, and spatial QA, where multi-step CoT-  
style solutions are essential [12, 22, 24, 28, 34].

In contrast, RL on perception-centric tasks (*e.g.*, recog-  
nition, OCR, and visual grounding) remains compara-  
tively under-explored, though recent works demonstrate  
that think-then-answer VLMs still benefit from RL even  
when the final objective is a perceptual metric [4, 21, 26,  
38]. However, these methods typically treat the reasoning  
trace as a hidden byproduct and supervise only the final pre-  
diction, leaving the content and faithfulness of the thinking  
process unconstrained. Explicitly labeling human chains of  
thought is costly and can overfit models to narrow, human-  
written templates, while distilling thoughts from stronger  
models risks propagating teacher biases and harming gen-  
eralization. Our work instead keeps the think-then-answer  
framework but introduces a label-free thinking reward that  
directly evaluates how much a rollout supports the ground-

147 truth answer, thereby regularizing the reasoning process  
148 itself for perception-heavy tasks such as visual intention  
149 grounding.

## 150 2.2. Data Filtering for RL

151 Early instruction-tuning studies [5, 11, 15, 16, 19, 36, 44]  
152 already demonstrate that carefully filtered, task-specialized  
153 subsets can outperform using the full corpus: LTD Instruc-  
154 tion Tuning [5], DEITA [19], and IFD-style selection [43]  
155 all prune redundant or low-quality instructions and retain  
156 only a small fraction of the original data while matching  
157 or surpassing larger baselines [5, 16, 19]. As RL becomes  
158 the final post-training stage for large (vision-)language  
159 models, analogous data-filtering and curriculum strategies  
160 have been explored: Vision-G1 [39], QwenLong-L1 [30],  
161 and ThinkLite-VL [33] adopt various difficulty-aware or  
162 influence-based criteria to prioritize “hard-but-solvable” or  
163 unsolved examples, under the assumption that such in-  
164 stances provide the strongest learning signal for reasoning  
165 tasks (*e.g.*, math, logic, long-context QA) [30, 33, 39].

166 Our experiments reveal that aggressively emphasizing  
167 hard, frequently unsolved samples can be harmful here, as  
168 persistent failure yields near-zero IoU rewards and noisy,  
169 uninformative gradients. Instead, we propose a selection  
170 scheme that favors informative samples of easy-to-medium  
171 difficulty, leading to more stable optimization and more ef-  
172 fective RL fine-tuning for visual intention grounding than  
173 the “hard-unsolved-first” heuristics prevalent in reasoning-  
174 focused RL works.

## 175 3. Method

176 In this section, we present a practical recipe for training a  
177 reliable and effective visual grounding model. Section 3.1  
178 introduces the visual intention grounding task and reviews  
179 GRPO, an RL fine-tuning method for vision–language mod-  
180 els. Section 3.2 details our core contribution, the free think-  
181 ing reward, which scores each sampled reasoning trace by  
182 the increase it induces in the likelihood of the ground-truth  
183 answer. Finally, Section 3.3 describes our data-filtering  
184 strategy: we estimate per-sample difficulty and informative-  
185 ness, and, by performing RL updates on only 5% of the su-  
186 pervised fine-tuning data, we enable the model to handle  
187 visual intention grounding effectively.

### 188 3.1. Preliminary

#### 189 3.1.1. Visual Intention Grounding

190 For AI assistant, VLMs bridge ambiguous human re-  
191 quests to embodied action; the first step is Visual Inten-  
192 tion Grounding [23, 27], which maps an intention sen-  
193 tence to the specific target object and its location. For-  
194 mally, given an input image  $I$  and a human intention query  
195  $Q$ , the model must infer the underlying intent and answer

the target object and its location by predicting its bounding  
box  $(x_1, y_1, x_2, y_2)$  in answer tokens  $A$ . Beyond predicting  
the box, recent RL-finetuned reasoning models additionally  
produce a thinking trace  $T$  that explains how the VLM links  
the intention query  $Q$  to the target object within answer  $A$ .  
Concretely,  $T$  should give an accurate justification that links  
the intended action to the affordances and visible attributes  
of target object.

#### 3.1.2. GRPO for Vision Language Models

Compared with traditional visual grounding, Visual Inten-  
tion Grounding (VIG) first requires the model to parse the  
user’s intention and reason about the object that satisfies  
it. Earlier approaches [1, 18, 27, 31] often adopt a two-  
stage pipeline, infer the target object name and then localize  
it. With the recent adoption of GRPO-style reinforcement  
learning, VLMs can think before answering, allowing the  
model to transform an implicit intention into an explicit tar-  
get and its location in a single pass.

**Training objective.** We fine-tune the VLM with Group  
Relative Policy Optimization (GRPO) to obtain a reliable  
and effective AI assistant for intention grounding. GRPO  
is a rule-based RL algorithm for post-training LLMs/VLMs  
that removes the need for a learned critic by using group-  
relative rewards. For each input  $q$ , the old policy samples  
 $N$  candidate outputs  $\{o_i\}_{i=1}^N$ . Their rewards  $\{r_i\}_{i=1}^N$  are  
normalized to form advantages

$$A_i = \frac{r_i - \text{mean}\{r_1, \dots, r_N\}}{\text{std}\{r_1, \dots, r_N\}}, \quad (1)$$

which measure each candidate’s quality relative to its sib-  
lings. The GRPO objective is

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\theta) &= \mathbb{E}_{\{o_i\} \sim \pi_{\theta_{\text{old}}}(q)} \left[ \frac{1}{N} \sum_{i=1}^N \left\{ \min[s_1 \cdot A_i, s_2 \cdot A_i] \right. \right. \\ &\quad \left. \left. - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}] \right\} \right], \\ s_1 &= \frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)}, \\ s_2 &= \text{clip} \left( \frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)}, 1 - \epsilon, 1 + \epsilon \right). \end{aligned} \quad (2)$$

**Rewards.** We next detail the reward design for visual in-  
tention grounding task. Following prior RL-for-reasoning  
setups, we use a verifiable IoU accuracy reward and a for-  
mat reward that enforces valid output structure. Beyond  
these, we introduce a new free thinking reward that im-  
proves intention reasoning without any reasoning trace an-  
notations, this component is described in Section 3.2.

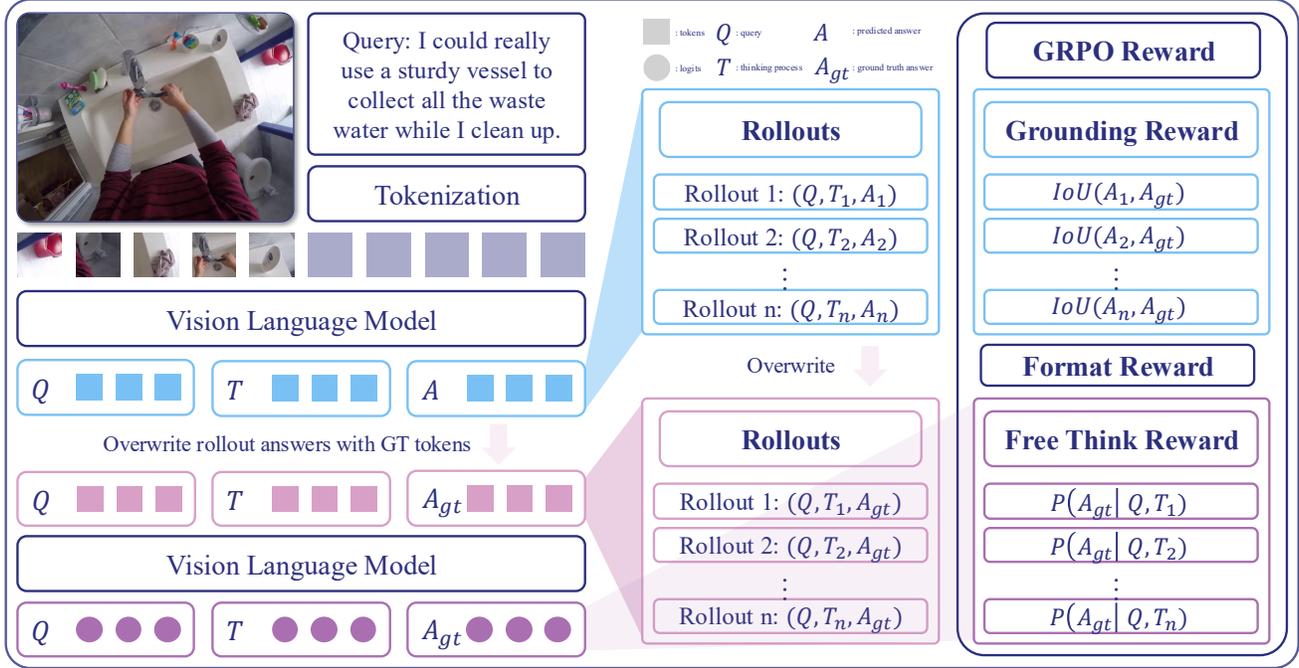


Figure 2. Overview of free thinking reward computation. Given an input image and human intention query, the VLM samples multiple think-then-answer rollouts  $(Q, T_k, A_k)$ , from which we derive a grounding reward via IoU with the ground-truth box and a format reward. We then overwrite  $A_k$  with the ground-truth answer  $A_{gt}$  and run the VLM in teacher-forced mode to obtain logits for the answer tokens, which define the free thinking reward  $P(A_{gt} | Q, T_k)$ . The grounding, free thinking, and format rewards are combined into the GRPO training signal.

233 **IoU reward.** We directly optimize spatial overlap:

234 
$$R_{\text{iou}}(I, Q, T) = \text{IoU}(\hat{b}, b^*) \in [0, 1], \quad (3)$$

235 where  $\hat{b}$  is the predicted box extracted from answer  $A$ ,  $b^*$  is  
236 the ground truth bounding box.

237 **Format reward.** We require a well-formed output con-  
238 taining a `<think>...</think>` trace  $T$  and a JSON box  
239 in `<answer>...</answer>`. Denote by  $\mathcal{F}$  the set of out-  
240 puts that satisfy the schema and yield a valid  $\hat{b}$ ; then

241 
$$R_{\text{fmt}}(I, Q, T) = \mathbb{I}\{(T, \hat{b}) \in \mathcal{F}\}. \quad (4)$$

242 **Overall reward.** Training uses a weighted sum

243 
$$R(I, Q, T) = \alpha R_{\text{iou}} + \beta R_{\text{fmt}} + \gamma R_{\text{think}}, \quad (5)$$

244 where  $R_{\text{think}}$  is our label-free thinking reward (Sec. 3.2) that  
245 evaluates the usefulness of the trace  $T$  for predicting the  
246 correct box, and  $\alpha, \beta, \gamma \geq 0$  are fixed weights, all set as 1  
247 in our framework.

### 248 3.2. Free Thinking Reward for VIG

249 Early RL approaches for VLM perception tasks predomi-  
250 nantly optimize verifiable rewards that directly align with  
251 evaluation metrics (e.g., IoU-based grounding accuracy),

252 which avoids relying on learned reward models but still  
253 leaves the model at risk of reward hacking, where the think-  
254 ing process becomes mismatched with the final answer.  
255 However, the content of the reasoning is critical: given an  
256 implicit intention query, both the inferred object semantics  
257 and the predicted box must be intention consistent, a cor-  
258 rect box paired with a wrong object interpretation is unre-  
259 liable. We therefore introduce a free thinking reward that  
260 improves the internal reasoning without any thinking trace  
261 annotations. It scores each sampled reasoning trace based  
262 on how much it raises the likelihood of the ground-truth  
263 answer, thereby rewarding reasoning paths that strengthen  
264 correct predictions. Given the input image  $I$  and intention  
265 query  $Q$ , the VLM first samples a trace  $T_{\text{pred}}$  and then an  
266 answer sequence  $A_{\text{pred}}$  that decodes to a box  $(x_1, y_1, x_2, y_2)$ .  
267 Under supervised finetuning setting, the standard answer  
268 cross-entropy is

$$\mathcal{L}_\theta(A | Q) = -\frac{1}{N} \sum_{i=1}^N \log P_\theta(w_i^A | Q, w_1^A, \dots, w_{i-1}^A), \quad (6)$$

269 where  $A = (w_1^A, \dots, w_N^A)$  is the ground-truth answer se-  
270 quence of length  $N$ ,  $w_i^A$  denotes its  $i$ -th token, and  $P_\theta(\cdot)$   
271 is the VLM's conditional next-token distribution parame-  
272

273 terized by  $\theta$  given the query  $Q$  and previously generated  
 274 answer tokens. While with RL-style generation the answer  
 275 tokens are conditioned on the sampled thinking trace,

$$276 P_\theta\left(w_i^{A_{\text{pred}}} \mid Q, T_{\text{pred}}, w_1^{A_{\text{pred}}}, \dots, w_{i-1}^{A_{\text{pred}}}\right). \quad (7)$$

277 **Hypothesis.** For a given sample, better (more intention-  
 278 faithful) traces make it easier for the model to produce the  
 279 correct answer.

280 Therefore, as shown in Figure 2, the likelihood of  $A_{\text{GT}}$   
 281 conditioned on the sampled trace can serve as an implicit  
 282 measure of trace quality, without any thinking process labels.  
 283 Thus, we evaluate the teacher-forced likelihood of the  
 284 ground-truth answer  $A_{\text{GT}}$  conditioned on the sampled trace  
 285 and define

$$\begin{aligned} 286 \mathcal{L}_\theta(A_{\text{GT}} \mid Q, T_{\text{pred}}) = \\ 287 - \frac{1}{N} \sum_{i=1}^N \log P_\theta\left(w_i^{A_{\text{GT}}} \mid Q, T_{\text{pred}}, w_1^{A_{\text{GT}}}, \dots, w_{i-1}^{A_{\text{GT}}}\right), \end{aligned} \quad (8)$$

288 which we convert to a scalar reward via  $r_{\text{think}} = -\mathcal{L}_\theta(A_{\text{GT}} \mid$   
 289  $Q, T_{\text{pred}})$ . Practically, we form the token sequence  
 290  $(Q_{\text{tok}}, T_{\text{tok}}, A_{\text{GT,tok}})$  by replacing the generated answer with  
 291 ground-truth tokens and perform a teacher-forced forward  
 292 pass to read logits; because  $T_{\text{pred}}$  is the only varying factor  
 293 per input, (8) acts as an implicit supervision signal on  
 294 the thinking process. Finally, we combine  $r_{\text{think}}$  with ver-  
 295 ifiable signals, which are IoU accuracy and output-format  
 296 validity, and optimize the policy using GRPO with group-  
 297 relative advantages, preferring traces that raises the likeli-  
 298 hood of ground truth answer without any thinking process  
 annotations.

### 299 3.3. Difficulty- and Variance-Aware Data Filtering

300 To efficiently train a VLM for egocentric visual intention  
 301 grounding, we identify training signals that are both learn-  
 302 able for the model and informative for policy improve-  
 303 ment. Concretely, we employ a think-then-answer VLM  
 304 and configure its exploration hyperparameters to stochasti-  
 305 cally sample multiple reasoning rollouts for each training  
 306 sample before RL finetuning.

307 Specifically, for each training sample  $i$  we execute  
 308  $R=8$  rollouts and compute an IoU-based verifiable re-  
 309 ward  $r_i^{(k)} \in [0, 1]$  for rollout  $k$ , with correctness indi-  
 310 cator  $c_i^{(k)} = \mathbb{I}[\text{IoU}_i^{(k)} \geq \tau]$  under a fixed threshold  $\tau$ .  
 311 We define a sample difficulty score as the number of cor-  
 312 rect rollouts  $d_i = \sum_{k=1}^R c_i^{(k)}$  (smaller  $d_i$  means harder),  
 313 and partition the dataset into three buckets: hard  $B_h =$   
 314  $\{i \mid d_i \in \{0, 1\}\}$ , medium  $B_m = \{i \mid d_i \in \{2, 3, 4, 5, 6\}\}$ ,  
 315 and easy  $B_e = \{i \mid d_i \in \{7, 8\}\}$ . In addition to diffi-  
 316 culty, we quantify a sample’s informativeness by the reward  
 317 spread  $\sigma_i = \text{Std}\left(\{r_i^{(k)}\}_{k=1}^R\right)$ , which upper-bounds the per-  
 318 sample advantage dispersion since for any baseline  $b_i$  the

rollout advantages  $A_i^{(k)} = r_i^{(k)} - b_i$  satisfy  $\text{Std}(\{A_i^{(k)}\}) =$  319  
 $\text{Std}(\{r_i^{(k)}\}) = \sigma_i$ . Intuitively, larger  $\sigma_i$  indicates a greater 320  
 separation between good and bad rollouts, yielding stronger 321  
 policy gradients for GRPO-style updates. Our filtering-and- 322  
 sampling scheme draws a portion vector  $\pi = (\pi_h, \pi_m, \pi_e)$  323  
 with  $\pi_h + \pi_m + \pi_e = 1$ , and samples  $N\pi_b$  items from bucket 324  
 $B_b$  according to a variance-biased distribution 325

$$P(i \mid B_b) = \frac{\sigma_i^\alpha}{\sum_{j \in B_b} \sigma_j^\alpha}, \quad \alpha > 0, \quad (9) \quad 326$$

where  $\alpha$  is a temperature controlling within-bucket sharp- 327  
 ness. To avoid misleading updates from unlearnable cases, 328  
 we discard degenerate samples with  $d_i = 0$  regardless of 329  
 high  $\sigma_i$ ; this prevents the optimizer from chasing high- 330  
 variance but consistently incorrect rollouts. In practice, 331  
 a balanced choice of  $\pi$  emphasizes  $B_m$  for learnability, 332  
 supplements with  $B_h$  for exploration and robustness, and 333  
 includes  $B_e$  for stabilization, while the within-bucket  $\sigma_i^\alpha$  334  
 weighting prioritizes samples whose rollouts provide larger 335  
 advantage gaps and thus more effective learning signals. 336

## 337 4. Experiments

### 338 4.1. Implementation Details

339 **Datasets.** Considering the most practical application sce- 340  
 nario, for example headmount camera, or smart glasses, we 341  
 focus on egocentric vision datasets [7]. We evaluate the 342  
 effectiveness of our frameworks on EgoIntention [27], a vi- 343  
 sual grounding benchmark pairing Ego4D images with in- 344  
 tention sentences. It contains two splits: *context* (typical 345  
 intents, e.g., “I need to sit down”  $\rightarrow$  chair) and *uncommon* 346  
 (atypical intents, e.g., “I need a boost to change the bulb” 347  
 $\rightarrow$  chair). To assess generalization, we construct RefEgo- 348  
 Int from RefEgo [14] by sampling frames (one frame per 349  
 video clip) and collecting context intention sentences fol- 350  
 lowing the EgoIntention dataset protocol. RefEgo-Int pre- 351  
 serves RefEgo key challenges: (i) frequent viewpoint mo- 352  
 tion, (ii) no-target cases, and (iii) long-tail object categories. 353  
 These challenges effectively extend the diversity of our test- 354  
 ing data, which in turn better assess model generalizability. 355

356 **Models.** Due to the strong vision-language under- 357  
 standing performance of Qwen2.5VL models [3], we use 358  
 Qwen2.5-VL-3B-Instruct as the main backbone and report 359  
 scalability with Qwen2.5-VL-7B-Instruct. GRPO settings 360  
 follow VLM-R1 [26]:  $N=8$  rollouts, temperature 0.9, it- 361  
 erations 1, KL ratio 0.04, and learning rate  $1e-6$ . The 362  
 grounding template and thinking prompts also match VLM- 363  
 R1. For the SFT baseline, we perform full finetuning with 364  
 LLaMA-Factory [42] on the EgoIntention training set for 365  
 1 epoch, whereas for RL finetuning we apply Low-Rank 366  
 Adaptation (LoRA) [9] for 500 update steps. 367

368 **Metrics.** We report *reasoning accuracy* (correct target 369  
 object category given image and intention) and *grounding* 370

Table 1. EgoIntention results split by Context / Uncommon and Overall. Metrics are  $P@0.5$  ( $\uparrow$ ) and  $mIoU$  ( $\uparrow$ ).

Method	Context		Uncommon		Overall	
	P@0.5	mIoU	P@0.5	mIoU	P@0.5	mIoU
Qwen-VL [2]	32.1	0.350	26.1	0.298	29.1	0.324
MiniGPT-v2 [6]	46.0	0.429	40.9	0.373	43.4	0.401
Qwen2.5-VL-3B-Instruct [3]	64.8	0.581	56.1	0.483	60.4	0.532
Ours-3B	<b>66.6</b>	<b>0.594</b>	<b>60.7</b>	<b>0.507</b>	<b>63.6</b>	<b>0.551</b>
Qwen2.5-VL-7B-Instruct [3]	67.2	0.594	61.8	0.514	64.5	0.554
Ours-7B	<b>68.9</b>	<b>0.605</b>	<b>64.3</b>	<b>0.527</b>	<b>66.6</b>	<b>0.566</b>

368 accuracy as Precision@0.5 using IoU between predicted  
369 and ground-truth boxes.

### 370 4.2. Main Results

371 We evaluated our method on EgoIntention context split  
372 and uncommon split, reporting overall accuracy as shown  
373 in Table 1. Compared with previous methods Qwen-VL  
374 and MiniGPT-v2, our base model Qwen2.5-VL-3B-Instruct  
375 achieves substantial improvements. Building on supervised  
376 finetuning models, for the 3B model, our method ground-  
377 ing precision on the context split increases by 1.8 points,  
378 while the uncommon split improves by 4.6 points. For 7B  
379 models, the context split increases by 1.7 points and the un-  
380 common split by 2.5 points. These results demonstrate that  
381 our framework is effective even when compared to strong  
382 supervised fine-tuning baselines at both 3B and 7B scales.  
383 The mIoU metric follows the same trend observed in precision@0.5.  
384

385 **Free thinking reward.** Here we use the 3B model with  
386 randomly selected (unfiltered) data for RL tuning to illus-  
387 trate the effectiveness of our thinking reward. As shown in  
388 Figure 3: First, RL finetuning with only the IoU reward  
389 causes a sharp early drop in performance due to reward  
390 hacking—object reasoning accuracy falls to 82.5% at step  
391 200. Second, adding our free-thinking reward stabilizes ob-  
392 ject reasoning accuracy at the start of RL tuning. Finally,  
393 the combined IoU + thinking reward yields consistent gains  
394 throughout training, reaching 86.4% at step 500, a 1.2% im-  
395 provement over the IoU-only model.

### 396 4.3. Ablation Study

397 **Does a logit-space reward alone help?** To understand  
398 whether a reward defined in logit space itself is helpful, we  
399 conduct an ablation where we compute a logit-based reward  
400 between the predicted answer distribution and a one-hot  
401 ground-truth label, without overwriting the answer tokens.  
402 Baseline model is trained with IoU and format rewards only.  
403 As shown in Table 2, adding this logit reward alone leads to  
404 a 0.6% drop in reasoning performance at 500 steps com-  
405 pared with the baseline. In contrast, our free thinking re-

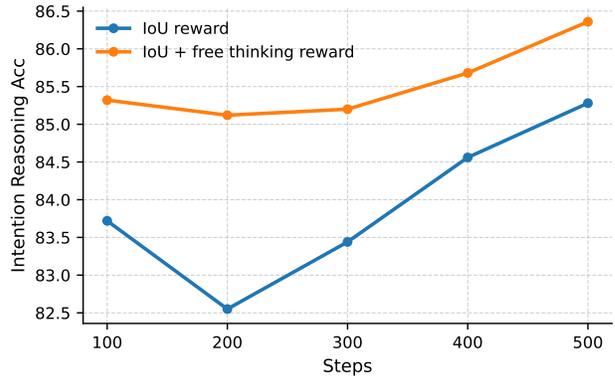


Figure 3. Intention reasoning accuracy (%) on EgoIntention context split. We show the intention reasoning accuracy trend during 500 steps RL finetuning.

Table 2. Intention reasoning accuracy (%) at different RL update steps on the EgoIntention context split. Step 0 corresponds to the supervised finetuned model before RL.

Method	0	100	200	300	400	500
IoU + format reward	86.6	83.7	82.5	83.4	84.5	85.2
+logit reward	86.6	85.2	84.6	85.2	83.7	84.6
+free thinking reward	86.6	85.3	85.1	85.2	85.7	<b>86.4</b>

ward replaces the predicted answer tokens with the ground-  
truth answer and measures how likely the model is to gener-  
ate the correct answer given the sampled thinking trace,  
*i.e.*, it constrains the probability  $P(A_{GT} | Q, T_{pred})$ . This  
explicitly ties the quality of the thinking process to answer  
correctness and yields a 1.2% improvement over the base-  
line model, demonstrating highly effectiveness.

**Data Filtering.** We broadly explore multiple data selection  
methods for RL on visual intention grounding tasks.

- Random. The baseline randomly samples the RL training set without any data filtering.
- Unsolved. Following prior RL-for-reasoning works, we keep only samples that the SFT model fails on, these form

Table 3. Overall grounding accuracy (Precision@0.5, %) across different training update steps on the EgoIntention dataset using five data filtering strategies.

Method	0	100	200	300	400	500
Random	60.4	62.7	62.8	62.9	63.0	63.1
Unsolved	60.4	60.0	60.3	60.1	59.2	59.5
Near threshold	60.4	62.8	62.5	62.6	62.4	62.9
Normal distribution	60.4	62.2	63.1	63.1	63.4	63.5
Ours	60.4	61.8	62.6	63.2	63.8	<b>64.0</b>

- 419 the hard unsolved RL set.
- 420 • Near threshold (IoU  $\approx 0.5$ ). We select samples whose
- 421 SFT IoU is near the 0.5 decision boundary, assuming they
- 422 are most correctable by RL.
- 423 • Normal distribution sampling. Using our difficulty
- 424 (wrong rollout count) and informativeness (reward vari-
- 425 ance) estimates, we sample with a normal distribution
- 426 centered at 4 wrong rollouts (midpoint of 0–8).
- 427 • Ours. Motivated by the drop with hard–unsolved data, we
- 428 downweight hard cases and exclude all-wrong samples,
- 429 focusing RL on informative easy–medium instances.

430 As shown in Table 3, we first observe that directly borrow-

431 ing strategies from previous RL work on reasoning tasks

432 does not work well. Using unsolved hard samples for RL

433 finetuning cannot even boost the final grounding precision,

434 performance drops from 60.4 at step 0 to 59.5 at step 500.

435 Compared with the baseline method without any data filter-

436 ing strategy, we see a clear gap at step 500, dropping from

437 63.1 to 59.5. For training samples with IoU near 0.5, there is

438 a consistent slight decrease across updating steps compared

439 with baseline results. For the normal distribution approach

440 based on the difficulty level we estimated for the training

441 set, our strategy prefers samples with higher informativ-

442 eness levels within each difficulty bin. We observe better

443 results compared with the baseline method. At step 500,

444 the overall grounding performance on EgoIntention outper-

445 forms the baseline method by 0.4. Finally, we exclude sam-

446 ples with all wrong rollouts and use a sample ratio of 1:6:3

447 for hard, medium, and easy level samples. This achieves the

448 best overall performance compared with baseline.

#### 449 4.4. Visualization

450 As shown in Figure 4, we visualize paired outputs to assess

451 how the free thinking reward shapes reasoning and ground-

452 ing. Human intention queries are shown in yellow. For clar-

453 ity, the system prompt is omitted. Red denotes the thinking

454 process and answer from the model trained only with the

455 IoU reward, while green shows the output from the model

456 trained with both IoU and free thinking rewards. Both mod-

457 els are finetuned on the same randomly sampled data

458 The six examples on the left illustrate cases where the

459 IoU-only model predicts correct bounding boxes but incor-

Table 4. Visual grounding accuracy (Precision@0.5, %) at different training update steps on the RefEgo-Int dataset.

Method	0	100	200	300	400	500
Naïve RL-3B	14.2	22.2	22.3	22.6	23.0	23.4
Ours-3B	14.2	23.2	23.3	23.5	24.5	<b>24.4</b>
Ours-7B	18.0	23.5	24.9	24.7	24.9	<b>25.5</b>

460 rect object reasoning. Models trained solely with IoU re-

461 ward often confuse visually or functionally similar objects

462 even when localization is correct. For instance, in the first

463 sample of the first row, the user asks for a container to hold

464 leftover paint; the model correctly localizes the paint can

465 but labels it as a bottle. In the second sample of the sec-

466 ond row, the user asks for a container for drinking coffee;

467 the model’s trace suggests it searches for a cup or mug, and

468 it correctly bounds the red and white cup on the counter,

469 but mislabels it as a mug. Other typical confusions include

470 bucket vs. tray, bucket vs. bowl, knife vs. scissors, and

471 sponge vs. tissue paper.

472 The two examples on the right show failures in both rea-

473 soning and grounding. In the first row, the model mentions

474 a ladder in its thought process, suggesting the user needs

475 one to reach the top shelf. However, it is a hallucination

476 that no ladder appears in the image. The model marks the

477 area around the stool as a ladder. In the second row, the

478 model’s reasoning fixates on the television’s prominent vi-

479 sual features. The thinking process doesn’t mention finding

480 a communication device like a cell phone or telephone. Our

481 model explicitly reasons that a phone is needed for mak-

482 ing calls or confirming plans, and successfully locates the

483 telephone in the image.

#### 484 4.5. Zero-shot Evaluation

485 To verify that our RL framework boosts the model’s gener-

486 alization ability, we test our models on RefEgo-Int, a chal-

487 lenging egocentric visual grounding dataset. First, we evalu-

488 ate the base models, which are supervised fine-tuned on

489 the EgoIntention dataset. The results are low, the 3B model

490 achieves 14.2% precision@0.5, and the 7B model achieves

491 18.0%. This is because the models are zero-shot tested

492 without training on RefEgo dataset, and because of chal-

493 lenges like blurry objects caused by view motion, long-tail

494 object categories, and cases where no ground truth bound-

495 ing boxes exist in the image. The model should be aware

496 of how to reject a query when there is no target object in

497 the image. However, the model trained on EgoIntention

498 doesn’t have such abilities. After applying our RL fine-

499 tuning framework on the EgoIntention dataset, we zero-shot

500 test our RL models on RefEgo-Int. We observe significant

501 improvements compared with the SFT baseline at step 0 for

502 both 3B (+10.2% P@0.5) and 7B (+7.5% P@0.5) models.

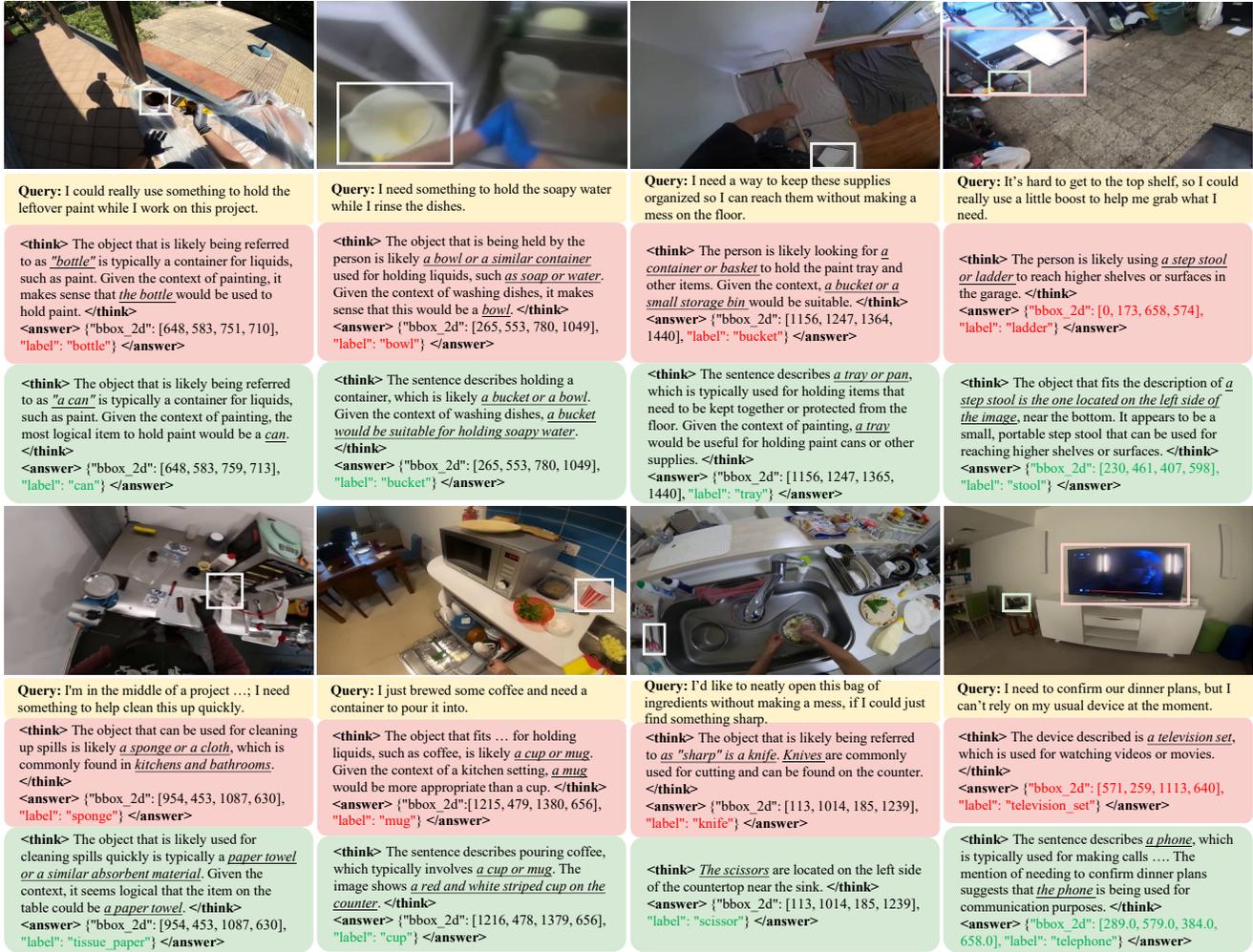


Figure 4. Visualization comparing the baseline VLM trained with IoU reward only and our model trained with both IoU and free thinking rewards. The human intention query is colored in yellow. The baseline model’s reasoning and answers are shown in red, while ours are shown in green, with key reasoning phrases underlined. For the six examples on the left, both methods produce the same grounding, shown with gray boxes. Please zoom in for better view.

503 We also compare our RL framework with the IoU reward  
 504 only RL method, and we achieve better performance, boost-  
 505 ing approximately 1.0% P@0.5.

## 506 5. Conclusion

507 In this paper, we identify a mismatch between incorrect  
 508 thinking processes and correct bounding boxes in visual  
 509 intention grounding. To address this, we propose a free  
 510 thinking reward that scores reasoning traces by how much  
 511 they boost the likelihood of the correct answer. Beyond this  
 512 thinking reward, we tailor a difficulty- and informativeness-  
 513 aware data filtering strategy for visual intention ground-  
 514 ing, which focuses GRPO training on informative easy-to-  
 515 medium samples. Extensive experiments on EgoIntention  
 516 and RefEgo-Int demonstrate that our approach consistently

improves both grounding accuracy and intention reasoning,  
 achieving state-of-the-art results over strong Qwen2.5-VL  
 baselines, and we hope our framework will inspire more  
 think-answer consistent VLMs in the community.

**Future Work.** Our label-free thinking reward is task-  
 agnostic, the final output format is not limited to bounding  
 box coordinates. We plan to extend it to reasoning tasks  
 such as visual question answering and vision-language-  
 action models. Moreover, beyond checking the correctness  
 of the inferred target object, we plan to design dedicated  
 metrics to evaluate the quality of the thinking process. For  
 example, we will check thinking process consistency across  
 rollouts during training and whether overthinking occurs.

530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
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

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2, 3
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond, 2023. arXiv:2308.12966 [cs]. 6
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 2, 5, 6
- [4] Sule Bai, Mingxing Li, Yong Liu, Jing Tang, Haoji Zhang, Lei Sun, Xiangxiang Chu, and Yansong Tang. UniVG-R1: Reasoning Guided Universal Visual Grounding with Reinforcement Learning, 2025. arXiv:2505.14231 [cs]. 2
- [5] Hao Chen, Yiming Zhang, Qi Zhang, Hantao Yang, Xiaomeng Hu, Xuetao Ma, Yifan Yanggong, and Junbo Zhao. Maybe Only 0.5% Data is Needed: A Preliminary Exploration of Low Training Data Instruction Tuning, 2023. arXiv:2305.09246 [cs]. 3
- [6] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. MiniGPT-v2: large language model as a unified interface for vision-language multi-task learning, 2023. arXiv:2310.09478 [cs]. 6
- [7] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Abraham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kotur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4D: Around the World in 3,000 Hours of Egocentric Video. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18973–18990, New Orleans, LA, USA, 2022. IEEE. 5
- [8] 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. 1, 2
- [9] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 5
- [10] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024. 1, 2
- [11] Naman Jain, Tianjun Zhang, Wei-Lin Chiang, Joseph E. Gonzalez, Koushik Sen, and Ion Stoica. LLM-Assisted Code Cleaning For Training Accurate Code Generators, 2023. arXiv:2311.14904 [cs]. 3
- [12] Chaoya Jiang, Yongrui Heng, Wei Ye, Han Yang, Haiyang Xu, Ming Yan, Ji Zhang, Fei Huang, and Shikun Zhang. VLM-R\$^3\$: Region Recognition, Reasoning, and Refinement for Enhanced Multimodal Chain-of-Thought, 2025. arXiv:2505.16192 [cs]. 2
- [13] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. ReferItGame: Referring to Objects in Photographs of Natural Scenes. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 787–798, Doha, Qatar, 2014. Association for Computational Linguistics. 2
- [14] Shuhe Kurita, Naoki Katsura, and Eri Onami. RefEgo: Referring Expression Comprehension Dataset from First-Person Perception of Ego4D. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15168–15178, Paris, France, 2023. IEEE. 2, 5
- [15] Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. Superfiltering: Weak-to-Strong Data Filtering for Fast Instruction-Tuning, 2024. arXiv:2402.00530 [cs]. 3
- [16] Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning, 2024. arXiv:2308.12032 [cs]. 3
- [17] You Li, Heyu Huang, Chi Chen, Kaiyu Huang, Chao Huang, Zonghao Guo, Zhiyuan Liu, Jinan Xu, Yuhua Li, Ruixuan Li, and Maosong Sun. Migician: Revealing the Magic of Free-Form Multi-Image Grounding in Multimodal Large Language Models, 2025. arXiv:2501.05767 [cs]. 2
- [18] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55. Springer, 2024. 3
- [19] Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What Makes Good Data for Alignment? A Com-

643	prehensive Study of Automatic Data Selection in Instruction Tuning, 2024. arXiv:2312.15685 [cs].	3	701
644			702
645	[20] Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visual-RFT: Visual Reinforcement Fine-Tuning, 2025. arXiv:2503.01785 [cs].	1	703
646			704
647			705
648			706
649	[21] Yan Ma, Linge Du, Xuyang Shen, Shaoxiang Chen, Pengfei Li, Qibing Ren, Lizhuang Ma, Yuchao Dai, Pengfei Liu, and Junjie Yan. One RL to See Them All: Visual Triple Unified Reinforcement Learning, 2025. arXiv:2505.18129 [cs].	2	707
650			708
651			709
652			710
653	[22] Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, Ping Luo, Yu Qiao, Qiaosheng Zhang, and Wenqi Shao. MM-Eureka: Exploring the Frontiers of Multimodal Reasoning with Rule-based Reinforcement Learning, 2025. arXiv:2503.07365 [cs].	2	711
654			712
655			713
656			714
657			715
658			716
659			717
660	[23] Mengxue Qu, Yu Wu, Wu Liu, Xiaodan Liang, Jingkuan Song, Yao Zhao, and Yunchao Wei. RIO: A Benchmark for Reasoning Intention-Oriented Objects in Open Environments, 2023. arXiv:2310.17290 [cs].	1, 3	718
661			719
662			720
663			721
664	[24] Gabriel Sarch, Snigdha Saha, Naitik Khandelwal, Ayush Jain, Michael J Tarr, Aviral Kumar, and Katerina Fragkiadaki. Grounded reinforcement learning for visual reasoning. <i>arXiv preprint arXiv:2505.23678</i> , 2025.	2	722
665			723
666			724
667			725
668	[25] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. <i>Advances in Neural Information Processing Systems</i> , 37:8612–8642, 2024.	2	726
669			727
670			728
671			729
672			730
673			731
674	[26] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jijia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. VLM-R1: A Stable and Generalizable R1-style Large Vision-Language Model, 2025. arXiv:2504.07615 [cs].	1, 2, 5	732
675			733
676			734
677			735
678			736
679	[27] Pengzhan Sun, Junbin Xiao, Tze Ho Elden Tse, Yicong Li, Arjun Akula, and Angela Yao. Visual intention grounding for egocentric assistants. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 2512–2522, 2025.	1, 2, 3, 5	737
680			738
681			739
682			740
683			741
684	[28] Huajie Tan, Yuheng Ji, Xiaoshuai Hao, Minglan Lin, Pengwei Wang, Zhongyuan Wang, and Shanghang Zhang. Reason-RFT: Reinforcement Fine-Tuning for Visual Reasoning, 2025. arXiv:2503.20752 [cs].	2	742
685			743
686			744
687			745
688	[29] Kimi Team, Angang Du, Bofei Gao, Bawei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. <i>arXiv preprint arXiv:2501.12599</i> , 2025.	2	746
689			747
690			748
691			749
692			750
693	[30] Fanqi Wan, Weizhou Shen, Shengyi Liao, Yingcheng Shi, Chenliang Li, Ziyi Yang, Ji Zhang, Fei Huang, Jingren Zhou, and Ming Yan. QwenLong-L1: Towards Long-Context Large Reasoning Models with Reinforcement Learning, 2025. arXiv:2505.17667 [cs].	2, 3	751
694			752
695			753
696			754
697			755
698	[31] Wenxuan Wang, Yisi Zhang, Xingjian He, Yichen Yan, Zijia Zhao, Xinlong Wang, and Jing Liu. Beyond Literal Descriptions: Understanding and Locating Open-World Objects		756
699			757
700			
	Aligned with Human Intentions, 2024. arXiv:2402.11265 [cs].	1, 3	
	[32] Xin Wang, Taein Kwon, Mahdi Rad, Bowen Pan, Ishani Chakraborty, Sean Andrist, Dan Bohus, Ashley Feniello, Bugra Tekin, Felipe Vieira Frujeri, et al. Holoassist: an egocentric human interaction dataset for interactive ai assistants in the real world. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 20270–20281, 2023.	1	
	[33] Xiyao Wang, Zhengyuan Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin, Furong Huang, and Lijuan Wang. SoTA with Less: MCTS-Guided Sample Selection for Data-Efficient Visual Reasoning Self-Improvement, 2025. arXiv:2504.07934 [cs].	2, 3	
	[34] Zhenhailong Wang, Xuehang Guo, Sofia Stoica, Haiyang Xu, Hongru Wang, Hyeonjeong Ha, Xiushi Chen, Yangyi Chen, Ming Yan, Fei Huang, et al. Perception-aware policy optimization for multimodal reasoning. <i>arXiv preprint arXiv:2507.06448</i> , 2025.	2	
	[35] 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. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.	2	
	[36] Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS: Selecting Influential Data for Targeted Instruction Tuning, 2024. arXiv:2402.04333 [cs].	3	
	[37] Jingkang Yang, Shuai Liu, Hongming Guo, Yuhao Dong, Xiamengwei Zhang, Sicheng Zhang, Pengyun Wang, Zitang Zhou, Binzhu Xie, Ziyue Wang, et al. Egolife: Towards egocentric life assistant. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 28885–28900, 2025.	1	
	[38] En Yu, Kangheng Lin, Liang Zhao, Jisheng Yin, Yana Wei, Yuang Peng, Haoran Wei, Jianjian Sun, Chunrui Han, Zheng Ge, Xiangyu Zhang, Daxin Jiang, Jingyu Wang, and Wenbing Tao. Perception-R1: Pioneering Perception Policy with Reinforcement Learning, 2025. arXiv:2504.07954 [cs].	2	
	[39] Yuheng Zha, Kun Zhou, Yujia Wu, Yushu Wang, Jie Feng, Zhi Xu, Shibo Hao, Zhengzhong Liu, Eric P. Xing, and Zhiting Hu. Vision-G1: Towards General Vision Language Reasoning with Multi-Domain Data Curation, 2025. arXiv:2508.12680 [cs].	2, 3	
	[40] Bolin Zhang, Jiahao Wang, Qianlong Du, Jiajun Zhang, Zhiying Tu, and Dianhui Chu. A Survey on Data Selection for LLM Instruction Tuning, 2025. arXiv:2402.05123 [cs].	2	
	[41] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative Verifiers: Reward Modeling as Next-Token Prediction, 2025. arXiv:2408.15240 [cs].	1	
	[42] Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyuan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System</i>		

758 *Demonstrations*), Bangkok, Thailand, 2024. Association for  
759 Computational Linguistics. 5  
760 [43] Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun,  
761 Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan  
762 Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and  
763 Omer Levy. LIMA: Less Is More for Alignment. 3  
764 [44] Terry Yue Zhuo, Arnel Zebaze, Nitchakarn Suppattarachai,  
765 Leandro von Werra, Harm de Vries, Qian Liu, and Niklas  
766 Muennighoff. Astraios: Parameter-Efficient Instruction Tun-  
767 ing Code Large Language Models, 2024. arXiv:2401.00788  
768 [cs]. 3