

# 000 TRACEABLE EVIDENCE ENHANCED VISUAL 001 GROUNDED REASONING: EVALUATION AND METHOD 002

003 **Anonymous authors**

004 Paper under double-blind review

## 005 ABSTRACT

006 Models like OpenAI-o3 pioneer visual grounded reasoning by dynamically ref-  
007 erencing visual regions, just like human “thinking with images”. However, no  
008 benchmark exists to evaluate these capabilities holistically. To bridge this gap, we  
009 propose **TreeBench** (Traceable Evidence Evaluation Benchmark), a diagnostic  
010 benchmark built on three principles: (1) *focused visual perception* of subtle targets  
011 in complex scenes, (2) *traceable evidence* via bounding box evaluation, and (3)  
012 *second-order reasoning* to test object interactions and spatial hierarchies beyond  
013 simple object localization. Prioritizing images with dense objects, we initially  
014 sample 1K high-quality images from SA-1B, and incorporate eight LMM experts  
015 to manually annotate questions, candidate options, and answers for each image.  
016 After three stages of quality control, **TreeBench** consists of 405 challenging vi-  
017 sual question-answering pairs, even the most advanced models struggle with this  
018 benchmark, where none of them reach 60% accuracy, *e.g.*, OpenAI-o3 scores only  
019 54.87. Furthermore, we introduce **TreeVGR** (Traceable Evidence Enhanced Visual  
020 Grounded Reasoning), a training paradigm to supervise localization and reasoning  
021 jointly with reinforcement learning, enabling accurate localizations and explainable  
022 reasoning pathways. Initialized from Qwen2.5-VL-7B, it improves V\* Bench  
023 (+16.8), MME-RealWorld (+12.6), and **TreeBench** (+13.4), proving traceability is  
024 key to advancing vision-grounded reasoning. The code and data will be released.  
025

## 026 1 INTRODUCTION

027 Recent breakthroughs in Large Language Models (LLMs) reasoning, such as OpenAI-o1 (OpenAI,  
028 2024b) and DeepSeek-R1 (Guo et al., 2025a) with remarkable test-time scaling properties, have  
029 motivated researchers to explore reasoning for Large Multimodal Models (LMMs) (Huang et al.,  
030 2025; Wei et al., 2025a;b; Chen et al., 2025). These models are typically remarkable in their  
031 mathematical and scientific reasoning, particularly through *text-space* reasoning. However, they  
032 exhibit critical limitations when applied to perception-heavy tasks (Jiang et al., 2025) or general  
033 multimodal benchmarks (Wang et al., 2024c), primarily due to accumulated language bias from their  
034 exclusive reliance on textual reasoning pathways. A paradigm shift toward *visual grounded reasoning*  
035 emerged with models like OpenAI-o3 (OpenAI, 2025), which is able to “think with images” by  
036 dynamically referencing and amplifying task-relevant regions during reasoning, resulting in *image-  
037 text interleaved* reasoning pathways. Yet, despite growing interest, the community currently lacks  
038 comprehensive evaluation benchmarks for assessing these capabilities.

039 Classical benchmarks like POPE (Li et al., 2023c), MMBench (Liu et al., 2023b), SEED-Bench (Li  
040 et al., 2023a), and MMMU (Yue et al., 2024) usually overlook fine-grained localization and verifiable  
041 reasoning chains. Others (Wu & Xie, 2024; Zhang et al., 2024a; Wang et al., 2025f; Dong et al.,  
042 2024; Wang et al., 2025b;a; Zhang et al., 2024b) partially address localization but lack traceability  
043 or complex reasoning: V\* Bench (Wu & Xie, 2024) is restricted to simple spatial queries (*e.g.*, “Is  
044 A left of B?”) and risks data contamination with COCO-derived images (Lin et al., 2014); MME-  
045 RealWorld (Zhang et al., 2024a), HR-Bench (Wang et al., 2025f), and document benchmarks (Biten  
046 et al., 2022; Mathew et al., 2021; Liu et al., 2023c) support high-resolution inputs but lack traceable  
047 evidence and second-order reasoning such as perspective shifts. In short, these benchmarks fail  
048 to adequately evaluate three key elements central to visual grounded reasoning: nuanced visual  
049 grounding, traceable multi-step reasoning, and dynamic cross-modal interaction through *interleaved*  
050 box-text reasoning pathways.

054	Attributes	Material	Physical State	Object Retrieval	OCR-Integrated QA
055					
056	<b>Question:</b> What is the girl wearing while sitting on the chair in the center of the image, which is partially obscured by a tall street light? <b>A. Pink shoes</b> B. White skirt C. Light blue skirt   D. Black shoes 	<b>Question:</b> What the materials for the bottles on the bike? <b>A. Plastic</b>   B. White Glass C. Bronze-aware product <b>D. Insulation material</b>	<b>Question:</b> What is the condition of the rear cargo door on the small white box-shaped truck parked in the leftmost lane? <b>A. Fully closed and latched</b>   B. Fully open upward C. Half-open (partially raised)  D. Missing entirely	<b>Question:</b> What is attached to the top of the hat in the middle of the image? <b>A. Flag</b>   <b>B. Sunglasses</b>  C. Flower D. Grass ring	<b>Question:</b> Recognize the question and options in the image and answer it. <b>A. Flag</b>   <b>B. Sunglasses</b>  C. Flower D. Grass ring
057	Target Instances: 	Target Instances: 	Target Instances: 	Target Instances: 	Target Instances: 
058	Perspective Transformation	Ordering	Contact and Occlusion	Spatial Containment	Comparison
059					
060	<b>Question:</b> From the perspective of the woman seated in a wheelchair, what is the relative direction of the signboard with "PROGRAMS"? <b>A. Front left</b> B. Front right  C. Left rear D. Right rear 	<b>Question:</b> Among the signs with text, which one (counting from left to right) has the most text? <b>A. The first one</b> B. The second one  C. The third one  D. The fourth one	<b>Question:</b> Considering the soccer player in red (number 22) in the foreground, is his left foot occluded with the soccer ball? <b>A. Yes, they are in direct contact</b>  B. No, they are separated by a gap  <b>C. It cannot be determined</b> D. They partially overlap	<b>Question:</b> Which of the following objects is inside the building on the right? <b>A. The Buddhist monk in orange</b>  B. The person in blue on the right  C. The couple in the middle of the image D. The black car in the middle	<b>Question:</b> Considering the relative distances in the image, which object is closer to the police officer holding the ice cream cone? <b>A. The tree on the bus stop</b> B. The black bollard on the sidewalk  C. The police van parked on the side of the road D. The bus stop sign 
061	Target Instances: 	Target Instances: 	Target Instances: 	Target Instances: 	Target Instances: 
062					
063					
064					
065					
066					
067					
068					
069					
070					
071					

Figure 1: Qualitative examples from **TreeBench** for each discipline. Each question requires focused visual parsing on mere objects, and some even request second-order reasoning beyond precise localization. Moreover, the bounding boxes of all target instances are provided, ensuring a traceable evaluation. All these questions are challenging, as OpenAI-o3 (OpenAI, 2025) and Gemini-2.5-Pro (DeepMind, 2025b) *cannot* answer them correctly simultaneously.

To bridge this gap, we propose **TreeBench** (Traceable Evidence Evaluation Benchmark), designed around three foundational principles essential for evaluating true “thinking with images” capabilities:

- **Focused Visual Perception.** It evaluates a model’s ability to identify *subtle* targets within cluttered, real-world scenes using detailed, precise, and unique textual descriptions, which requires hierarchical scene understanding and the discrimination of extremely similar distractors.
- **Traceable Evidence.** It not only evaluates the final accuracy but also pioneers quantifiable evaluation of reasoning chains, resulting in an explainable, reliable, and transparent evaluation.
- **Vision-Centric Second-Order Reasoning Capabilities.** It moves beyond simple object localization and primitive “what/where” queries. It focuses on complex physical interactions between objects (such as contact and occlusion), as well as spatial containment (inside/outside, above/below) and relative relationships with perspective transformation.

To conduct **TreeBench**, we sample 1K images from SA-1B (Kirillov et al., 2023), prioritizing images with dense objects, as SA-1B (Kirillov et al., 2023) offers high-resolution, real-world scenes with many small and varied objects, making it particularly suitable for evaluating visual grounded reasoning. Subsequently, 8 experts with solid technical backgrounds are involved in hand-crafted annotation for 10 sub-tasks, as demonstrated in Figure 1. In particular, we present a semi-automated pipeline. Each of OpenAI-o3 (OpenAI, 2025) and Gemini-2.5-Pro (DeepMind, 2025b) is required to create three distinct questions belonging to a specific subtask, accompanied by multiple-choice options and the respective correct answers. Subsequently, experts curated or replaced these to ensure quality and difficulty. We additionally incorporate a cross-verification stage for further quality control. Finally, **TreeBench** incorporates 405 high-quality and extremely challenging VQA pairs with accurate bounding boxes of target instances. A comprehensive comparison between **TreeBench** and other related benchmarks is provided in Table 1. Key advantages are:

- **Annotation Quality.** Unlike benchmarks relying on LMM-generated labels such as MMT-Bench (Ying et al., 2024) and SEED-Bench (Li et al., 2023a), our expert-driven process ensures correctness and extreme difficulty. However, relying on models would inevitably introduce significant noise, compromising the quality of the annotations. On the contrary, our **TreeBench** is manually designed by 8 LMM experts, ensuring the annotation correctness and ensuring the difficulty of all questions.

108  
109  
110  
111 Table 1: Comparison between benchmarks related to “thinking with images”. **TreeBench** features  
112 traceable evidence annotations, as well as high input resolution and challenging questions.  
113  
114

Benchmark	Resolution	Traceable Evidence Annotation	Mean Area of Target Objects (↓)	Qwen2.5-VL-72B Performance (↓)
V* Bench	2,246×1,583	✗	—	85.9
HR-Bench-4K	4,023×3,503	✗	—	79.3
HR-Bench-8K	5,727×4,430	✗	—	76.0
MME-RealWorld	2,076×1,434	✗	—	62.9
<b>TreeBench</b>	2,152×1,615	✓	<b>3.05%</b>	<b>42.2</b>

115  
116  
117  
118 • **Small Target Objects.** All questions in **TreeBench** focus on extremely small objects in complex  
119 real-world scenes, where target instances occupy an average of 3.05% of the image.

120  
121 • **Traceable Evidence Evaluation.** Our **TreeBench** provides bounding box annotations of each  
122 target instance. It not only evaluates the final answer, but also reveals the quality of *intermediate*  
123 *reasoning steps*. Those predicted bounding boxes serve as a window into its process, helping to  
124 diagnose the source of errors, *i.e.*, whether the model misunderstood the question or failed to  
125 locate the relevant object.

126 • **Task Difficulty.** While models approach saturation (>90%) on benchmarks like V\* Bench (Wu  
127 & Xie, 2024), even open-sourced state-of-the-art performers like Qwen2.5-VL-72B (Bai et al.,  
128 2025a) achieve only 42.2 on our **TreeBench**, implying a large potential improvement.

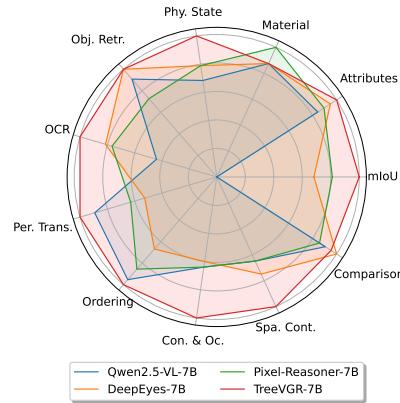
129 Beyond evaluation, we further introduce **TreeVGR** (Trace-  
130 able Evidence for Visual Grounded Reasoning), a training  
131 paradigm enhancing localization-driven visual reasoning.  
132 Previous attempts like (Wang et al., 2025e; Zheng et al.,  
133 2025b; Cao et al., 2025; Fan et al., 2025; Shao et al., 2024a;  
134 Qi et al., 2024; Su et al., 2025; Liu et al., 2025a) solely  
135 supervise final answers and neglect intermediate region-of-  
136 interest generation processes. It becomes hard to quantify  
137 the actual contribution of the “grounding-then-answering”  
138 framework. On the contrary, we propose **TreeVGR**, a  
139 novel training methodology emphasizing traceable evi-  
140 dence through reinforcement learning (RL), which *explicitly*  
141 supervises bounding box generation.

142 Building on RL with conventional accuracy-based and  
143 formatting rewards, **TreeVGR** leverages a novel *dual* IoU  
144 reward to ensure both precision and recall in localizing  
145 the ground-truth bounding boxes for each target instance.  
146 To implement this, we curate 37K samples for RL train-  
147 ing, each comprising an image, a question, an answer,  
148 and corresponding *bounding box annotations* for all target  
149 instances. Empirically, initialized from  
150 Qwen2.5-VL-7B (Bai et al., 2025a), **TreeVGR** brings significant  
151 improvements on various benchmarks, *i.e.*, +16.8 on V\* Bench (Wu & Xie, 2024), +12.6 on MME-RealWorld-Lite (Zhang et al.,  
152 2024a), and +13.4 on our **TreeBench**. Moreover, as illustrated in Figure 2, compared with related  
153 approaches, our **TreeVGR** enables traceable and explainable  
154 reasoning pathways with more accurate  
155 localizations (mIoU), and finally contributes to bootstrapped overall performance.

156 In conclusion, **TreeBench** pioneers the evaluation of how models “think with images”, while  
157 **TreeVGR** establishes a blueprint for training them. Together, they significantly advance the depth  
158 and utility of multimodal reasoning assessment with *traceable evidence*.

## 2 RELATED WORK

159 **Large Multimodal Models.** Initial breakthroughs in Large Multimodal Models (LMMs), such  
160 as Flamingo (Alayrac et al., 2022) and BLIP-2 (Li et al., 2023b), achieved this by integrating  
161 visual features into the LLM backbone via cross-attention mechanisms. A significant shift towards  
efficiency emerged with LLaVA (Liu et al., 2023a), which introduced a much more efficient approach.



162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
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  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
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  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
13

162 It projects visual features from a pre-trained encoder (*e.g.*, CLIP (Radford et al., 2021)) directly  
 163 into the LLM’s semantic space using a simple two-layer MLP. This paradigm of feature projection  
 164 catalyzed rapid advancement. Subsequent research has dramatically scaled LMM capabilities and  
 165 tackled increasingly complex tasks (Wang et al., 2025c; Liu et al., 2024a; Li et al., 2024; Wang  
 166 et al., 2025d; Bai et al., 2025a; Wang et al., 2024b; Lei et al., 2025; Zhu et al., 2025; Wu et al.,  
 167 2024; Wang et al., 2024a). A critical frontier has been handling high-resolution inputs. Models  
 168 like LLaVA-NeXT (Liu et al., 2024a) and InternVL-1.5 (Chen et al., 2024b) adopt any resolution  
 169 strategy. Qwen2-VL (Wang et al., 2024b) and Qwen2.5-VL (Bai et al., 2025a) introduce multimodal  
 170 Rotary Position Embedding (mROPE) to support arbitrary resolution inputs. Beyond resolution,  
 171 scaling pretraining with high-quality data is also vital, as demonstrated by InternVL3 (Zhu et al.,  
 172 2025). Collectively, these models represent the state-of-the-art, forming robust baselines for diverse  
 173 real-world multimodal applications. Our work builds upon these advances by leveraging their strong  
 174 native visual grounding capabilities. However, existing LMMs do not naturally perform an explicit  
 175 “grounding-then-answering” process, often resulting in misaligned or incomplete responses. By  
 176 explicitly modeling this sequential process, our approach ensures more accurate and interpretable  
 177 answers through grounded reasoning.

178 **Reasoning LMMs.** The groundbreaking reasoning capabilities of LLMs, exemplified by systems  
 179 like OpenAI-o1 (OpenAI, 2024b) and DeepSeek-R1 (Guo et al., 2025a) have motivated efforts to  
 180 extend similar competencies to multimodal settings using reinforcement learning (RL) (Sutton et al.,  
 181 1998). Early approaches primarily focused on equipping LMMs to solve complex math and science  
 182 problems involving image inputs (Huang et al., 2025; Wei et al., 2025a;b; Chen et al., 2025). Other  
 183 approaches (Shen et al., 2025; Liu et al., 2025b; Bai et al., 2025b) directly adopt GRPO (Shao et al.,  
 184 2024b) to open-ended visual grounding. Moreover, some attempts (Liu et al., 2024b; Mondal et al.,  
 185 2024; Shao et al., 2024a; Qi et al., 2024) focus on regions-of-interest localization before actually  
 186 answering the question. A recent milestone, OpenAI-o3 (OpenAI, 2025), advanced multimodal  
 187 reasoning by enabling dynamic image manipulation, *e.g.*, cropping and zooming into regions of  
 188 interest, to emulate human-like “thinking with images.” Subsequent research has sought to replicate  
 189 this capability through diverse strategies: constructing SFT data (Wang et al., 2025e), vanilla RL (Fan  
 190 et al., 2025), framing grounding as a function (Zheng et al., 2025b), decoupling grounding and  
 191 answering (Cao et al., 2025), multi-task reinforcement learning (Liu et al., 2025a), and curiosity-  
 192 driven reasoning (Su et al., 2025). Critically, these RL-based methods supervise *only* the final answer.  
 193 In contrast, our **TreeVGR** emphasizes *traceable evidence* during RL training, *i.e.*, supervising  
 194 generated bounding boxes to ensure precise localization throughout the reasoning process. By doing  
 195 so, **TreeVGR** enables more transparent, reliable, and fine-grained control over the reasoning pipeline.

196 **Benchmarks for LMMs.** Current benchmarks lack comprehensive evaluation of multimodal models’  
 197 ability to “think with images”, a capability demanding three core competencies: (1) focused visual  
 198 perception (identifying small targets in large scenes), (2) traceable evidence (evaluating generated  
 199 bounding boxes for explainability), and (3) second-order reasoning (deriving insights *beyond* precise  
 200 instance localization). Some benchmarks may *partially* satisfy the first condition. While some  
 201 benchmarks address isolated aspects, critical gaps persist. Classical benchmarks like POPE (Li et al.,  
 202 2023c), MMBench (Liu et al., 2023b), SEED-Bench (Li et al., 2023a), and MMMU (Yue et al., 2024)  
 203 usually overlook fine-grained localization and verifiable reasoning chains. V\* (Wu & Xie, 2024)  
 204 evaluates detailed attributes and spatial relationships (*e.g.*, “Is A left of B?”) but relies on COCO-  
 205 derived images (Lin et al., 2014), introducing high contamination risk. MME-RealWorld (Zhang et al.,  
 206 2024a) and HR-Bench (Wang et al., 2025f) support high-resolution inputs but lack traceable evidence,  
 207 and their questions often become easy when grounded precisely. Crucially, no benchmark integrates  
 208 all three requirements, particularly the need for complex reasoning conditional on precise grounding,  
 209 *e.g.*, perspective transform: “*From the perspective of person A, what is the relative direction of*  
 210 *object B?*”. To bridge this gap, we propose **TreeBench**, the first benchmark designed explicitly for  
 211 “thinking with images” with *traceable*, multistep evaluation. Beyond accuracy, **TreeBench** assesses:  
 212 (1) region quality, *i.e.*, faithfulness of generated regions-of-interest in visual reasoning chains, and  
 213 (2) second-order reasoning, *i.e.*, capabilities requiring inference *beyond* localization. State-of-the-art  
 214 models, Gemini-2.5-Pro (DeepMind, 2025b) and OpenAI-o3 (OpenAI, 2025), perform poorly on  
 215 **TreeBench** (<60%), underscoring its rigor and the unmet challenges in multimodal reasoning.

216 **3 TREEBENCH**  
 217

218 **TreeBench** is designed to address a critical gap in multimodal evaluation by establishing the first  
 219 comprehensive benchmark for assessing “thinking with images” capabilities. Specifically, it mainly  
 220 evaluates (1) the ability of identifying small target objects with long, detailed, and unique text  
 221 captions in large, complex, and real-world scenes, (2) the explainability of reasoning pathways and  
 222 traceable evidence, and (3) second-order reasoning beyond precise localization. Our **TreeBench**  
 223 systematically evaluates 10 core competencies through 405 distinct questions, organized into two  
 224 progressive protocols, *i.e.*, “Perception” and “Reasoning”, with representative examples in Figure 1.  
 225 In the following, we provide a detailed exploration of task definitions. The annotation pipeline and  
 226 the final statistics of **TreeBench** can be found in Appendix B and Appendix C, respectively.  
 227

228 **Perception** evaluates the model’s ability to accurately “see” and “identify” specific content, which  
 229 is one of the basic capabilities of directly extracting and interpreting visual information from every  
 230 detail of the provided image. These tasks primarily evaluate *first-order* visual reasoning capabilities,  
 231 where correct answers usually depend on the accurate localization of target questions (*e.g.*, objects,  
 232 regions, or text) and directly recognize their explicit attributes *without* requiring higher-level logical  
 233 inference or abstract conceptualization. It includes:

- 234 1. **Attributes** evaluates the ability to identify and describe specific visual properties (*e.g.*, color,  
 235 shape, material, or precise classification) of objects or elements within images, particularly  
 236 requiring attention to fine details, subtle distinctions, and accurate recognition of small-scale or  
 237 context-dependent features.
- 238 2. **Material** measures the ability to analyze and distinguish material properties (*e.g.*, texture, surface  
 239 finish, composition, or physical state) through visual cues such as light reflection, transparency,  
 240 wear patterns, or microscopic structural characteristics, requiring precise reasoning about tactile  
 241 qualities and material-specific visual indicators.
- 242 3. **Physical State** assesses the ability to assess structural integrity (*e.g.*, damage, wear, or breakage),  
 243 detect positional states (*e.g.*, open/closed, bent/straight), and interpret age-related features (*e.g.*,  
 244 freshness, decay) through precise analysis of visual cues like cracks, alignment anomalies,  
 245 lighting/shadow patterns, or contextual degradation markers.
- 246 4. **Object Retrieval** probes the ability to interpret linguistically complex, spatially explicit de-  
 247 scriptions and map them to visually subtle or contextually embedded targets in images, testing  
 248 the integration of natural language understanding, spatial grounding, and discriminative object  
 249 recognition under high specificity constraints.
- 250 5. **OCR-Integrated Question-Answering** evaluates the ability to extract text-based questions  
 251 and answer options from images, requiring seamless integration of OCR, natural language  
 252 understanding, and multimodal alignment to produce accurate responses grounded in both  
 253 textual and visual modalities.

254 **Reasoning** evaluates the ability to analyze and infer meaningful conclusions beyond recognition.  
 255 These tasks demand *second-order* visual reasoning capabilities, where correct answers require not  
 256 only accurate localization but also higher-level cognitive operations over aggregated visual evidence.  
 257 Precise perceptual grounding is just the first step for these tasks. It includes:

- 258 1. **Perspective Transform** measures the capacity to perform viewpoint transformations (*e.g.*,  
 259 aligning viewer-centric and agent-centric frames of reference) and interpret spatial relations  
 260 under mirror-reversed or perspective-shifted conditions, testing the ability to disambiguate  
 261 directional relationships that depend on the visualized entity’s orientation rather than the image’s  
 262 literal pixel layout.
- 263 2. **Ordering** evaluates the ability to analyze linearly ordered arrangements of objects (*e.g.*, left-to-  
 264 right, front-to-back, or depth-based sequences) and resolve ordinal relationships by integrating  
 265 spatial context with discriminative feature recognition, requiring precise localization within  
 266 continuous layouts and contextual comparison of positional cues (*e.g.*, adjacency, centrality, or  
 267 extremity) to answer questions dependent on sequential alignment and relative placement.
- 268 3. **Contact and Occlusion** measures the ability to analyze physical interactions between multiple  
 269 objects (*e.g.*, direct contact, occlusion layers, or shadow-based overlaps) and resolve ambiguities  
 270 in object identification by leveraging spatial dependencies, requiring precise parsing of contact

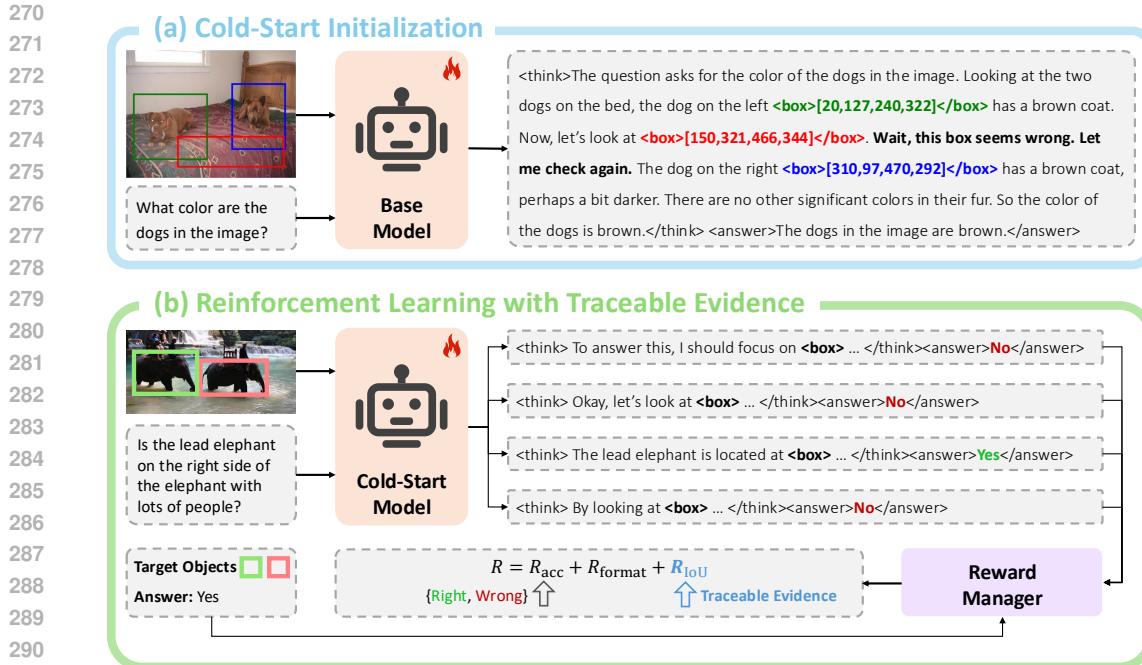


Figure 3: Training pipeline of **TreeVGR**, including (a) a cold-start initialization stage and (b) a reinforcement learning with traceable evidence post-training stage.

cues (*e.g.*, alignment, boundary fusion), occlusion boundaries (*e.g.*, partial/full coverage, layer stacking), and contextual constraints to answer questions that hinge on understanding how objects physically coexist and obscure one another in complex scenes.

4. **Spatial Containment** benchmarks the ability to analyze hierarchical spatial relationships (*e.g.*, containment, surface attachment, or regional boundaries) by parsing visual cues like object boundaries, spatial context, and contextual containment rules, requiring precise interpretation of containment hierarchies, surface dependencies, and regional constraints to resolve questions dependent on explicit spatial membership rather than isolated positional attributes.
5. **Comparison** assesses to compare attributes across multiple objects (*e.g.*, distance, size, color) and resolve spatial or perceptual differences, requiring precise parsing of attribute discrimination and contextual distance estimation to answer questions demanding explicit comparison of visually co-present entities.

## 4 TREEVGR

In this section, we introduce our **TreeVGR**. Specifically, we leverage the native grounding capabilities of pre-trained LMMs and unlock *visual grounded reasoning* capabilities, *i.e.*, localizing regions-of-interest first and answering the question next, through a two-stage training pipeline shown in Figure 3, *i.e.*, cold initialization introduced in Section 4.1 and reinforcement learning with traceable evidence elaborated in Section 4.2.

Notably, our **TreeVGR** does *not* require actually replaying cropped images as previous approaches (Wang et al., 2025e; Zheng et al., 2025b; Su et al., 2025) do, as *text-space grounding* is already effective. It leads to much more efficient training and inference procedures.

### 4.1 COLD-START INITIALIZATION

While end-to-end reinforcement learning (RL) has demonstrated validity by (Zheng et al., 2025b) for visual grounded reasoning (VGR) tasks, its practical deployment remains hindered by *extreme computational demands*. Specifically, DeepEyes-7B (Zheng et al., 2025b) requests RL training on 47K samples across 32 episodes, a process requiring 32 H100 (80GB) GPUs operating continuously for 48 hours. Such resource intensity creates barriers to broader accessibility.

To address these limitations, we investigate a computationally efficient alternative. Initial attempts revealed significant training inefficiencies when applying direct RL to VGR: models required extensive iterations to autonomously identify task-relevant visual regions before generating answers. This bottleneck motivates our adoption of a cold initialization strategy as illustrated in Figure 3a. Specifically, we introduce a supervised fine-tuning (SFT) phase using a curated dataset comprising multimodal samples: each sample includes an image, a question, reasoning trajectories with corresponding bounding boxes, and a final answer. This structured initialization ensures VGR capabilities are established prior to RL. Details of data construction and optimization can be found in Appendix E.1.

## 4.2 REINFORCEMENT LEARNING WITH TRACEABLE EVIDENCE

We proceed to reinforcement learning (RL) to refine reasoning trajectories through *traceable evidence supervision* as demonstrated in Figure 3b. Specifically, the bounding boxes generated are evaluated using a box intersection-over-union (IoU) reward, a precise and interpretable metric that measures the alignment between predicted and ground-truth regions. This reward ensures explicit accountability to human-annotated visual evidence, guiding the policy toward spatially accurate and logically coherent reasoning pathways.

**Reward Design.** The total reward consists of three parts: an accuracy reward  $R_{\text{acc}} \in \{0, 1\}$ , a formatting reward  $R_{\text{format}} \in \{0, 1\}$ , and a *dual* Intersection-over-Union (IoU) reward  $R_{\text{IoU}} \in [0, 1]$ :

$$R = R_{\text{acc}} + R_{\text{format}} + R_{\text{IoU}}, \quad (1)$$

where the accuracy reward assesses whether the final answer is correct. We utilize exact-matching for multiple-choice questions, and leverage an online reward model, *i.e.*, Qwen2.5-72B-Instruct (Team, 2024), to judge whether the prediction is correct given the question and the ground-truth answer. The formatting reward ensures the reasoning process and the final answer must be enclosed between `<think>` and `</think>`, and `<answer>` and `</answer>`, respectively. The *dual* IoU reward measures the quality of predicted boxes against ground-truths. Specifically, for  $N$  predicted bounding boxes  $\{\hat{b}_i\}_{i=1}^N$ , where  $\hat{b}_i = [x_1^i, \hat{y}_1^i, \hat{x}_2^i, \hat{y}_2^i]$  and  $M$  ground-truths  $\{b_k\}_{k=1}^M$ , where  $b_k = [x_1^k, y_1^k, x_2^k, y_2^k]$ , the *dual* IoU is an average of a *recall* term and a *precision* term.

$$R_{\text{IoU}} = \frac{1}{2}(R_{\text{IoU}}^{\text{R}} + R_{\text{IoU}}^{\text{P}}), \quad (2)$$

where the  $R_{\text{IoU}}^{\text{R}}$  indicates the *recall* and  $R_{\text{IoU}}^{\text{P}}$  means the *precision*. Specifically, the *recall* term ensures that each ground-truth bounding box  $b_k$  is matched with at least one prediction.

$$R_{\text{IoU}}^{\text{R}} = \frac{1}{M} \sum_{k=1}^M \text{IoU} \left[ \{\hat{b}_i\}_{i=1}^N, b_k \right], \quad (3)$$

where  $\text{IoU} \left[ \{\hat{b}_i\}_{i=1}^N, b_k \right] = \max_i \text{IoU}(\hat{b}_i, b_k)$  indicates the maximum IoU between *all* predictions  $\{\hat{b}_i\}_{i=1}^N$  and *each* ground-truth  $b_k$ . Maximizing this term ensures each ground-truth  $b_k$  is matched with *at least* one prediction. However, we empirically find that the policy model tends to *enumerate all possible boxes* to obtain a larger recall. Therefore, we introduce a dual term, *i.e.*,  $R_{\text{IoU}}^{\text{P}}$ , to ensure the *precision* and discourage “empty” boxes that do not match with any ground-truths:

$$R_{\text{IoU}}^{\text{P}} = \frac{1}{N} \sum_{i=1}^N \text{IoU} \left[ \{b_k\}_{k=1}^M, \hat{b}_i \right]. \quad (4)$$

Similarly,  $\text{IoU} \left[ \{b_k\}_{k=1}^M, \hat{b}_i \right] = \max_k \text{IoU}(b_k, \hat{b}_i)$  indicates the maximum IoU between *all* ground-truths  $b_k$  and *each* prediction  $\{\hat{b}_i\}_{i=1}^N$ . Maximizing this term encourages each prediction  $\hat{b}_i$  to be matched with *at least* one ground-truth. Therefore, simultaneous optimization of *both* recall and precision eliminates the need for exhaustive enumeration of bounding boxes, thereby contributing to more accurate reasoning pathways. Details of data and optimization can be found in Appendix E.2.

## 5 EXPERIMENTS

**Baselines.** We include four state-of-the-art private models, GPT-4o-1120 (OpenAI, 2024a) and o3-0416 (OpenAI, 2025) from OpenAI, and Gemini-2.5-Flash-0520 (DeepMind, 2025a) and Gemini-2.5-Pro-0605 (DeepMind, 2025b) from Google. Additionally, representative open-source general

378  
 379 Table 2: Selected results of different models on **TreeBench**. Evaluations of open-source general  
 380 models are implemented using VLMEvalKit (Duan et al., 2024), while evaluations of visual grounded  
 381 reasoning models are conducted by us. <sup>†</sup>Reasoning pathways of o3 (OpenAI, 2025) are unavailable,  
 382 and thus traceable evaluations are *not* valid. Best performances for open-source models are highlighted  
 383 in **bold**. *Our TreeVGR-7B achieves comparable performance with InternVL3-78B (Zhu et al., 2025).*

		Overall	mIoU	Attributes	Material	Phy. State	Obj. Relt.	OCR	Per. Trans.	Ordering	Con. & Oc.	Spa. Cont.	Comparison
				Perception				Reasoning					
				Private Models									
<b>Open-source General Models</b>													
Gemini-2.5-Flash-0520	45.9	–	48.3	53.9	69.6	68.8	75.0	15.3	19.3	56.1	72.4	43.2	
GPT-4o-1120	46.9	–	51.7	61.5	65.2	43.8	69.1	18.8	38.6	48.8	72.4	43.2	
Gemini-2.5-Pro-0605	54.1	–	51.7	61.5	56.5	75.0	83.8	20.0	36.8	65.9	86.2	54.6	
o3-0416	54.8	– <sup>†</sup>	69.0	69.2	65.2	68.8	79.4	22.4	38.6	61.0	86.2	50.0	
<b>Open-source Visual Grounded Reasoning Models</b>													
LLaVA-OneVision-7B	37.3	–	55.2	53.8	56.5	50.0	32.4	21.2	22.8	41.5	72.4	36.4	
LLaVA-OneVision-72B	40.5	–	62.1	53.8	65.2	62.3	36.8	12.9	28.1	53.7	65.5	<b>47.7</b>	
Qwen2.5-VL-7B	37.0	–	55.2	53.8	56.5	62.5	27.9	20.0	35.1	39.0	44.8	43.2	
Qwen2.5-VL-72B	42.2	–	65.5	<b>69.2</b>	56.5	56.3	48.5	11.8	33.3	51.2	72.4	38.6	
InternVL3-8B	38.8	–	51.7	<b>69.2</b>	56.5	56.3	33.7	21.2	24.6	39.0	72.4	43.2	
InternVL3-78B	46.4	–	62.1	61.5	52.2	<b>68.8</b>	52.9	16.5	33.3	61.0	<b>86.2</b>	45.5	
DeepEyes-7B	37.5	30.0	62.1	53.8	65.2	68.8	51.5	11.8	24.6	36.6	51.7	<b>47.7</b>	
Pixel-Reasoner-7B	39.0	35.7	58.6	61.5	65.2	50.0	48.5	14.1	31.6	39.0	44.8	40.9	
<b>TreeVGR-7B</b>	<b>50.4</b>	<b>44.0</b>	<b>65.5</b>	53.8	<b>82.6</b>	<b>68.8</b>	<b>63.3</b>	<b>22.4</b>	<b>36.8</b>	<b>61.0</b>	69.0	45.5	
Δ v.s. Qwen2.5-VL-7B	↑ 13.4	–	↑ 11.7	– 0.0	↑ 26.1	↑ 6.3	↑ 35.4	↑ 2.2	↑ 1.7	↑ 22.0	↑ 24.2	↑ 2.3	

403 models are incorporated, including LLaVA-OneVision series (Li et al., 2024), Qwen2.5-VL series (Bai  
 404 et al., 2025a), and InternVL3 series (Zhu et al., 2025). Furthermore, two very recent visual grounded  
 405 reasoning models are also included, *i.e.*, DeepEyes (Zheng et al., 2025b) and Pixel-Reasoner (Su  
 406 et al., 2025), as both of them follow a “grounding then answering” pipeline, with the capability of  
 407 “thinking with images”. Evaluations are mainly conducted on **TreeBench**, V\* Bench (Wu & Xie,  
 408 2024), HR-Bench (Wang et al., 2025f), and MME-RealWorld-Lite (Zhang et al., 2024a).

409 **Results on TreeBench.** Table 2 presents per per-category performance of different models. Overall,  
 410 OpenAI’s o3-0416 (OpenAI, 2025), the state-of-the-art visual grounded reasoning model, demon-  
 411 strates the strongest perception abilities, as expected. Larger models usually contribute to better  
 412 performance. Notably, our **TreeVGR-7B** even achieves comparable performance with InternVL3-  
 413 78B (Zhu et al., 2025), demonstrating the effectiveness of the visual grounded reasoning pipeline.  
 414 Moreover, compared with visual grounded reasoning models, our **TreeVGR** not only achieves a  
 415 higher overall performance, but also obtains a larger mIoU, indicating its effectiveness in precisely  
 416 localizing target objects. More in-depth analysis on **TreeBench** can be found in Appendix D.

417 **Results on High-Resolution Benchmarks.** In Table 3, **TreeVGR** achieves open-source state-of-  
 418 the-art on V\* Bench (Wu & Xie, 2024). On HR-Bench (Wang et al., 2025f) and MME-RealWorld-  
 419 Lite (Zhang et al., 2024a) illustrated in Table 3 and Table 4, respectively, our **TreeVGR** brings  
 420 significant improvements over our base model, Qwen2.5-VL-7B (Bai et al., 2025a). Results on other  
 421 general benchmarks can be found in Appendix F.1.

422 **Ablation Studies.** The core contribution of **TreeVGR** is the *traceable* training pipeline, where  $R_{IoU}$  is  
 423 incorporated in conventional RL training. The effectiveness of this design is ablated in Appendix F.2.

## 425 6 CONCLUSION

427 This paper introduces **TreeBench**, a benchmark designed to rigorously evaluate visual grounded  
 428 reasoning (VGR) or “thinking with images” in large multimodal models, and **TreeVGR**, a two-stage  
 429 training framework that enhances VGR methods through traceable evidence supervision.

431 **TreeBench** addresses critical gaps in existing benchmarks by focusing on three principles: focused  
 432 visual perception (identifying subtle targets in cluttered scenes), traceable evidence (quantifiable

432 Table 3: Comparison with state-of-the-art alternatives on V\* Bench (Wu & Xie, 2024) and  
 433 HRBench (Wang et al., 2025f). All results are self-collected. Best performances of visual grounded  
 434 reasoning models are highlighted in **bold**.

	V* Bench			HR-Bench-4K			HR-Bench-8K		
	Overall	Attr.	Spatial	Overall	Single	Cross	Overall	Single	Cross
<b>Private Models</b>									
GPT-4o-1120	66.0	–	–	–	–	–	–	–	–
o3-0416	95.7	–	–	–	–	–	–	–	–
<b>Open-source General Models</b>									
LLaVA-OneVision-7B	70.7	73.0	60.5	64.3	74.8	53.8	59.8	65.3	54.3
LLaVA-OneVision-72B	73.8	80.9	63.2	66.3	76.5	56.0	60.9	68.8	53.0
InternVL3-8B	72.3	73.0	71.1	70.8	79.3	62.3	62.0	64.3	59.8
InternVL3-78B	76.4	75.7	77.6	75.5	84.5	66.5	67.3	71.8	62.8
Qwen2.5-VL-7B	74.3	77.4	69.7	72.1	88.8	55.5	68.8	83.5	54.0
Qwen2.5-VL-72B	84.8	90.8	80.9	79.4	88.8	70.0	76.3	84.3	68.3
<b>Open-source Visual Grounded Reasoning Models</b>									
Pixel-Reasoner-7B	80.6	83.5	76.3	72.9	86.0	60.3	66.9	80.0	54.3
DeepEyes-7B	90.0	92.1	86.8	75.1	<b>91.3</b>	59.0	72.6	<b>86.8</b>	58.5
<b>TreeVGR-7B</b>	<b>91.1</b>	<b>94.0</b>	<b>87.0</b>	<b>77.1</b>	90.3	<b>64.0</b>	<b>73.1</b>	86.5	<b>59.8</b>
Δ v.s. Qwen2.5-VL-7B	↑ 16.8	↑ 16.6	↑ 17.3	↑ 5.0	↑ 1.5	↑ 8.5	↑ 4.3	↑ 3.0	↑ 5.8

454 Table 4: Comparison with state-of-the-art alternatives on MME-RealWorld-Lite (Zhang et al., 2024a).  
 455 All results are self-collected. The best performance is highlighted in **bold**.

	Perception					Reasoning				
	Overall	OCR	RS	DT	MO	AD	OCR	DT	MO	AD
<b>General Models</b>										
Qwen2.5-VL-7B	42.3	87.6	32.7	83.0	27.3	30.0	72.0	62.0	28.7	23.0
Qwen2.5-VL-72B	43.7	<b>90.8</b>	34.0	87.0	27.9	30.6	74.0	61.0	26.7	25.5
LLaVA-OneVision-7B	43.7	80.0	40.0	56.0	31.7	39.4	65.0	33.0	38.0	32.0
LLaVA-OneVision-72B	48.7	79.2	50.7	67.0	37.9	40.0	76.0	41.0	38.7	39.3
InternVL3-8B	47.9	83.6	49.3	75.0	34.5	36.9	70.0	44.0	40.0	37.0
InternVL3-78B	52.3	87.6	<b>54.7</b>	77.0	42.6	36.6	76.0	56.0	46.0	<b>40.3</b>
<b>Visual Grounded Reasoning Models</b>										
Pixel-Reasoner-7B	49.7	89.6	52.0	86.0	38.9	30.9	71.0	<b>72.0</b>	46.0	32.5
DeepEyes-7B	53.2	90.0	52.7	<b>89.0</b>	43.3	33.4	76.0	69.0	44.0	35.0
<b>TreeVGR-7B</b>	<b>54.9</b>	87.6	50.7	83.0	<b>47.0</b>	<b>43.4</b>	74.0	66.0	<b>51.3</b>	39.0
Δ v.s. Qwen2.5-VL-7B	↑ 12.6	– 0.0	↑ 18.0	– 0.0	↑ 19.7	↑ 13.4	↑ 2.0	↑ 4.0	↑ 22.6	↑ 16.0

472 reasoning chains via bounding box annotations), and vision-centric second-order reasoning. Constructed through expert-driven annotation and multi-stage quality control, **TreeBench** features 405  
 473 high-difficulty visual question-answer pairs with precise bounding boxes, emphasizing small objects  
 474 in real-world scenarios. It reveals the limitations of state-of-the-art models, *e.g.*, OpenAI-o3 (OpenAI,  
 475 2025) scores 54.8%, while setting a new standard for assessing nuanced visual grounding, multi-step  
 476 reasoning transparency, and cross-modal interaction.

477 **TreeVGR** advances VGR training through reinforcement learning guided by dual IoU rewards, which  
 478 explicitly supervise bounding box generation to ensure both precision and recall. This approach  
 479 enables explainable reasoning pathways and achieves significant improvements across benchmarks.

480 **Limitation and future works.** The current implementation of **TreeVGR** is based on a 7B parameter  
 481 model, which may limit scalability compared to larger architectures. **TreeBench** contains only 405  
 482 rigorously curated question-answer pairs. Expanding the benchmark with additional samples across  
 483 broader domains would further challenge model capabilities. Scaling up would be future work.

484

486 ETHICS STATEMENT  
487

488 Our research is grounded in ethical practices, with particular attention paid to the responsible use  
489 of data. All datasets employed in this study are publicly available and well-established within the  
490 computer vision community. Specifically, our benchmarking was conducted on SA-1B (Kirillov  
491 et al., 2023). Our use of this data is in accordance with their provided licenses and intended academic  
492 purpose.

493  
494 REPRODUCIBILITY STATEMENT  
495

496 We are committed to ensuring the reproducibility of the research presented in this paper. To this end,  
497 comprehensive implementation details for our models and experiments are provided in Appendix E,  
498 including the training procedures and all hyperparameters used. Furthermore, upon acceptance of  
499 this paper, all source code, datasets, and trained model checkpoints will be made publicly available.

500 REFERENCES  
501

502 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur  
503 Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot  
504 learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:23716–23736, 2022.

505 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang,  
506 Jun Tang, et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025a.

507 Sule Bai, Mingxing Li, Yong Liu, Jing Tang, Haoji Zhang, Lei Sun, Xiangxiang Chu, and Yansong Tang.  
508 Univg-r1: Reasoning guided universal visual grounding with reinforcement learning. *arXiv preprint  
509 arXiv:2505.14231*, 2025b.

511 Ali Furkan Biten, Ron Litman, Yusheng Xie, Srikanth Appalaraju, and R Manmatha. Latr: Layout-aware  
512 transformer for scene-text vqa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
513 Recognition (CVPR)*, 2022.

514 Meng Cao, Haoze Zhao, Can Zhang, Xiaojun Chang, Ian Reid, and Xiaodan Liang. Ground-r1: Incentivizing  
515 grounded visual reasoning via reinforcement learning. *arXiv preprint arXiv:2505.20272*, 2025.

516 Liang Chen, Lei Li, Haozhe Zhao, and Yifan Song. R1-v: Reinforcing super generalization ability in vision-  
517 language models with less than \$3, 2025.

519 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao,  
520 Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint  
521 arXiv:2403.20330*, 2024a.

522 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu,  
523 Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models  
524 with open-source suites. *Science China Information Sciences*, 67(12):220101, 2024b.

525 Google DeepMind. Gemini-2.5-flash. <https://deepmind.google/models/gemini/flash/>,  
526 2025a.

527 Google DeepMind. Gemini-2.5-pro. <https://deepmind.google/models/gemini/pro/>, 2025b.

529 Hongyuan Dong, Jiawen Li, Bohong Wu, Jiacong Wang, Yuan Zhang, and Haoyuan Guo. Benchmarking and  
530 improving detail image caption. *arXiv preprint arXiv:2405.19092*, 2024.

531 Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang,  
532 Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models.  
533 In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 11198–11201, 2024.

535 Yue Fan, Xuehai He, Diji Yang, Kaizhi Zheng, Ching-Chen Kuo, Yuting Zheng, Sravana Jyothi Narayananaraju,  
536 Xinze Guan, and Xin Eric Wang. Grit: Teaching mllms to think with images. *arXiv preprint arXiv:2505.15879*,  
537 2025.

538 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,  
539 Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.  
540 *arXiv preprint arXiv:2501.12948*, 2025a.

540 Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang,  
 541 Jiawei Wang, et al. Seed1.5-vl technical report. *arXiv preprint arXiv:2505.07062*, 2025b.

542

543 Tuomo Hiippala, Malihe Alikhani, Jonas Haverinen, Timo Kalliokoski, Evanfiya Logacheva, Serafina Orekhova,  
 544 Aino Tuomainen, Matthew Stone, and John A Bateman. Ai2d-rst: A multimodal corpus of 1000 primary  
 545 school science diagrams. *Language Resources and Evaluation*, 55:661–688, 2021.

546

547 Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui  
 548 Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint  
 arXiv:2503.06749*, 2025.

549

550 Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhi Wang, Jianhan Jin, Claire  
 551 Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning  
 552 quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025.

553

554 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer  
 555 Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF  
 International Conference on Computer Vision (ICCV)*, pp. 4015–4026, 2023.

556

557 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gon-  
 558 zalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with  
 559 pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

560

561 Weixian Lei, Jiacong Wang, Haochen Wang, Xiangtai Li, Jun Hao Liew, Jiashi Feng, and Zilong Huang. The  
 562 scalability of simplicity: Empirical analysis of vision-language learning with a single transformer. *arXiv  
 preprint arXiv:2504.10462*, 2025.

563

564 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei  
 565 Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

566

567 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking  
 568 multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023a.

569

570 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training  
 571 with frozen image encoders and large language models. In *International Conference on Machine Learning  
 (ICML)*, pp. 19730–19742, 2023b.

571

572 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucina-  
 573 tion in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023c.

574

575 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and  
 576 C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer  
 Vision (ECCV)*, pp. 740–755, 2014.

576

577 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in Neural  
 578 Information Processing Systems (NeurIPS)*, 36:34892–34916, 2023a.

578

579 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-  
 580 next: Improved reasoning, ocr, and world knowledge. <https://llava-vl.github.io/blog/2024-01-30-llava-next/>, 2024a.

580

581 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang,  
 582 Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint  
 arXiv:2307.06281*, 2023b.

583

584 Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On  
 585 the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023c.

586

587 Yuqi Liu, Tianyuan Qu, Zhisheng Zhong, Bohao Peng, Shu Liu, Bei Yu, and Jiaya Jia. Visionreasoner: Unified  
 588 visual perception and reasoning via reinforcement learning. *arXiv preprint arXiv:2505.12081*, 2025a.

589

590 Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang.  
 591 Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025b.

591

592 Zuyan Liu, Yuhao Dong, Yongming Rao, Jie Zhou, and Jiwen Lu. Chain-of-spot: Interactive reasoning improves  
 593 large vision-language models. *arXiv preprint arXiv:2403.12966*, 2024b.

593

Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint  
 arXiv:1608.03983*, 2016.

594 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*,  
 595 2017.

596

597 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang,  
 598 Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in  
 599 visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.

600 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for  
 601 question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.

602

603 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In  
 604 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021.

605 Debjyoti Mondal, Suraj Modi, Subhadarshi Panda, Rituraj Singh, and Godawari Sudhakar Rao. Kam-cot:  
 606 Knowledge augmented multimodal chain-of-thoughts reasoning. In *Proceedings of the AAAI Conference on*  
 607 *Artificial Intelligence (AAAI)*, volume 38, pp. 18798–18806, 2024.

608 OpenAI. Openai-gpt-4o. <https://openai.com/index/gpt-4o-system-card/>, 2024a.

609

610 OpenAI. Openai-o1. <https://openai.com/o1/>, 2024b.

611

612 OpenAI. Openai-o3. <https://openai.com/index/introducing-o3-and-o4-mini/>, 2025.

613 Ji Qi, Ming Ding, Weihan Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao  
 614 Dong, et al. Cogcom: Train large vision-language models diving into details through chain of manipulations.  
 615 *arXiv preprint arXiv:2402.04236*, 2024.

616 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
 617 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language  
 618 supervision. In *International Conference on Machine Learning (ICML)*, pp. 8748–8763, 2021.

619

620 Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng  
 621 Li. Visual-cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for  
 622 chain-of-thought reasoning. *Advances in Neural Information Processing Systems*, 37:8612–8642, 2024a.

623 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang,  
 624 YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models.  
 625 *arXiv preprint arXiv:2402.03300*, 2024b.

626 Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia  
 627 Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv*  
 628 *preprint arXiv:2504.07615*, 2025.

629

630 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin,  
 631 and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.

632 Alex Su, Haozhe Wang, Weimin Ren, Fangzhen Lin, and Wenhui Chen. Pixel reasoner: Incentivizing pixel-space  
 633 reasoning with curiosity-driven reinforcement learning. *arXiv preprint arXiv:2505.15966*, 2025.

634

635 Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press  
 636 Cambridge, 1998.

637 Qwen Team. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

638

639 Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri IYER, Sai Charitha Akula,  
 640 Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. Cambrian-1: A fully open, vision-  
 641 centric exploration of multimodal llms. *Advances in Neural Information Processing Systems (NeurIPS)*, 37:  
 642 87310–87356, 2024a.

643 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring  
 644 the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
 645 and Pattern Recognition (CVPR)*, pp. 9568–9578, 2024b.

646 Fengxiang Wang, Mingshuo Chen, Yueying Li, Di Wang, Haotian Wang, Zonghao Guo, Zefan Wang, Boqi Shan,  
 647 Long Lan, Yulin Wang, et al. Geollava-8k: Scaling remote-sensing multimodal large language models to 8k  
 648 resolution. *arXiv preprint arXiv:2505.21375*, 2025a.

648 Fengxiang Wang, Hongzhen Wang, Zonghao Guo, Di Wang, Yulin Wang, Mingshuo Chen, Qiang Ma, Long  
 649 Lan, Wenjing Yang, Jing Zhang, et al. Xlrs-bench: Could your multimodal llms understand extremely large  
 650 ultra-high-resolution remote sensing imagery? In *Proceedings of the IEEE/CVF Conference on Computer  
 651 Vision and Pattern Recognition (CVPR)*, pp. 14325–14336, 2025b.

652 Haochen Wang, Yucheng Zhao, Tiancai Wang, Haoqiang Fan, Xiangyu Zhang, and Zhaoxiang Zhang. Ross3d:  
 653 Reconstructive visual instruction tuning with 3d-awareness. *Proceedings of the IEEE/CVF International  
 654 Conference on Computer Vision (ICCV)*, 2025c.

655 Haochen Wang, Anlin Zheng, Yucheng Zhao, Tiancai Wang, Zheng Ge, Xiangyu Zhang, and Zhaoxiang Zhang.  
 656 Reconstructive visual instruction tuning. In *International Conference on Learning Representations (ICLR)*,  
 657 2025d.

658 Jiacong Wang, Bohong Wu, Haiyong Jiang, Zhou Xun, Xin Xiao, Haoyuan Guo, and Jun Xiao. World to code:  
 659 Multi-modal data generation via self-instructed compositional captioning and filtering. In *Proceedings of the  
 660 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 4608–4623, 2024a.

661 Jiacong Wang, Zijiang Kang, Haochen Wang, Haiyong Jiang, Jiawen Li, Bohong Wu, Ya Wang, Jiao Ran, Xiao  
 662 Liang, Chao Feng, et al. Vgr: Visual grounded reasoning. *arXiv preprint arXiv:2506.11991*, 2025e.

663 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang,  
 664 Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution.  
 665 *arXiv preprint arXiv:2409.12191*, 2024b.

666 Weiyun Wang, Zhe Chen, Wenhui Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Jinguo Zhu, Xizhou Zhu,  
 667 Lewei Lu, Yu Qiao, et al. Enhancing the reasoning ability of multimodal large language models via mixed  
 668 preference optimization. *arXiv preprint arXiv:2411.10442*, 2024c.

669 Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, Wei Yu, and Dacheng Tao. Divide,  
 670 conquer and combine: A training-free framework for high-resolution image perception in multimodal large  
 671 language models. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 39, pp.  
 672 7907–7915, 2025f.

673 Lai Wei, Yuting Li, Chen Wang, Yue Wang, Linghe Kong, Weiran Huang, and Lichao Sun. Unsupervised  
 674 post-training for multi-modal llm reasoning via grpo. *arXiv preprint arXiv:2505.22453*, 2025a.

675 Lai Wei, Yuting Li, Kaipeng Zheng, Chen Wang, Yue Wang, Linghe Kong, Lichao Sun, and Weiran Huang.  
 676 Advancing multimodal reasoning via reinforcement learning with cold start. *arXiv preprint arXiv:2505.22334*,  
 677 2025b.

678 Penghao Wu and Saining Xie. V\*: Guided visual search as a core mechanism in multimodal llms. In *Proceedings  
 679 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13084–13094, 2024.

680 Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma,  
 681 Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for  
 682 advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024.

683 Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin,  
 684 Shuo Liu, et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language  
 685 models towards multitask agi. *arXiv preprint arXiv:2404.16006*, 2024.

686 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu,  
 687 Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint  
 688 arXiv:2503.14476*, 2025.

689 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang,  
 690 Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning  
 691 benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
 692 Recognition (CVPR)*, 2024.

693 Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun Wang,  
 694 Qingsong Wen, Zhang Zhang, et al. Mme-realworld: Could your multimodal llm challenge high-resolution  
 695 real-world scenarios that are difficult for humans? *arXiv preprint arXiv:2408.13257*, 2024a.

696 Yuan Zhang, Tao Huang, Chun-Kai Fan, Hongyuan Dong, Jiawen Li, Jiacong Wang, Kuan Cheng, Shanghang  
 697 Zhang, Haoyuan Guo, et al. Unveiling the tapestry of consistency in large vision-language models. *Advances  
 698 in Neural Information Processing Systems*, 37:118632–118653, 2024b.

702 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyen Luo, Zhangchi Feng, and Yongqiang Ma.  
703 Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual*  
704 *Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok,  
705 Thailand, 2024. Association for Computational Linguistics.

706 Yaowei Zheng, Junting Lu, Shenzhi Wang, Zhangchi Feng, Dongdong Kuang, and Yuwen Xiong. Easyrl: An  
707 efficient, scalable, multi-modality rl training framework. <https://github.com/hiyouga/EasyR1>,  
708 2025a.

709 Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing Yu.  
710 Deepeyes: Incentivizing "thinking with images" via reinforcement learning. *arXiv preprint arXiv:2505.14362*,  
711 2025b.

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

715 Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. Detection and  
716 tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*,  
717 44(11):7380–7399, 2021.

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

## APPENDIX

## A OVERVIEW

Here, we provide a table of contents:

- First, in Appendix B, we provide the annotation pipeline in detail, which includes three rounds of quality control.
- In Appendix C, we introduce statistics of our **TreeBench**.
- In Appendix D, we perform in-depth analysis on our **TreeBench**.
- In Appendix E, we provide implementation details of our two-stage training pipeline, including cold-start initialization and reinforcement learning with traceable evidence.
- In Appendix F, we provide more experiments of our **TreeVGR**, including results on general multimodal benchmarks and ablation studies.
- In Appendix G, we discuss our limitations in detail.
- Finally, in Appendix H, we provide qualitative examples and failure cases of our **TreeVGR**.

## B ANNOTATION PIPELINE

**TreeBench** was constructed through a systematic pipeline combining automated sampling, LMM-assisted generation, and three rounds of human validation. The annotation team contains eight human experts in LMMs, including six Ph.D candidates and two senior research scientists.

**1. Image Selection.** A total of 1K images are initially sampled from the SA-1B (Kirillov et al., 2023), with deliberate prioritization of images containing high-density objects (*e.g.*, scenes with overlapping or clustered items), as it offers high-resolution, real-world scenes with a large number of small and varied objects, making it particularly suitable for evaluating visual grounded reasoning. To ensure balanced representation across categories, 100 images are initially allocated per category.

**2. First Round Quality Control.** The annotation team manually evaluates the relevance and quality of each image for its assigned category. This step is critical for addressing category-specific requirements, *e.g.*, the “Ordering” category necessitates images with visually similar or repetitive objects for practical reasoning tasks. Following this review, 647 images meet the criteria.

**3. Automated Question Generation.** Question-option-answer trios are then generated using two advanced LMMs, *i.e.*, OpenAI-o3 (OpenAI, 2025) and Gemini-2.5-Pro (DeepMind, 2025b), each tasked with producing three diverse, high-quality questions per image. Prompts are designed to emphasize task-specific complexity and visual-semantic alignment.

**4. Second Round Quality Control.** Human experts then manually review all six model-generated questions per image. For each image, annotators selected the most semantically coherent and task-relevant question from the pool of six, prioritizing: (1) alignment with the target subtask, (2) avoidance of trivial or ambiguous object referring, and (3) clarity and unambiguous answerability. If none of the six questions met these criteria, annotators manually constructed a new question. This step ensures that only high-quality, human-vetted questions advance to the next stage.

**5. Difficulty Filtering.** Questions deemed insufficiently challenging are removed through model-based consensus screening. Specifically, any question answered correctly by all four state-of-the-art vision-language models (Qwen2.5-VL-72B (Bai et al., 2025a), InternVL3-78B (Zhu et al., 2025), GPT-4o (OpenAI, 2024a), Gemini-2.5-Flash (DeepMind, 2025a)) was excluded to ensure the benchmark retained meaningful difficulty.

**6. Third Round Quality Control.** The final cross-verification phase engages independent human annotators to cross-validate the accuracy and relevance of each question-option-answer pair. The final dataset comprised 405 rigorously validated questions.

## C BENCHMARK STATISTICS

**Distribution of Each Subtask.** As demonstrated in Figure 4, **TreeBench** emphasizes advanced reasoning tasks, accounting for 63% of the total subtasks (256 questions), while basic perception-

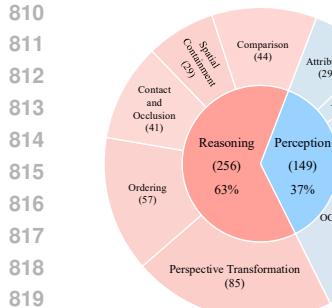


Figure 4: Distribution of each discipline in **TreeBench**, which prioritizes reasoning over perception.

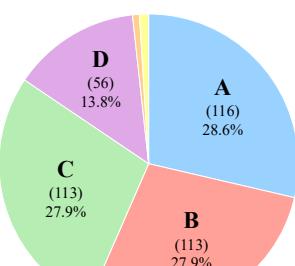


Figure 5: The ground-truth distribution of **TreeBench** with 3 instances of E and 4 instances of F.

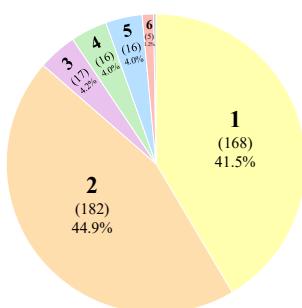


Figure 6: Distribution of the number of instances in **TreeBench**, with one question with 8 target instances.

related tasks constitute 37% (149 questions). Within the reasoning category, key subtasks reflect a focus on complex spatial and relational understanding. This structure underscores a deliberate prioritization of higher-order reasoning over foundational perceptual tasks, aligning with the goal of challenging models to process nuanced relationships and transformations rather than mere object recognition or attribute detection.

**Distribution of Answers.** As illustrated in Figure 5, the ground-truth distribution of **TreeBench** is dominated by four main categories: A (28.6%, 116 instances), B (27.9%, 113 instances), C (27.9%, 113 instances), and D (13.8%, 56 instances). These account for 98.2% of the total 405 instances. The remaining 1.8% (7 instances) includes E (3 instances) and F (4 instances). This structure highlights a balanced emphasis on categories A, B, and C, with D as a notable secondary group, while E and F represent minor but distinct components.

**Distribution of the Number of Target Instances.** Figure 6 shows the distribution of the number of target instances per question. The majority of questions in **TreeBench** require identifying 1 or 2 target instances, accounting for 41.5% (168 questions) and 44.9% (182 questions) of the total, respectively. Questions requiring 3, 4, 5, or 6 targets constitute smaller fractions: 4.2% (17 questions), 4.0% (16 questions), 4.0% (16 questions), and 1.2% (5 questions), respectively. Notably, a single question (highlighted in gray) demands 8 target instances, representing an extreme case. Overall, 86.4% of questions focus on 1–2 targets, suggesting a balance between simplicity and complexity in task design while incorporating rare multi-target scenarios for comprehensive evaluation.

**Distribution of Target Instance Area.** We compute the *relative* area for each target instance using its bounding box, *i.e.*,  $\text{area} = \frac{1}{HW}(y_2 - y_1)(x_2 - x_1)$ , where  $H$  and  $W$  are the input resolution. Figure 7 is the histogram of the mean area for each question. It illustrates that the majority of target instances in **TreeBench** are extremely small, with a sharp peak near 0.0 and a long tail extending to larger areas (up to 0.7). The mean area across all questions is 0.0305, confirming that targets are predominantly tiny. Most questions (highest frequency bin) involve target instances with areas clustered around 0.0 to 0.05, while only a small fraction require identifying larger objects. This distribution highlights the importance of addressing challenging scenarios where small-scale object detection and reasoning are crucial, potentially compromising model performance.

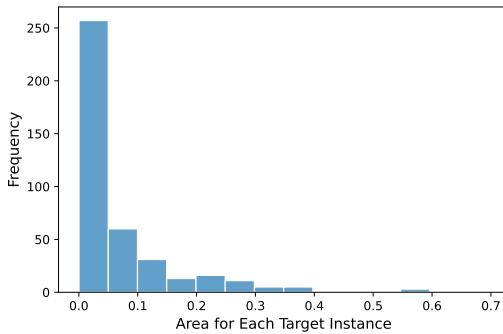


Figure 7: The histogram of mean target instance areas per question with a low average of 0.0305 (indicating small target instances).

## D ANALYSIS OF TREEBENCH

**Correlation between Localization and Performance.** Importantly, for visual grounded reasoning models, our traceable evaluation demonstrates a *positive correlation* between localization preci-

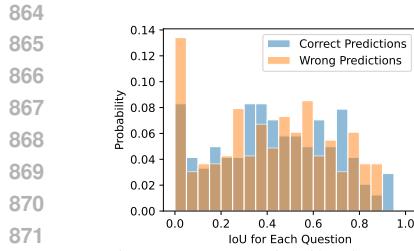


Figure 8: Distribution of IoU for each question in **TreeBench**.

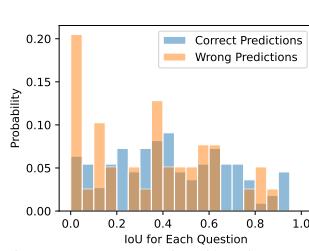


Figure 9: Distribution of IoU for each question in **TreeBench-Perception**.

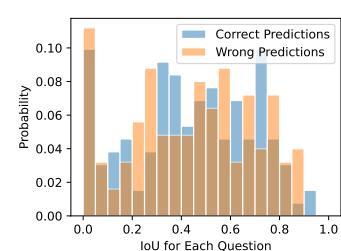


Figure 10: Distribution of IoU for each question in **TreeBench-Reasoning**.

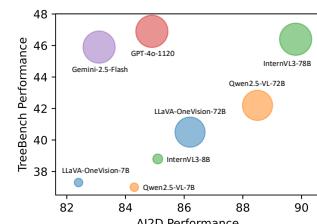


Figure 11: Performance decoupling with AI2D (Hiippala et al., 2021).

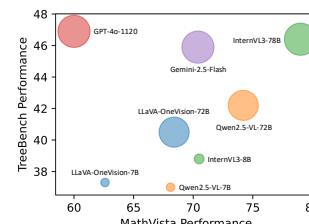


Figure 12: Performance decoupling with MathVista (Lu et al., 2023).

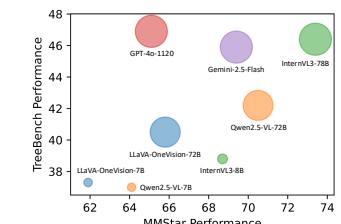


Figure 13: Performance decoupling with MMStar (Chen et al., 2024a).

889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
18010  
18011  
18012  
18013  
18014  
18015  
18016  
18017  
18018  
18019  
18020  
18021  
18022  
18023  
18024  
18025  
18026  
18027  
18028  
18029  
18030  
18031  
18032  
18033  
18034  
18035  
18036  
18037  
18038  
18039  
18040  
18041  
18042  
18043  
18044  
18045  
18046  
18047  
18048  
18049  
18050  
18051  
18052  
18053  
18054  
18055  
18056  
18057  
18058  
18059  
18060  
18061  
18062  
18063  
18064  
18065  
18066  
18067  
18068  
18069  
18070  
18071  
18072  
18073  
18074  
18075  
18076  
18077  
18078  
18079  
18080  
18081  
18082  
18083  
18084  
18085  
18086  
18087  
18088  
18089  
18090  
18091  
18092  
18093  
18094  
18095  
18096  
18097  
18098  
18099  
180100  
180101  
180102  
180103  
180104  
180105  
180106  
180107  
180108  
180109  
180110  
180111  
180112  
180113  
180114  
180115  
180116  
180117  
180118  
180119  
180120  
180121  
180122  
180123  
180124  
180125  
180126  
180127  
180128  
180129  
180130  
180131  
180132  
180133  
180134  
180135  
180136  
180137  
180138  
180139  
180140  
180141  
180142  
180143  
180144  
180145  
180146  
180147  
180148  
180149  
180150  
180151  
180152  
180153  
180154  
180155  
180156  
180157  
180158  
180159  
180160  
180161  
180162  
180163  
180164  
180165  
180166  
180167  
180168  
180169  
180170  
180171  
180172  
180173  
180174  
180175  
180176  
180177  
180178  
180179  
180180  
180181  
180182  
180183  
180184  
180185  
180186  
180187  
180188  
180189  
180190  
180191  
180192  
180193  
180194  
180195  
180196  
180197  
180198  
180199  
180200  
180201  
180202  
180203  
180204  
180205  
180206  
180207  
180208  
180209  
180210  
180211  
180212  
180213  
180214  
180215  
180216  
180217  
180218  
180219  
180220  
180221  
180222  
180223  
180224  
180225  
180226  
180227  
180228  
180229  
180230  
180231  
180232  
180233  
180234  
180235  
180236  
180237  
180238  
180239  
180240  
180241  
180242  
180243  
180244  
180245  
180246  
180247  
180248  
180249  
180250  
180251  
180252  
180253  
180254  
180255  
180256  
180257  
180258  
180259  
180260  
180261  
180262  
180263  
180264  
180265  
180266  
180267  
180268  
180269  
180270  
180271  
180272  
180273  
180274  
180275  
180276  
180277  
180278  
180279  
180280  
180281  
180282  
180283  
180284  
180285  
180286  
180287  
180288  
180289  
180290  
180291  
180292  
180293  
180294  
180295  
180296  
180297  
180298  
180299  
180300  
180301  
180302  
180303  
180304  
180305  
180306  
180307  
180308  
180309  
180310  
180311  
180312  
180313  
180314  
180315  
180316  
180317  
180318  
180319  
180320  
180321  
180322  
180323  
180324  
180325  
180326  
180327  
180328  
180329  
180330  
180331  
180332  
180333  
180334  
180335  
180336  
180337  
180338  
180339  
180340  
180341  
180342  
180343  
180344  
180345  
180346  
180347  
180348  
180349  
180350  
180351  
180352  
180353  
180354  
180355  
180356  
180357  
180358  
180359  
180360  
180361  
180362  
180363  
180364  
180365  
180366  
180367  
180368  
180369  
180370  
180371  
180372  
180373  
180374  
180375  
180376  
180377  
180378  
180379  
180380  
180381  
180382  
180383  
180384  
180385  
180386  
180387  
180388  
180389  
180390  
180391  
180392  
180393  
180394  
180395  
180396  
180397  
180398  
180399  
180400  
180401  
180402  
180403  
180404  
180405  
180406  
180407  
180408  
180409  
180410  
180411  
180412  
180413  
180414  
180415  
180416  
180417  
180418  
180419  
180420  
180421  
180422  
180423  
180424  
180425  
180426  
180427  
180428  
180429  
180430  
180431  
180432  
180433  
180434  
180435  
180436  
180437  
180438  
180439  
180440  
180441  
180442  
180443  
180444  
180445  
180446  
180447  
180448  
180449  
180450  
180451  
180452  
180453  
180454  
180455  
180456  
180457  
180458  
180459  
180460  
180461  
180462  
180463  
180464  
180465  
180466  
180467  
180468  
180469  
180470  
180471  
180472  
180473  
180474  
180475  
180476  
180477  
180478  
180479  
180480  
180481  
180482  
180483  
180484  
180485  
180486  
180487  
180488  
180489  
180490  
180491  
180492  
180493  
180494  
180495  
180496  
180497  
180498  
180499  
180500  
180501  
180502  
180503  
180504  
180505  
180506  
180507  
180508  
180509  
180510  
180511  
180512  
180513  
180514  
180515  
180516  
180517  
180518  
180519  
180520  
180521  
180522  
180523  
180524  
180525  
180526  
180527  
180528  
180529  
180530  
180531  
180532  
180533  
180534  
180535  
180536  
180537  
180538  
180539  
180540  
180541  
180542  
180543  
180544  
180545  
180546  
180547  
180548  
180549  
180550  
180551  
180552  
180553  
180554  
180555  
180556  
180557  
180558  
180559  
180560  
180561  
180562  
180563  
180564  
180565  
180566  
180567  
180568  
180569  
180570  
180571  
180572  
180573  
180574  
180575  
180576  
180577  
180578  
180579  
180580  
180581  
180582  
180583  
180584  
180585  
180586  
180587  
180588  
180589  
180590  
180591  
180592  
180593  
180594  
180595  
180596  
180597  
180598  
180599  
180600  
180601  
180602  
180603  
1

918  
 919 **Table 5: Performance comparison with *masked* target instances.** When masking out all target  
 920 instances on **TreeBench**, we observe a significant performance drop across all models, confirming  
 921 that the annotated bounding boxes are not only high-quality but also indispensable for accurate visual  
 922 grounded reasoning.

Masking	Qwen2.5-VL-7B	InternVL3-8B	GPT-4o	o3	Gemini-2.5-Flash	Gemini-2.5-Pro
	37.0	38.8	46.9	54.8	45.9	54.1
✓	31.8 $\downarrow$ 5.2	29.6 $\downarrow$ 9.2	29.1 $\downarrow$ 17.8	33.8 $\downarrow$ 21.0	29.9 $\downarrow$ 16.0	33.1 $\downarrow$ 21.0

923  
 924 **Table 6: Performance comparison with explicit bounding boxes-based textual hints.** When  
 925 we provide ground-truth bounding boxes as explicit evidence hints to models, all models achieve  
 926 consistent performance gains.

Textual Boxes	Qwen2.5-VL-7B	InternVL3-8B	GPT-4o	o3	Gemini-2.5-Flash	Gemini-2.5-Pro
	37.0	38.8	46.9	54.8	45.9	54.1
✓	43.7 $\uparrow$ 6.7	43.5 $\uparrow$ 4.7	49.4 $\uparrow$ 2.5	58.3 $\uparrow$ 3.5	51.9 $\uparrow$ 6.0	61.0 $\uparrow$ 6.9

## 934 E IMPLEMENTATION DETAILS

### 935 E.1 COLD-START INITIALIZATION

936 **Data Construction.** We base our supervised fine-tuning (SFT) dataset on VGR-158K (Wang et al.,  
 937 2025e), which provides pseudo-chain-of-thought annotations paired with bounding boxes for visual  
 938 reasoning tasks. However, to align with the grounding capabilities of our base model (Qwen2.5-VL  
 939 series (Bai et al., 2025a)), which outputs *absolute* coordinates rather than the normalized coordinates  
 940 (ranging from 0 to 1) used by LLaVA-NeXT (Liu et al., 2024a) in (Wang et al., 2025e), we perform  
 941 coordinate system conversion. Specifically, for each bounding box, we transform normalized coordinates  
 942  $[r_{x_1}, r_{y_1}, r_{x_2}, r_{y_2}]$  into *absolute* coordinates via  $[x_1, y_1, x_2, y_2] = [Wr_{x_1}, Hr_{y_1}, Wr_{x_2}, Hr_{y_2}]$ ,  
 943 where  $H \times W$  is the resolution of the input image. Next, we filter samples to prioritize complex  
 944 reasoning pathways, retaining only entries with multiple bounding boxes (*i.e.*, more than one box  
 945 per reasoning trajectory). This yields 35K samples, as multi-box interactions demand stronger  
 946 spatial-temporal reasoning compared to single-box tasks. Subsequently, we construct a reflective  
 947 subset of 4.7K samples among them by introducing controlled perturbations: for each sample, we (1)  
 948 inject a synthetic error by inserting a randomly generated incorrect bounding box into the reasoning  
 949 sequence, and (2) append the meta-cognitive prompt “Wait, this box seems to be wrong” immediately  
 950 afterward, resulting in our **TreeVGR-SFT-35K**. This design explicitly trains the model to detect and  
 951 correct erroneous visual grounding, which is a critical skill for robust real-world deployment.

952 **Optimization.** Initialized from Qwen2.5-VL-7B-Instruct (Bai et al., 2025a), we train **TreeVGR-7B-CI** (“CI” here stands for Cold Initialization) with 8 GPUs using LLaMA-Factory (Zheng et al., 2024),  
 953 where the AdamW optimizer (Loshchilov & Hutter, 2017) with a learning rate of 5e-6 and a global  
 954 batch size of 256 is utilized. The learning rate is decayed following a cosine schedule (Loshchilov &  
 955 Hutter, 2016) with a warmup ratio of 0.1.

### 956 E.2 REINFORCEMENT LEARNING

957 **Data Construction.** **TreeVGR** incorporates a novel dual IoU reward, which means each sample  
 958 should contain ground-truth bounding boxes during the RL phase. To this end, we filter *hard* samples  
 959 from the original 191K training set of V\* (Wu & Xie, 2024) using Qwen2.5-VL-7B-Instruct (Bai  
 960 et al., 2025a), resulting in 30K samples. Additionally, we incorporate the VisDrone dataset (Zhu  
 961 et al., 2021), which is originally designed for detection and tracking under UAV images, which offers  
 962 extremely high-resolution, real-world scenes with a large number of small and varied objects and  
 963 their corresponding bounding box annotations. We reformulate the training set and the validation set  
 964 into 38K multiple-choice counting problems, and only retain samples with the ground-truth number  
 965 ranging from 5 to 10, contributing to the final 7K samples. Finally, our **TreeVGR-RL-37K** consists  
 966 of 30K open-ended question-answering samples from V\* (Wu & Xie, 2024) and 7K multiple-choice  
 967 problems from VisDrone (Zhu et al., 2021).

968 **Optimization.** Initialized from **TreeVGR-7B-CI**, we train our final **TreeVGR-7B** with 8 GPUs,  
 969 with another 8 GPUs serving the reward model, *i.e.*, Qwen2.5-72B-Instruct (Team, 2024), using  
 970 vLLM (Kwon et al., 2023). We adopt Group Relative Policy Optimization (GRPO) (Shao et al.,

972  
 973 Table 7: Comparison with state-of-the-art alternatives on other multimodal benchmarks, including  
 974 CV-Bench (Tong et al., 2024a), MMVP (Tong et al., 2024b), MMBench (Liu et al., 2023b), POPE (Li  
 975 et al., 2023c), AI2D (Hiippala et al., 2021), and ChartQA (Masry et al., 2022).  $\dagger$ Results are obtained  
 from (Guo et al., 2025b), otherwise are self-collected.

976 977 978 979 980 981 982 983 984	Capability	Benchmark	Qwen2.5-VL-7B	TreeVGR-7B	Qwen2.5-VL-72B
question answering	Vision-centric	CV-Bench-2D	74.1	<b>76.9</b> $\uparrow$ 2.8	77.7
	CV-Bench-3D	72.6	<b>77.6</b> $\uparrow$ 5.0	87.0	
	MMVP	66.7	<b>75.3</b> $\uparrow$ 8.6	66.7 $\dagger$	
General VQA	MMBench <sup>en</sup> <sub>dev</sub>	83.1	<b>84.4</b> $\uparrow$ 1.3	88.6 $\dagger$	
	POPE	86.7	<b>87.2</b> $\uparrow$ 0.5	84.9	
Document and chart	AI2D <sub>test</sub>	<b>84.9</b>	84.8 $\downarrow$ 0.1	88.7 $\dagger$	
	ChartQA <sub>test</sub>	85.6	<b>85.8</b> $\uparrow$ 0.2	89.5 $\dagger$	

985  
 986 Table 8: Ablations of each component of our **TreeVGR**. “MME-RW” stands for MME-RealWorld-  
 987 Lite (Zhang et al., 2024a), and “Acc” represents the multiple-choice accuracy.  $\dagger$ This improvement  
 988 mainly comes from the training set, as many training samples from V\* (Wu & Xie, 2024) are included  
 989 in RL.  $\ddagger$ The model *enumerates* boxes to obtain larger IoU recall, and fails to produce final answers.

990 991 992		Rewards			TreeBench		V*	MME-RW	
		Cold-Start	$R_{\text{acc}} + R_{\text{format}}$	$R_{\text{IoU}}^{\text{R}}$	$R_{\text{IoU}}^{\text{P}}$	Acc	mIoU	Acc	
①	Qwen2.5-VL-7B					37.0	–	71.2	42.3
②	Cold-Start	✓				39.0	23.4	76.4	48.4
③	<b>TreeVGR</b>	✓	✓	✓	✓	<b>50.4</b>	<b>44.0</b>	<b>91.1</b>	<b>54.9</b>
④	w/o Traceable Evidence	✓	✓			38.0	27.2	87.9 $\dagger$	51.6
⑤	w/o Precision $\ddagger$	✓	✓		✓	0.0	78.3	0.0	0.0
⑥	w/o Recall	✓	✓		✓	45.4	20.6	89.5	52.6
⑦	Text-Only RL		✓			39.0	–	86.9 $\dagger$	46.3

1000 2024b), which has been proved to be effective and efficient for diverse tasks. We have also tried  
 1001 DAPO (Yu et al., 2025), but we find it unstable compared with GRPO. Therefore, we simply utilize  
 1002 the original GRPO (Shao et al., 2024b). We implement using EasyR1 (Zheng et al., 2025a), which is  
 1003 a clean fork of veRL (Sheng et al., 2024). We train our **TreeVGR-7B** with 5 epochs on **TreeVGR-  
 1004 RL-37K**, which is significantly less than DeepEyes-7B (Zheng et al., 2025b) (which is trained on  
 1005 47K samples with 32 epochs).

## 1007 F MORE EXPERIMENTS

### 1008 F.1 RESULTS ON OTHER MULTIMODAL BENCHMARKS

1010 In Table 7, we compare our **TreeVGR** with its base model Qwen2.5-VL-7B (Bai et al., 2025a) on  
 1011 a variety of conventional multimodal benchmarks. Specifically, we select CV-Bench (Tong et al.,  
 1012 2024a) and MMVP (Tong et al., 2024b) to evaluate vision-centric question-answering capabilities.  
 1013 MMBench (Liu et al., 2023b) and POPE (Li et al., 2023c) are selected for evaluating general VQA  
 1014 capabilities, and AI2D (Hiippala et al., 2021) and ChartQA (Masry et al., 2022) for comprehension  
 1015 with document and chart. We observe significant improvements in most cases, especially for vision-  
 1016 centric benchmarks. Notably, **TreeVGR-7B** achieves 75.3 on MMVP (Tong et al., 2024b), even  
 1017 surpasses Qwen2.5-VL-72B (Bai et al., 2025a) by a significant margin.

### 1018 F.2 ABLATION STUDIES

1019 The core contribution of **TreeVGR** is the *traceable* training pipeline, where the dual IoU reward  
 1020  $R_{\text{IoU}}$  is incorporated in conventional RL training. Therefore, we aim to evaluate the effectiveness  
 1021 of including this traceable term. As demonstrated in Table 8, we ablate each component of our  
 1022 **TreeVGR**, including the cost-start initialization and reward functions.

1024 The cold-start stage is quite beneficial for visual grounded reasoning, when compared with ①  
 1025 and ②. This means the formatting of outputting bounding boxes of target instances is useful for  
 conventional visual grounded reasoning benchmarks like V\* Bench (Wu & Xie, 2024) and MME-

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
RealWorld-Lite (Zhang et al., 2024a). Note that these benchmarks can be regarded as Out-of-Domain (OOD) samples for the SFT dataset.

**Traceable visual grounded reasoning is more effective than untraceable one**, when compared with ③ and ④. Starting from the *same* cold-start checkpoint, integrating dual IoU rewards into the RL framework yields substantial performance gains, particularly on our **TreeBench** and MME-RealWorld-Lite (Zhang et al., 2024a), which represent out-of-distribution (OOD) scenarios relative to the RL training data. Notably, on **TreeBench**, our **TreeVGR** demonstrates significant enhancements in both overall accuracy and mIoU. This dual improvement suggests that precise and interpretable reasoning pathways are critical for achieving optimal performance, underscoring the value of structured reward design in complex, real-world tasks.

**The precision term is crucial for alleviating the repetition problem**, when compared with ③ and ⑤. As illustrated in Figure 14, without precision, the mean response length grows rapidly. When evaluating this model, we find that it tends to *enumerate* candidate bounding boxes to obtain larger IoU recall and thus always fails to produce final answers.

**The recall term is crucial for precise and complete localization**, when compared with ③ and ⑥. On **TreeBench**, without the recall term, the model achieves significant accuracy improvements, but the localization accuracy (mIoU) remains limited, usually grounding *incomplete* target instances.

**Vanilla text-only RL is not so effective as visual grounded reasoning**, when compared with ③ and ⑦. Vanilla RL in text-based tasks demonstrates value through its text-space reasoning capabilities. However, when integrating visual grounded reasoning with traceable evidence, the performance gains become more significant. This highlights the critical role of two factors: (1) pre-answer contextual grounding to anchor responses in multimodal evidence, and (2) accurate spatial localization to refine decision-making precision.

## G LIMITATIONS AND FUTURE WORKS

One possible limitation of **TreeVGR** is the model scale and architecture, which is limited to Qwen2.5-VL-7B (Bai et al., 2025a). Experiments with other base models and larger model scales could be future work. Furthermore, **TreeVGR** is *not* a general multimodal reasoner, as it is not designed to perform ultra-long reasoning processes in math, sciences, and coding. How to effectively unify vision-centric reasoning models with standard text-centric models could be a future work.

As for **TreeBench**, we find that the “perspective transform” protocol becomes one of the major bottlenecks, which means current state-of-the-art multimodal models, even including visual grounded reasoning models, have *not* effectively modeled the ego-view 3D awareness. Moreover, basic perception capabilities *under complex scenes* are also limited, leading to relatively low scores on “attributes” and “material”. How to effectively let LMMs perceive *any* details of the dense visual world becomes a critical challenge.

## H QUALITATIVE EXAMPLES OF TREEVGR

We first provide qualitative examples of the correlation between the predicted bounding boxes and the model’s internal attention maps in Figure 15, where the attention values are the *mean* attention of generated tokens related to bounding boxes (highlighted in red). Strong correlations are observed across examples, indicating that the model’s focus is truly guided by those textual bounding boxes-based hints. This explains why visual grounded reasoning pipelines could be such effective. Note that sometimes the visual attention is still *being overwhelmed by textual domain knowledge* even if the grounding is accurate (see the failure case in Figure 15). Explicitly enforcing the model’s focus beyond attention may be one potential solution.

Next, we provide qualitative examples and failure cases of our **TreeVGR** on each category of **TreeBench** in Figures 16 to 25.

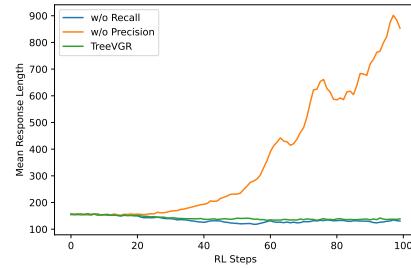


Figure 14: Mean response length with different IoU rewards. The precision term is crucial for alleviating the repetition problem.



Figure 15: Qualitative examples (first three rows) and failure cases (the last row) on V\* Bench (Wu & Xie, 2024) with attention maps.

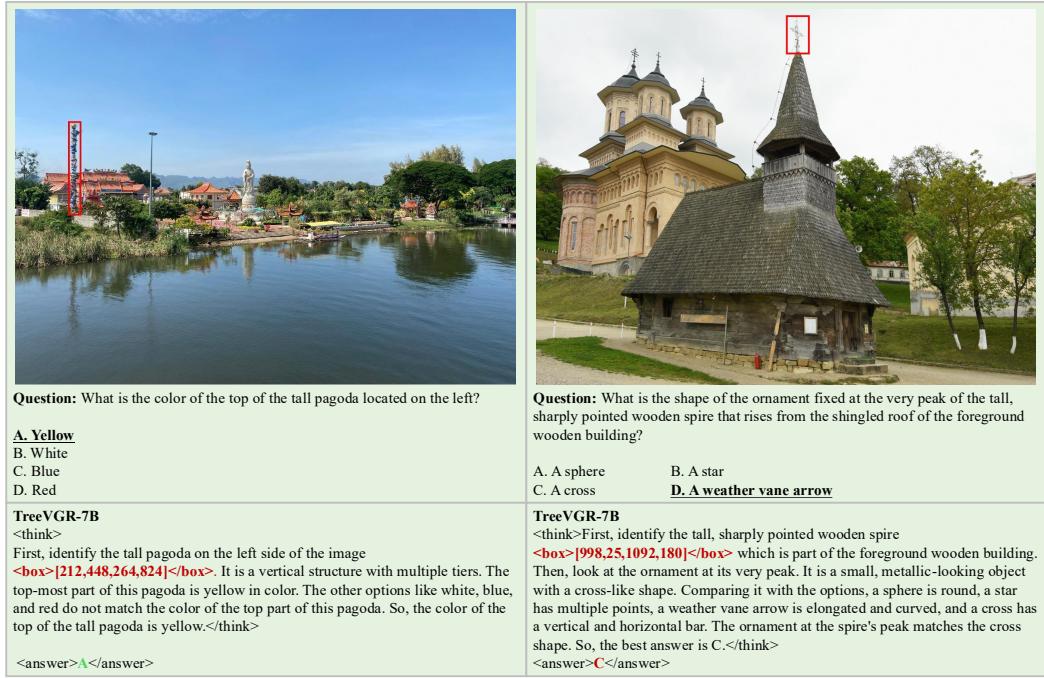


Figure 16: Qualitative examples (left) and failure cases (right) on the “Attributes” protocol of TreeBench.

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158



**Question:** Based on the reflections and sheen, what is the most likely primary material and finish combination of the tall, gooseneck kitchen faucet prominently positioned in the center foreground, above the stainless steel sink?

A. Polished brass  
C. Matte black coated alloy

**B. Brushed stainless steel**  
D. Oil-rubbed bronze



**Question:** The white plate in the foreground, adorned with a light brown geometric pattern and holding the falafel, most likely has which type of surface finish?

**A. High-gloss.**  
C. Textured, non-slip silicone

B. Matte.

D. Glass-like.

**Target Instances:**

**TreeVGR-7B**

<think>

First, identify the tall, gooseneck kitchen faucet **<box>[1804,704,1128,1507]</box>** in the center-foreground above the stainless-steel sink. The faucet has a shiny, reflective surface. The stainless-steel sink also has a reflective surface. The lighting in the image is bright, which enhances the reflective quality. Considering the options: Polished brass would have a warm, golden color and a different sheen. Brushed stainless steel has a light, metallic sheen which matches the faucet's appearance. Matte black coated alloy would have a dark, non-reflective finish, which doesn't match. Oil-rubbed bronze has a dark, matte finish, also not matching. The faucet's finish is consistent with brushed stainless steel, which is a common and durable material for kitchen faucets. </think>

<answer>**B**</answer>

<answer>**B**</answer>

Figure 17: Qualitative examples (left) and failure cases (right) on the “Material” protocol of **TreeBench**.

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175



**Question:** What is the physical posture of the boy in the red shirt on the grassy hill at the far top left?

A. He is standing and watching the performance.  
B. He is lying down flat on the grass.  
C. He is in the middle of running down the hill.  
**D. He is sitting with his knees bent.**

**TreeVGR-7B**

<think>

First, identify the boy in the red shirt on the grassy hill at the far top left. **<box>[120,492,175,571]</box>** He is sitting on the grass with his legs folded under him, which is a typical sitting posture with knees bent. There is no indication that he is standing, lying down, or running. So, the physical posture that matches the description is sitting with his knees bent. </think>

<answer>**D**</answer>



**Question:** What is the current state of the soccer ball located near the lower-right portion of the field?

**A. Rolling along the grass** B. Resting motionless on the ground  
C. Suspended in mid-air D. Trapped under the white-uniformed player's boot

**TreeVGR-7B**

<think>

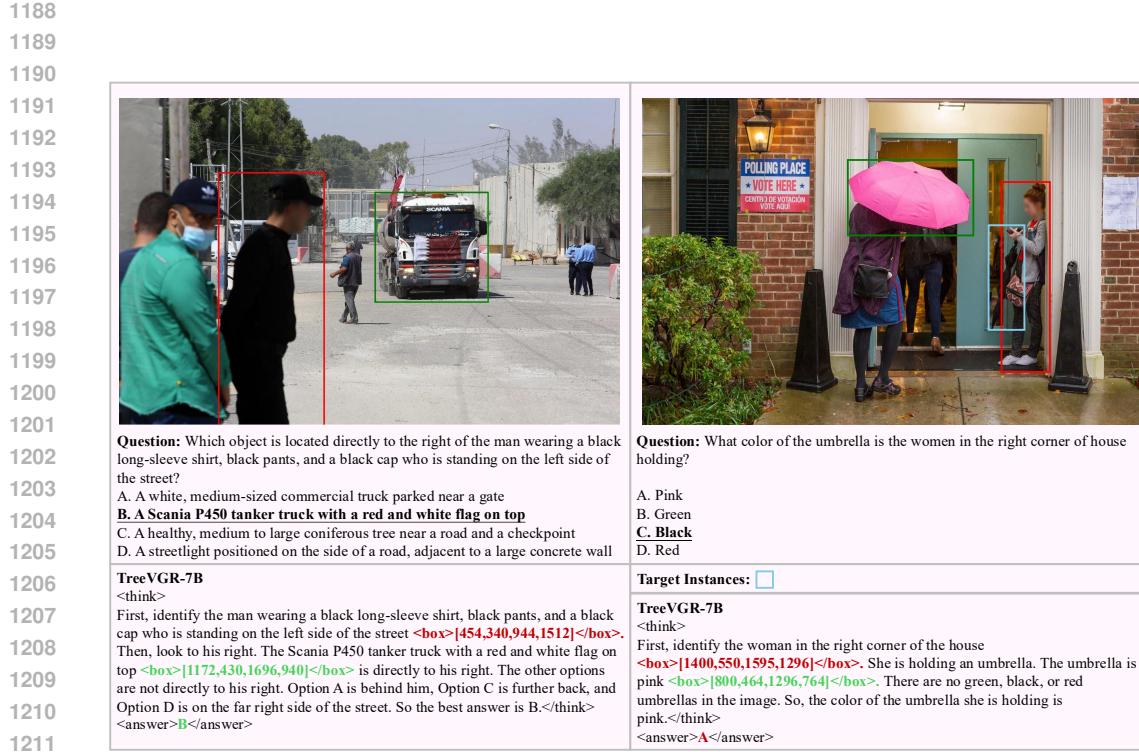
First, identify the soccer ball near the lower-right portion of the field **<box>[1520,1008,1620,1112]</box>**. It is a white ball with blue patterns. Then, observe its position relative to the players. The player in the white uniform (number 89) is in the process of kicking it. The ball is clearly on the ground and in motion as it is being kicked. There is no indication that it is suspended in mid-air, resting, or trapped under the player's boot. So, the ball is rolling along the grass. </think>

<answer>**A**</answer>

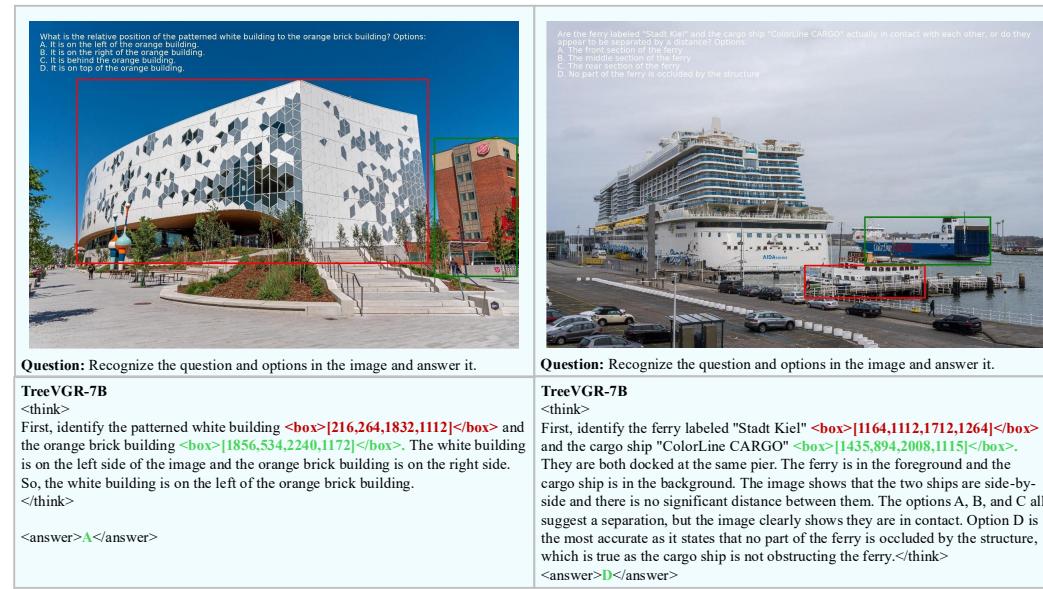
1186

1187

Figure 18: Qualitative examples on the “Physical State” protocol of **TreeBench**.



1212 Figure 19: Qualitative examples (left) and failure cases (right) on the “Object Retrieval” protocol of  
1213 **TreeBench**.



1238 Figure 20: Qualitative examples on the “OCR-Integrated Question-Answering” protocol of  
1239 **TreeBench**.

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265



**Question:** From the perspective of the man wearing a mask and a brown-and-white striped shirt, in which direction is the woman, wearing a white long-sleeve top, dark brown trousers, and carrying a bag on her left shoulder, located?

A. Front left  
 B. Front right  
**C. Left rear**  
 D. Right rear

**TreeVGR-7B**

<think>

First, identify the man wearing a mask and a brown-and-white striped shirt `<box>[584,1014,656,1140]</box>`. Then, find the woman in a white long-sleeve top, dark brown trousers, and carrying a bag on her left shoulder `<box>[724,884,784,1020]</box>`. The woman is positioned behind and to the left of the man. So, from the man's perspective, she is in the rear and to his left. Thus, the woman is in the left-rear direction relative to the man.</think>

<answer>**C**</answer>



**Question:** From the perspective of the third person riding an electric scooter from left to right, in which direction is the red compact SUV located?

A. Front right  
**B. Front left**  
 C. Left rear  
 D. Right rear

**Target Instances:**

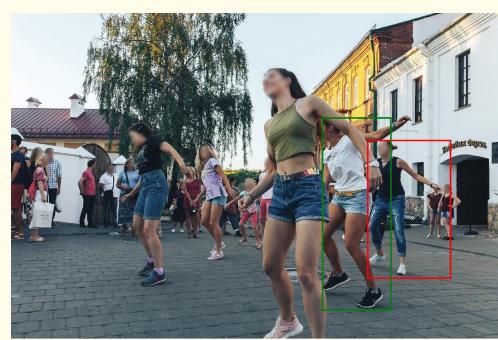
**TreeVGR-7B**

<think>

First, identify the red compact SUV `<box>[1084,1054,1164,1096]</box>` which is parked near the building on the right side of the image. Then, consider the perspective of a third-person riding an electric scooter from left to right. The SUV is positioned behind and to the right of the scooter rider. So, relative to the scooter rider, the SUV is in the rear and to the right. Thus, the SUV is in the right-rear direction.</think>

<answer>**D**</answer>

Figure 21: Qualitative examples (left) and failure cases (right) on the “Perspective Transform” protocol of **TreeBench**.



**Question:** Counting from right to left, which person has a black watch on his/her left hand?

A. The first one  
 B. The second one  
**C. The third one**  
 D. The fourth one

**TreeVGR-7B**

<think>

First, identify the person on the far right in the image `<box>[1640,584,2027,1224]</box>`. This person is wearing a black watch on their left hand. Then, check the person to their left `<box>[1432,480,1752,1371]</box>`, who also has a black watch on their left hand. Since the question asks for the first one, the person on the far right is the first one in the sequence from right to left. So the answer is C.</think>

<answer>**C**</answer>



**Question:** Among the trees behind the pink sign, which one is the shortest when counted from left to right?

A. The first one  
 B. The second one  
 C. The third one  
**D. The fourth one**

**Target Instances:**

**TreeVGR-7B**

<think>

First, identify the pink sign `<box>[844,832,1104,912]</box>` in the image. Then, look at the trees behind it. The first tree on the left is a tall, narrow evergreen tree. The second tree is also tall and narrow. The third tree is shorter than the first two and has a more rounded shape. The fourth tree is taller than the third one. So, the third tree is the shortest among them when counted from left to right.</think>

<answer>**C**</answer>

Figure 22: Qualitative examples (left) and failure cases (right) on the “Ordering” protocol of **TreeBench**.

1293

1294

1295

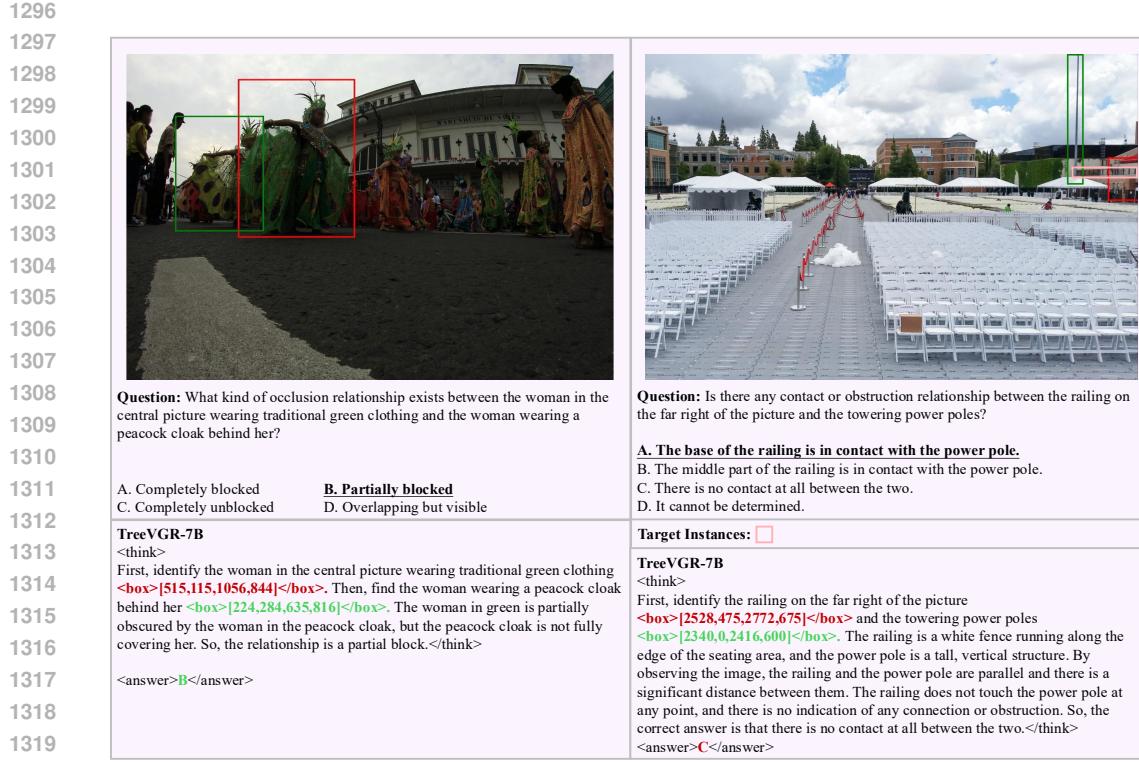


Figure 23: Qualitative examples (left) and failure cases (right) on the “Contact and Occlusion” protocol of **TreeBench**.

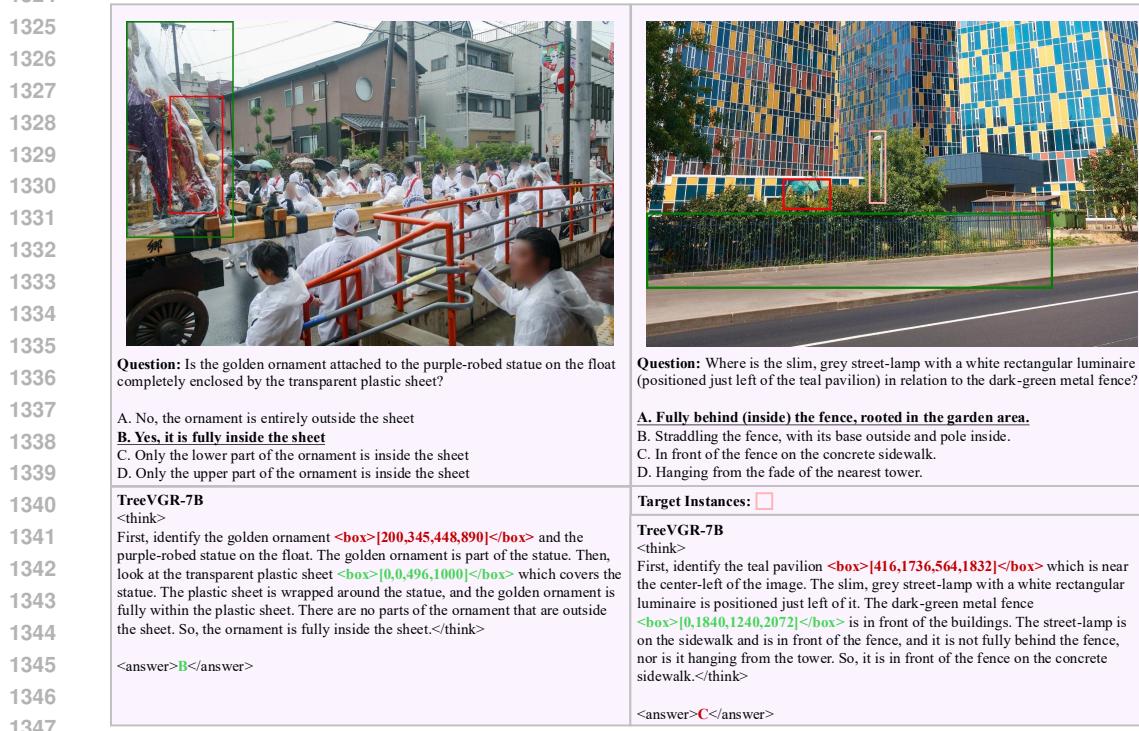


Figure 24: Qualitative examples (left) and failure cases (right) on the “Spatial Containment” protocol of **TreeBench**.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364



**Question:** In the distant skyline, the actual height of the slender, blue-glass, cone-shaped skyscraper towering in the center of the frame, the narrow, dark gray chimney not far to its right, and the milky white building to the left of the frame. Which object in the photo is actually taller?

A. The smokestack appears tallest among the 3 buildings  
B. The milky white building appears tallest among the 3 buildings  
**C. The blue-glass skyscraper appears tallest among the 3 buildings**  
D. All objects are not visible clearly enough to judge height.

**TreeVGR-7B**

<think>

First, identify the three key buildings in the image. The slender, blue-glass, cone-shaped skyscraper `<box>[1144,534,1283,1224]</box>` is the tallest among the three. The narrow, dark gray chimney `<box>[1688,592,1752,816]</box>` is the shortest. The milky white building `<box>[0,0,731,1512]</box>` is in the left-hand side of the frame and is shorter than the skyscraper. So, the skyscraper is the tallest among the three. </think>

<answer>**C**</answer>



**Question:** In the picture, the left side of the bus extends from the back of the car to the front of the car, a total of 10 hand straps, counting from 1 to 10, which is the actual distance and the little girl's head is the closest?

A. 2  
**C. 4**  
D. 8

**TreeVGR-7B**

<think>

First, identify the left-side hand straps in the image.  
`<box>[192,75,288,450]</box> <box>[400,160,475,475]</box>`  
`<box>[550,220,611,490]</box> <box>[772,285,824,525]</box>`  
`<box>[904,360,965,550]</box> <box>[965,385,1032,570]</box>`  
`<box>[1024,410,1075,575]</box> <box>[1352,335,1432,515]</box>`  
`<box>[1424,355,1500,570]</box>` Counting these, there are 10 hand straps. The little girl's head is the closest among them. So, the left-side of the bus from the back of the car to the front of the car is 10 units long, and the little girl's head is the closest. The question asks for the option that matched this distance and the position of the little girl's head. The option 6 is the closest match as it is the closest among the given options. </think>

<answer>**B**</answer>

Figure 25: Qualitative examples (left) and failure cases (right) on the “Comparison” protocol of **TreeBench**.

1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403