

DiG: Differential Grounding for Enhancing Fine-Grained Perception in Multimodal Large Language Models

Zhou Tao^{1,2*}, Shida Wang^{1,2*}, YongXiang Hua^{1,2},
Haoyu Cao^{1,2}, Linli Xu^{1,2†}

¹University of Science and Technology of China

²State Key Laboratory of Cognitive Intelligence

CVPR 2026

- 1 **Motivation** — Fine-grained perception gap in MLLMs
- 2 **Method**
 - DiG task formulation
 - 3D rendering data pipeline
 - Reward design & GRPO optimization
 - Curriculum learning strategy
- 3 **Experiments**
 - Perception, general, and grounding benchmarks
- 4 **Ablation & Analysis**
- 5 **Conclusion**

Motivation: Fine-Grained Perception Gap in MLLMs


Current MLLMs:

- Excel at holistic scene understanding and semantic alignment
- Struggle with **fine-grained visual perception** and **precise spatial reasoning**
- Often overlook subtle visual cues:
 - Small object changes
 - Color variations
 - Missing elements

Root Cause:

Existing pretraining paradigms provide insufficient supervision for developing fine-grained perceptual sensitivity.


Visual Perception QA



Is there a **bottle** in the image?

Qwen3-VL: No. 😞

Ours: Yes. 😊



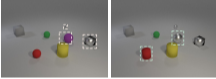
What is the color of the **clock**?

A.black B.yellow C.green D.red

Qwen3-VL: D. red 😞

Ours: C. green 😊

Differential Grounding



Identify all visual differences between the two images and localize them with bounding boxes $[x_1, y_1, x_2, y_2]$ on the before image.

Qwen3-VL: [544, 437, 660, 579], [585, 341, 622, 393], [766, 466, 894, 638]

Ours: [575, 421, 687, 592], [277, 612, 387, 788]

Limitations of Existing Approaches

Direct Fine-tuning on Grounding Tasks:

- Task-specific, poor generalization
- Costly manual annotations (e.g., RefCOCO)
- Limited in diversity and scalability

Alternative Proxy Tasks:

- Visual jigsaw puzzles
- Caption hallucination detection
- Not explicitly designed for fine-grained perception

Our Solution: DiG

Key Insight

Use **differential image comparison** as a proxy task to train fine-grained perception — identify and localize all differences between similar image pairs.

Differential Grounding Task:

Given an input tuple $X = (I_a, I_b, P)$:

- I_a : reference image
- I_b : modified image containing M subtle discrepancies
- P : textual instruction prompting the model to find all differences

Output: A set of predicted bounding boxes $B_{\text{pred}} = \{b_i\}_{i=1}^N$

Objective: Align B_{pred} with ground-truth $B_{\text{gt}} = \{b_j^{\text{gt}}\}_{j=1}^M$ in both **completeness** and **spatial precision**.

Challenge

The model does NOT know the number of differences in advance — it must discover them all.

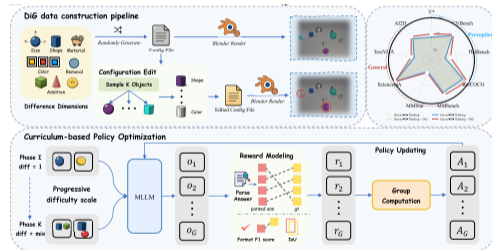
Data Construction Pipeline (3D Rendering)

Automated Pipeline based on Blender:

- 1 Generate base 3D scene via JSON config
- 2 Render reference image I_a
- 3 Apply stochastic modifications:
 - Shape, color, size, material changes
 - Object addition / removal
- 4 Render modified image I_b
- 5 Extract 2D bounding-box annotations from 3D projections

Advantages:

- Fully automated, no manual labeling
- Controllable difficulty (number of differences)
- Perfect ground-truth annotations



Three-component reward for GRPO optimization:

1. **Format Reward** — ensures valid, parseable output:

$$r_{\text{format}} = \begin{cases} 1, & \text{if } \mathcal{P}(O) \neq \emptyset \\ 0, & \text{otherwise} \end{cases}$$

2. **Accuracy Reward** — detection + localization quality:

- Hungarian matching between B_{pred} and B_{gt}
- F1 score: detection completeness
- IoU: spatial alignment precision

$$r_{\text{acc}} = \lambda_1 \cdot \text{F1} + \lambda_2 \cdot \overline{\text{IoU}}$$

3. **Overall Reward:**

$$r = (1 - \alpha) r_{\text{acc}} + \alpha r_{\text{format}}$$

Group Relative Policy Optimization (GRPO):

For each input, sample G candidate responses $\{O^{(g)}\}_{g=1}^G$:

Group-level advantage:

$$A_g = \frac{r_g - \text{mean}(\{r_i\}_{i=1}^G)}{\text{std}(\{r_i\}_{i=1}^G)}$$

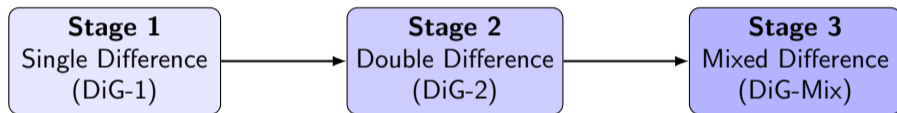
GRPO Loss:

$$\mathcal{L}_{\text{GRPO}}(\theta) = -\mathbb{E} \left[\frac{1}{G} \sum_{g=1}^G \frac{1}{|o^{(g)}|} \sum_{t=1}^{|o^{(g)}|} \mathcal{L}_{g,t}^{\text{CLIP}}(\theta) \right] + \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}})$$

- Clipped surrogate objective for stable updates
- KL regularization to prevent deviation from reference policy

Challenge: Reward sparsity at initialization → unstable training

Solution: Progressive difficulty scheduling



- **Stage 1:** Learn basic difference localization with dense feedback
- **Stage 2:** Develop compositional reasoning and multi-object perception
- **Stage 3:** Generalize across diverse configurations; model infers the number of differences

Training:

- Backbone: Qwen3-VL-4B-Thinking & Qwen3-VL-8B-Thinking
- Dataset: \sim 4.8K image pairs (1.6K per stage)
- Framework: EasyR1 with GRPO

Evaluation Benchmarks:

- **Perception:** HalBench, HRB8K, POPE, V*, VSR, CV-Bench, MMVP
- **General:** MMBench, MM-Vet, MMStar, ScienceQA, TextVQA, MME, AI2D
- **Grounding:** RefCOCO, RefCOCO+, RefCOCOg

Results: Perception Benchmarks

Model	HalBench	HRB8K	POPE	V*	VSR	CV-Bench	MMVP	AVG
Qwen2.5-VL-7B	50.1	65.3	86.4	77.0	77.7	64.9	54.7	68.0
Qwen2.5-VL-72B	55.2	77.1	–	85.9	–	83.0	81.7	–
Qwen3-VL-4B	70.1	69.3	88.3	79.1	82.5	82.8	80.0	78.9
+ DiG (Ours)	73.8 ^{↑3.8}	71.0 ^{↑1.7}	88.5 ^{↑0.2}	79.1	83.9 ^{↑1.4}	83.8 ^{↑1.0}	79.0	79.9 ^{↑1.0}
Qwen3-VL-8B	73.3	69.9	87.9	79.1	82.6	85.0	77.3	79.3
+ DiG (Ours)	76.7 ^{↑3.4}	70.4 ^{↑0.5}	88.0 ^{↑0.1}	81.2 ^{↑2.1}	83.0 ^{↑0.4}	85.8 ^{↑0.8}	78.7 ^{↑1.4}	80.5 ^{↑1.2}

- Consistent improvements on both 4B and 8B models
- Notable gains: **+3.8 HalBench** (4B), **+2.1 V*** (8B)
- Reduced hallucination & enhanced perceptual fidelity

Results: General Multimodal Benchmarks

Model	LLM	MMBench	MMVet	MMStar	SQA ^I	VQA ^T	MME	AI2D
InternVL3.5-38B	Qwen3-32B	87.3	82.2	75.3	–	82.7	–	87.8
Qwen2.5-VL-72B	Qwen2.5-72B	88.6	76.2	70.8	–	84.9	–	83.9
Qwen3-VL	Qwen3-4B	83.7	76.0	66.7	86.9	76.2	1592.4	81.0
+ DiG	Qwen3-4B	84.5 ^{↑0.8}	77.3 ^{↑1.3}	70.2 ^{↑3.5}	89.2 ^{↑2.3}	76.5 ^{↑0.3}	1643 ^{↑51}	82.2 ^{↑1.2}
Qwen3-VL	Qwen3-8B	85.5	76.4	71.7	90.3	75.5	1648.4	82.5
+ DiG	Qwen3-8B	87.2 ^{↑2.2}	76.8 ^{↑0.4}	72.7 ^{↑1.0}	90.5 ^{↑0.2}	77.5 ^{↑2.0}	1666 ^{↑18}	84.1 ^{↑1.6}

- Strong transferability to general multimodal tasks
- 4B model: **+3.5 MMStar**, **+2.3 SQA**, **+51 MME**
- Performance comparable to much larger models (72B)

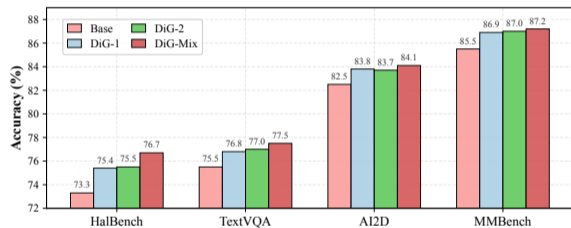
Results: Visual Grounding (RefCOCO)

Method	Size	RefCOCO @50			RefCOCO avg		
		val	testA	testB	val	testA	testB
Qwen3-VL	4B	85.7	90.1	79.1	64.0	67.6	57.1
+ DiG	4B	88.6	91.5	83.5	66.8	70.1	60.8
Qwen3-VL	8B	86.9	90.6	81.7	65.0	68.1	58.8
+ DiG	8B	90.0	91.9	86.4	67.0	69.9	62.2

Method	Size	RefCOCO+ @50			RefCOCOg @50	
		val	testA	testB	val	test
Qwen3-VL	8B	81.8	86.6	73.9	86.3	86.5
+ DiG	8B	84.2	87.7	77.6	87.7	88.0

- Average improvement of **2–3 points** across all grounding benchmarks
- Fine-grained spatial reasoning transfers effectively

Ablation: Curriculum Scheduling



Key Findings:

- **DiG-1** (single diff): Already clear gains over baseline
- **DiG-2** (double diff): Additional improvements on most benchmarks
- **DiG-Mix** (mixed): Best overall performance

⇒ Progressive exposure to complexity yields more robust visual understanding.

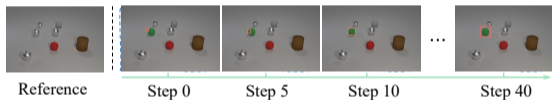
Ablation: Reward Components

Format	IoU	F1	RefCOCO val@50
✓			–
✓	✓		84.2
✓		✓	87.8
✓	✓	✓	88.6





Analysis:

- Format-only: model cannot converge
- IoU-only: insufficient — high IoU without correct detection
- F1-only: better detection, but lacks spatial precision
- **IoU + F1**: best combination — balanced spatial + categorical signals

Training Dynamics & Case Study



Model progressively refines discrepancy localization during RL training.

	<p>Q:Is the right orange circle smaller than the left orange circle? GT: No</p>	
	<p>Qwen3-VL: Yes. Ours: No.</p>	
	<p>Q:Which object is closer to the camera taking this photo, the shelves (highlighted by a red box) or the table (highlighted by a blue box)? GT: A (A) shelves (B) table</p>	
	<p>Qwen3-VL: B. Ours: A.</p>	
	<p>Q:What is the breed of the dog? GT: D (A) Husky (B) corgi (C) Dalmatian (D) golden retriever</p>	
	<p>Qwen3-VL: C. Ours: D.</p>	
		
		<p>Q:What is the color of the man's cap? GT: D (A) white (B) blue (C) red (D) black</p>
		<p>Qwen3-VL: C. Ours: D.</p>

DiG produces more accurate answers for subtle spatial and attribute-based questions.

Contributions:

- 1 **DiG** — a novel proxy task that enhances fine-grained perception via differential image grounding
- 2 **Scalable data pipeline** — automated 3D rendering with controllable discrepancies and perfect annotations
- 3 **Curriculum RL training** — overcomes reward sparsity for stable optimization

Key Results:

- Consistent improvements on perception, grounding, and general benchmarks
- Strong transferability with only $\sim 4.8\text{K}$ training pairs
- Comparable to or exceeds larger models (72B) with 4B/8B scale

Takeaway

Differential grounding offers a **scalable and principled** paradigm for bridging high-level reasoning and low-level perception in MLLMs.

Thank You!

Questions & Discussion

{zhoutao24, wangshida}@mail.ustc.edu.cn