

# Unified Generation and Self-Verification for VLMs via Advantage Decoupled Preference Optimization

# Motivation & Overview: Parallel Test-Time Scaling

## • Motivation

- Serial scaling (o1, R1): long chains, small multimodal gains.
- Parallel scaling: best-of- $N$  works but needs separate generator + verifier.
- Pain: two models  $\Rightarrow$  double data, training, and inference cost.
- Goal: one policy that generates and self-verifies for best-of- $N$ .

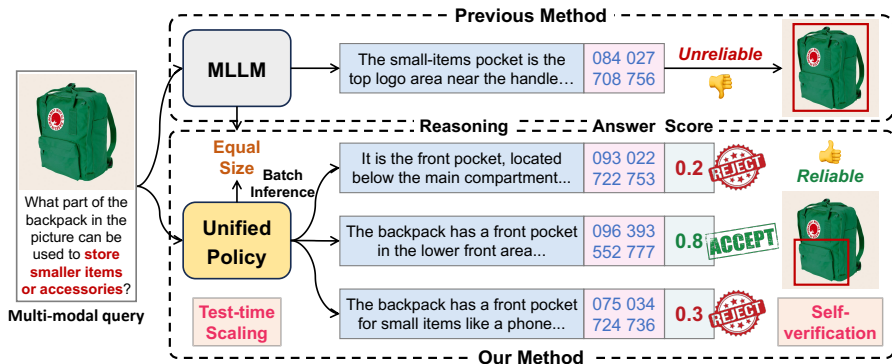
## • Method: Advantage Decoupled Preference Optimization (ADPO)

- Unified framework: one policy that both generates answers and self-verifies.
- Preference verification reward: cast verification as ranking to stay informative under severe class imbalance.
- Advantage decoupled optimization: separate advantages and token masks to avoid reward hacking and disentangle generation vs. verification.

## • Contributions

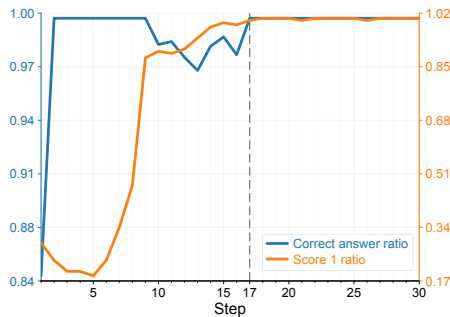
- Preference verification reward
- Advantage decoupled optimization
- Comprehensive evaluation: +2.8/+1.4 acc. on MathVista/MMMU, +1.9 cloU on ReasonSeg, and +1.7/+1.0 step success on AndroidControl/GUI Odyssey.

# ADPO Overview



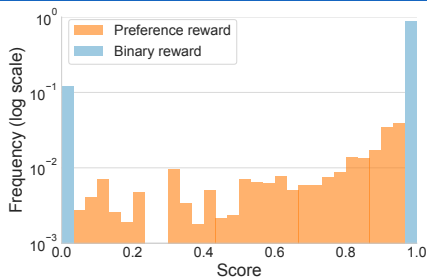
- Extends GRPO to unify answer generation and self-verification within one policy.
- Model outputs an answer plus a verification score  $s \in [0, 1]$  per query.
- Inference: batch decode multiple candidates; pick the answer with the highest self-score.

# Challenges for Unified Self-Verification



- **Class imbalance:** binary verification reward collapses as more answers become correct; scores drift toward a single value, killing gradients.
- **Reward hacking:** summing answer and verification rewards lets the model output bad answers with low self-scores yet still get high total reward.

# Preference Verification Reward



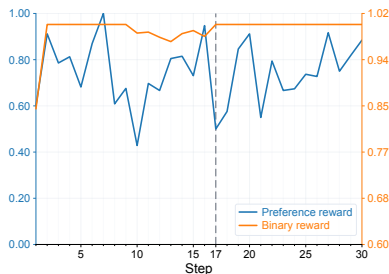
Score distribution collapse without preference reward

Positive/negative ratio imbalance over training

- Reframes verification as a ranking task to avoid collapse under imbalance.
- For sample  $i$ , contrastive set  $\mathcal{C}_i = \{j \mid R_j^a \neq R_i^a\}$  (or  $|R_j^a - R_i^a| > \gamma$  for continuous tasks).
- Reward:

$$R_i^p = \frac{1}{\max(|\mathcal{C}_i|, 1)} \sum_{j \in \mathcal{C}_i} \mathbf{1}\{(s_i - s_j)(R_i^a - R_j^a) > 0\}.$$

- Encourages higher scores for better answers and lower scores for worse ones; works for discrete and continuous rewards.



# Advantage-Decoupled Optimization

- **Entangled advantage**

Sum rewards  $R^a + R^p$  and compute one advantage over all tokens.

Verification gradients leak into generation tokens, enabling reward hacking: bad answers + low self-scores can still get positive total signal.

- **Decoupled advantage**

Compute two advantages  $\hat{A}^{(a)}$  (answer) and  $\hat{A}^{(p)}$  (verification).

- **Unified training objective**

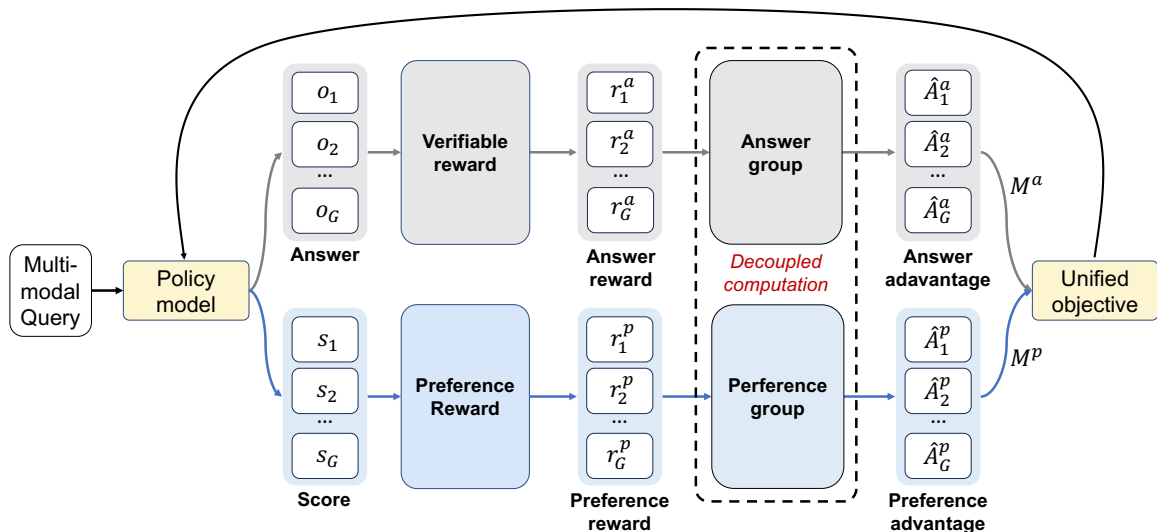
Apply disjoint masks:  $M^a$  on answer tokens,  $M^p$  on the score token.

$$\mathcal{J}(\theta) = M^a \odot \mathcal{J}_\theta(\hat{A}^{(a)}) + M^p \odot \mathcal{J}_\theta(\hat{A}^{(p)}).$$

- **Effect**

Isolates gradients, preserves pass@1 quality, and calibrates self-scores for best-of- $N$ .

# Training Pipeline



# Main Results: Math & Grounding

Table 2. **Performance on MathVista [22] and MMMU [44].** We adopt Qwen2-VL-7B [34] as the base model and use majority voting for both the base and GRPO models. We report accuracy (%) and highlight the best results in **bold**.

Method	MathVista (In-domain)			MMMU (OOD)							
	GVQA	MVQA	ALL	ARD	BUS	HEM	HSS	SCI	TEN	ALL	
<i>Sample 1</i>											
R1-VL-7B [46]	-	-	63.5	-	-	-	-	-	-	-	-
Base	68.9	48.5	57.9	<b>67.5</b>	39.1	49.3	69.0	33.9	36.7	47.1	
GRPO	<b>69.8</b>	55.7	62.2	65.0	45.9	48.2	68.2	<b>35.9</b>	<b>39.8</b>	<b>48.7</b>	
ADPO	68.7	<b>57.0</b>	<b>62.4</b>	63.1	<b>46.2</b>	<b>50.2</b>	<b>71.1</b>	33.3	35.3	47.7	
<i>Sample 4</i>											
MM-Verifier [33]	67.0	53.7	59.8	-	-	-	-	-	-	-	
Base	65.7	51.9	58.2	66.7	47.3	50.7	65.8	34.0	38.1	48.6	
GRPO	69.8	58.0	63.4	65.8	44.7	50.0	<b>70.0</b>	<b>42.0</b>	36.7	49.4	
ADPO	<b>71.3</b>	<b>59.3</b>	<b>64.8</b>	<b>68.3</b>	<b>48.0</b>	<b>52.0</b>	69.2	39.3	<b>39.5</b>	<b>50.8</b>	
<i>Sample 8</i>											
MM-Verifier [33]	68.5	57.4	62.5	-	-	-	-	-	-	-	
Base	68.0	53.3	60.1	<b>68.3</b>	50.0	53.3	68.3	32.7	36.7	49.4	
GRPO	70.4	56.5	62.9	66.7	48.7	51.3	<b>74.2</b>	<b>42.7</b>	36.7	51.1	
ADPO	<b>72.2</b>	<b>58.9</b>	<b>65.0</b>	65.8	<b>54.0</b>	<b>54.7</b>	66.7	40.7	<b>41.0</b>	<b>52.1</b>	
<i>Sample 12</i>											
MM-Verifier [33]	70.4	58.7	64.1	-	-	-	-	-	-	-	
Base	67.4	55.0	60.7	<b>69.2</b>	52.0	50.7	70.8	38.0	36.7	50.7	
GRPO	70.7	57.2	63.4	64.2	50.0	51.3	73.3	<b>43.3</b>	39.5	51.7	
ADPO	<b>71.7</b>	<b>59.8</b>	<b>65.3</b>	67.5	<b>53.3</b>	<b>54.0</b>	<b>71.7</b>	38.7	<b>40.5</b>	<b>52.3</b>	

Table 3. **Performance on ReasonSeg [11].** We use Qwen2.5-VL-7B [1] as the base model.

Method	Short query			Long query			Overall		
	gIoU	cIoU	ACC	gIoU	cIoU	ACC	gIoU	cIoU	ACC
<i>Sample 1</i>									
LISA-7B [11]	47.1	48.5	-	49.2	48.9	-	48.7	48.8	-
SegLLM [37]	-	-	-	-	54.2	-	-	48.4	-
Seg-Zero-7B [18]	-	-	-	-	-	-	57.5	52.0	-
VLM-R1 [32]	-	-	-	-	-	-	-	-	63.1
Base	49.5	53.0	67.0	56.8	57.5	68.5	56.3	57.2	68.4
GRPO	51.8	<b>55.5</b>	67.9	59.1	<b>59.7</b>	71.3	<b>58.6</b>	<b>59.5</b>	71.1
ADPO	<b>51.7</b>	54.8	<b>68.0</b>	<b>60.2</b>	59.4	<b>71.9</b>	58.1	59.1	<b>71.7</b>
<i>Sample 4</i>									
Base	47.8	52.0	66.0	57.3	57.9	69.3	56.7	57.5	69.1
GRPO	<b>54.5</b>	<b>57.0</b>	<b>68.0</b>	58.8	59.5	72.1	58.5	59.4	71.8
ADPO	52.2	55.1	67.0	<b>61.0</b>	<b>61.5</b>	<b>73.3</b>	<b>60.5</b>	<b>61.1</b>	<b>72.9</b>
<i>Sample 8</i>									
Base	47.8	51.4	63.1	57.2	57.8	69.2	56.6	57.4	68.8
GRPO	52.0	55.6	<b>68.0</b>	59.2	59.9	72.0	58.7	59.6	71.7
ADPO	<b>53.2</b>	<b>56.0</b>	67.0	<b>60.9</b>	<b>61.5</b>	<b>73.7</b>	<b>60.4</b>	<b>61.2</b>	<b>73.5</b>
<i>Sample 12</i>									
Base	50.2	53.7	66.0	57.2	57.8	69.3	56.8	57.6	69.1
GRPO	<b>55.6</b>	<b>58.1</b>	<b>69.9</b>	58.8	59.5	72.2	58.6	59.4	72.0
ADPO	53.9	56.2	67.0	<b>61.3</b>	<b>62.0</b>	<b>73.6</b>	<b>60.9</b>	<b>61.6</b>	<b>73.2</b>



# Main Results: Agents & Verifiers

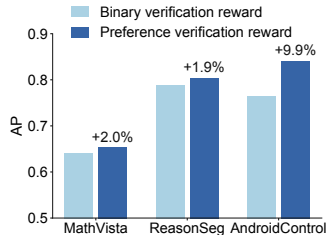
Table 4. **Performance on AndroidControl [13] and GUI Odyssey [23].** We adopt Qwen2.5-VL-7B as base model and report type accuracy, grounding accuracy and step success rate (SR).

Generator	AndroidControl			GUI Odyssey		
	Type	Grounding	SR	Type	Grounding	SR
<i>Sample 1</i>						
UI-TARS-7B [29]	83.7	-	72.5	86.1	-	67.9
SpiritSight-8B [10]	-	-	68.1	-	-	75.8
AgentCPM-GUI-8B [49]	77.7	-	69.2	90.8	-	75.0
Base	82.2	73.6	61.3	81.1	61.4	52.8
GRPO	<b>86.0</b>	<b>76.9</b>	<b>71.0</b>	93.1	<b>83.9</b>	<b>79.8</b>
ADPO	85.8	76.2	70.9	<b>94.2</b>	82.5	79.7
<i>Sample 4</i>						
Base	76.3	68.1	56.0	76.9	55.3	46.5
GRPO	85.5	77.2	71.0	94.7	83.9	81.3
ADPO	<b>86.3</b>	<b>79.5</b>	<b>72.7</b>	<b>94.7</b>	<b>84.5</b>	<b>81.6</b>
<i>Sample 8</i>						
Base	78.7	68.8	58.3	76.7	55.4	46.6
GRPO	85.6	76.9	70.8	94.6	84.4	81.5
ADPO	<b>86.4</b>	<b>78.7</b>	<b>72.7</b>	<b>94.8</b>	<b>84.7</b>	<b>81.7</b>
<i>Sample 12</i>						
Base	78.9	68.7	58.3	76.9	55.5	46.9
GRPO	85.6	77.4	71.1	<b>94.5</b>	84.0	81.1
ADPO	<b>86.3</b>	<b>78.9</b>	<b>72.9</b>	94.4	<b>84.5</b>	<b>81.4</b>

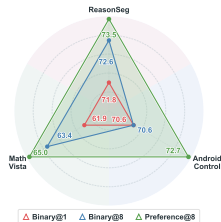
Table 5. **Performance of different generator-verifier settings on MathVista [22], ReasonSeg [11] and AndroidControl [13].**

Generator \ Verifier	MathVista			ReasonSeg			AndroidControl		
	Base	GRPO	ADPO	Base	GRPO	ADPO	Base	GRPO	ADPO
<i>Sample 4</i>									
Base	55.7	55.5	56.4	57.1	57.7	57.7	52.5	57.7	60.7
GRPO	62.4	62.1	62.0	60.2	59.5	60.9	71.0	70.8	71.2
ADPO	61.5	62.1	<b>64.8</b>	59.6	60.3	<b>61.1</b>	71.0	72.0	<b>72.7</b>
<i>Sample 8</i>									
Base	57.0	56.4	56.5	56.9	57.0	57.9	54.3	61.0	64.7
GRPO	60.7	60.8	60.5	60.4	60.4	61.1	71.0	70.9	71.4
ADPO	62.3	62.3	<b>65.0</b>	59.9	60.5	<b>61.2</b>	70.8	71.4	<b>72.7</b>
<i>Sample 12</i>									
Base	56.9	56.3	55.0	57.4	57.6	57.8	53.6	60.7	64.5
GRPO	62.5	62.5	61.8	59.7	59.6	61.3	71.4	70.9	71.5
ADPO	63.0	63.5	<b>65.3</b>	60.7	60.7	<b>61.6</b>	71.6	71.9	<b>72.9</b>

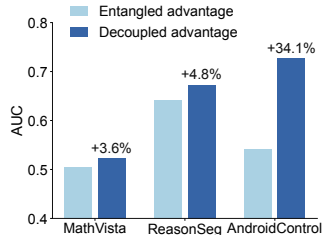
# Ablations: Reward and Advantage



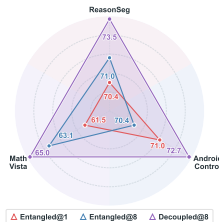
AP improvement (preference reward)



Preference reward vs. binary reward



AUC improvement (decoupled advantage)



Decoupled vs. entangled advantage

# Additional Comparisons

Table 6. Ablation of the margin  $\gamma$  for preference verification reward on ReasonSeg.

$\gamma$	Short query			Long query			Overall		
	gIoU	cIoU	ACC	gIoU	cIoU	ACC	gIoU	cIoU	ACC
0.025	<b>53.7</b>	56.5	<b>69.9</b>	58.1	58.9	71.3	57.8	58.8	71.1
0.050	52.6	54.4	63.1	60.2	61.0	73.3	59.8	60.5	72.7
<b>0.100</b>	53.2	56.0	67.0	<b>60.9</b>	<b>61.5</b>	<b>73.7</b>	<b>60.4</b>	<b>61.2</b>	<b>73.5</b>
0.200	53.2	55.7	66.0	59.9	60.7	72.7	59.6	60.4	72.3
0.250	<b>53.7</b>	<b>56.8</b>	68.9	59.7	60.4	72.5	59.3	60.2	72.3

Table 7. Comparison of unified and separate verification. *GRPO*: GRPO post-trained model as generator. +*Major*: majority voting as verifier. +*Judge*: GRPO post-trained model as verifier.

Method	MathVista Acc. $\uparrow$	Latency (s) $\downarrow$
GRPO+Major	62.9	<b>2.1</b>
GRPO+Judge	60.8	5.6
ADPO	<b>65.0</b>	2.6

# Takeaways

- ADPO unifies generation and verification, enabling reliable parallel test-time scaling with one policy.
- Preference Verification Reward delivers stable, informative gradients under severe class imbalance.
- Advantage Decoupled Optimization isolates gradients and prevents reward hacking while preserving pass@1 quality.
- Stronger best-of- $N$  performance and better-calibrated self-scores across math reasoning, visual grounding, and mobile agents with lower deployment overhead.