



# Steer Where It Matters: Token-Level Visual-Sensitivity Steering for LLMs Hallucination Mitigation



Ruipeng Zhang<sup>1,3</sup>, Zhihao Li<sup>1,2,3</sup>, C.L. Philip Chen<sup>1,2,3</sup>, Tong Zhang<sup>1,2,3\*</sup>

<sup>1</sup>Guangdong Provincial Key Laboratory of Computational AI Models and Cognitive Intelligence, School of Computer Science & Engineering, South China University of Technology

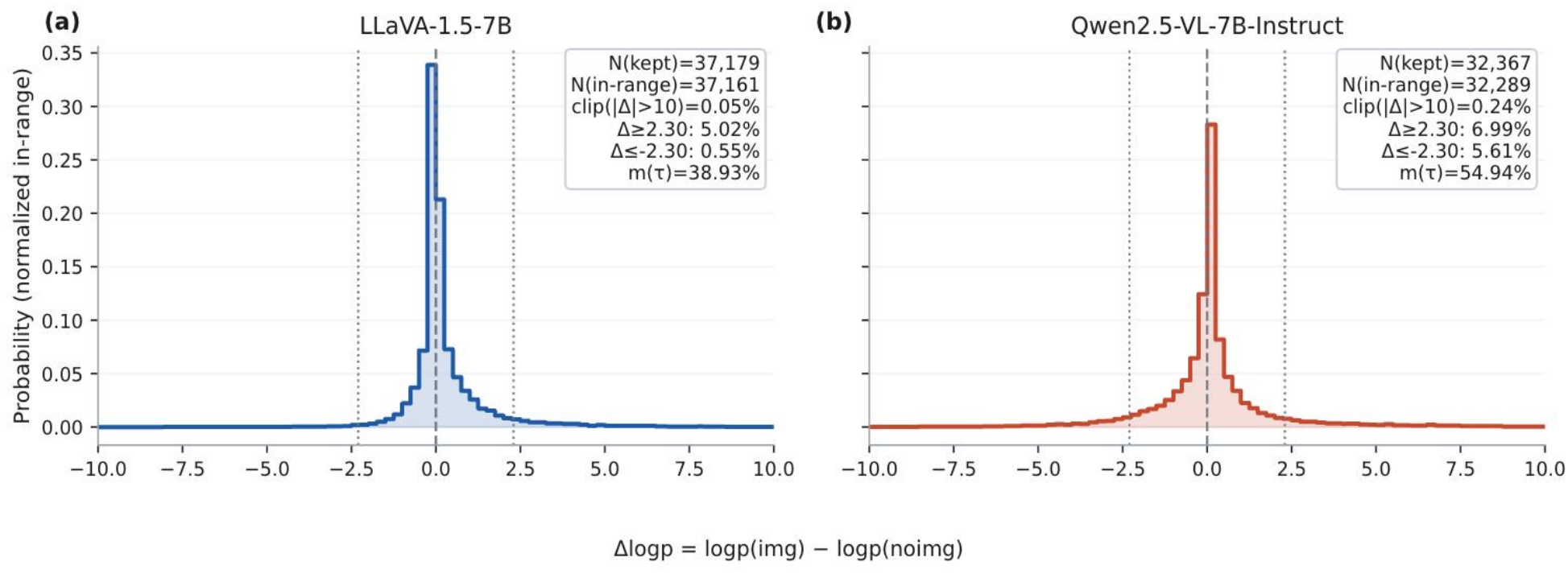
<sup>2</sup>Pazhou Lab, Guangzhou, China

<sup>3</sup>Engineering Research Center of the Ministry of Education on Health Intelligent Perception and Paralleled Digital-Human, Guangzhou, China

## Motivation

## Method

## Partial Experimental Results



Our motivation comes from three observations.

- First, visual information only affects a small number of tokens during generation.
- Second, traditional activation steering averages visual signals across all tokens, which weakens the important image-related information.
- Third, applying the same steering signal to every token may change many tokens that do not need adjustment, so a token-level adaptive method is needed.

## Problem discussion

### Observation 1

We design a token-level diagnosis using teacher forcing with and without the image. For each token, we compute

$$\Delta_t = \log p(y_t | x, I, y_{<t}) - \log p(y_t | x, \emptyset, y_{<t}).$$

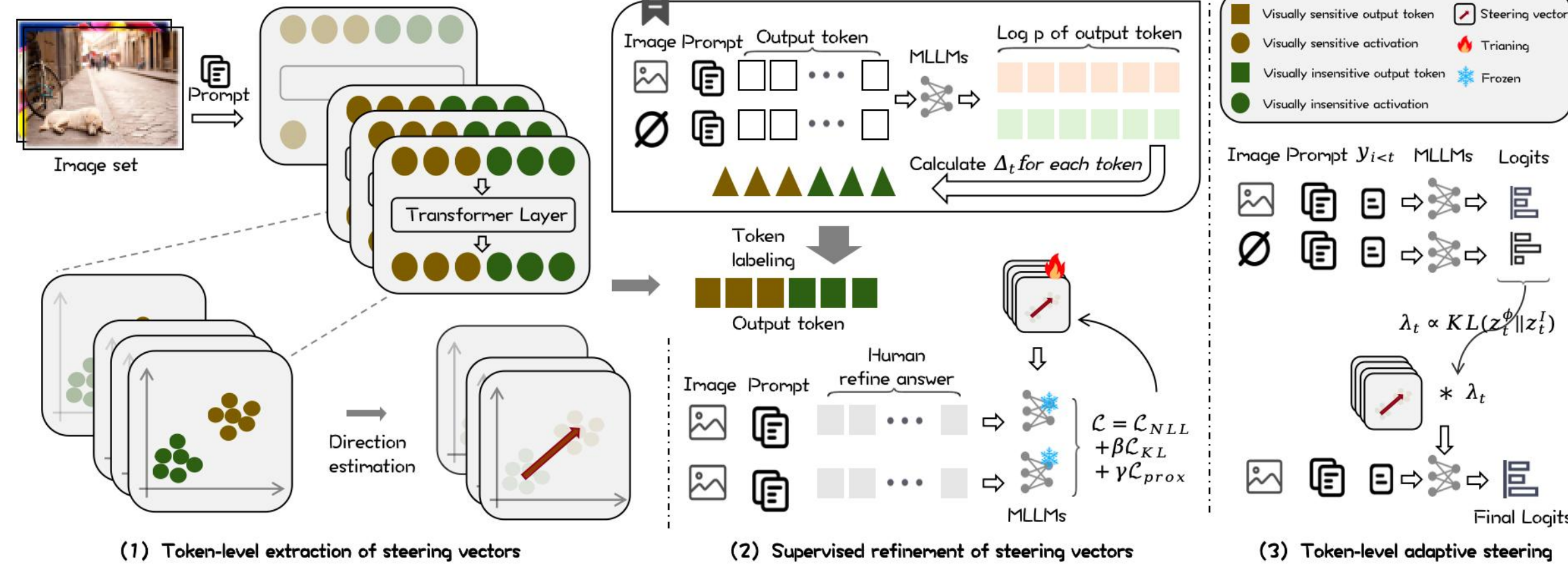
We find that most tokens have ( $\Delta_t \approx 0$ ), while only a few tokens are strongly affected by the image.

### Observation 2

We further test how global steering affects each token using teacher forcing. For each token, we measure

$$\Delta NLL_t(\text{method} - \text{vanilla}) = NLL_t^{\text{method}} - NLL_t^{\text{vanilla}}$$

We find that global steering changes many non-critical tokens, while token-aware steering mainly affects visually important tokens. This shows that global steering misallocates the intervention budget and motivates adaptive token-level steering.



### Step 1: Token-level extraction of steering vectors.

We compare token probabilities with and without the image:

$$\Delta_t = \log p_{\theta}(\hat{y}_t | x, I, y_{<t}) - \log p_{\theta}(\hat{y}_t | x, \emptyset, y_{<t}). \quad \mathcal{T}^+ = \text{TOP}_p(\{|\Delta_t|\}_{t \in \mathcal{T}}), \quad \mathcal{T}^- = \text{BOTTOM}_p(\{|\Delta_t|\}_{t \in \mathcal{T}}).$$

Tokens with large ( $|\Delta_t|$ ) are treated as visually sensitive tokens ( $\mathcal{T}^+$ ), while tokens with small ( $|\Delta_t|$ ) are treated as visually insensitive tokens ( $\mathcal{T}^-$ ). For each layer, we contrast their normalized hidden states to obtain an initial steering direction:

$$d^{(\ell)} = \frac{1}{|\mathcal{T}^+|} \sum_{t \in \mathcal{T}^+} \tilde{h}_t^{(\ell)} - \frac{1}{|\mathcal{T}^-|} \sum_{t \in \mathcal{T}^-} \tilde{h}_t^{(\ell)}. \quad W_{\text{init}}^{(\ell)} = \text{PCA}_1(\{\tilde{d}_i^{(\ell)}\}_{i=1}^N), \quad W_{\text{init}}^{(\ell)} \leftarrow \frac{W_{\text{init}}^{(\ell)}}{\|W_{\text{init}}^{(\ell)}\|_2}.$$

### Step 2: Supervised refinement of steering vectors

We then keep the LLM frozen and only optimize the steering vectors ( $W^{(\ell)}$ ) using correction data:

$$\mathcal{L} = \mathcal{L}_{\text{NLL}} + \beta \mathcal{L}_{\text{KL}} + \gamma \mathcal{L}_{\text{prox}}. \quad \mathcal{L}_{\text{NLL}} = -\frac{1}{|\mathcal{A}|} \sum_{t \in \mathcal{A}} \log p_{\theta, W}(y_t | x, I, y_{<t}).$$

This lightweight refinement improves the steering direction without updating the model parameters.

### Step 3: Adapt steering during decoding.

At each decoding step, we measure visual sensitivity by comparing image-conditioned and image-ablated logits:

$$VS_t = \text{KL} \left( \text{softmax} \left( \frac{z_t^I}{T_{KL}} \right) \parallel \text{softmax} \left( \frac{z_t^{\emptyset}}{T_{KL}} \right) \right), \quad \bar{VS}_t = \frac{VS_t - \mu}{\sigma + \epsilon}.$$

We then map ( $VS_t$ ) to a token-specific steering strength ( $\lambda_t$ ), and apply

$$\tilde{\lambda}_t = \lambda_{\min} + (\lambda_{\max} - \lambda_{\min}) g_t. \quad \hat{\lambda}_t = \alpha \hat{\lambda}_{t-1} + (1 - \alpha) \tilde{\lambda}_t. \quad \lambda_t = \min(\hat{\lambda}_t, \lambda_t^{\text{cap}}).$$

Thus, TLVS applies stronger steering to visually important tokens and weaker steering to non-critical tokens.

Table 1. Main results. TLVS consistently improves hallucination-related metrics across POPE, HallusionBench, CHAIR (COCO), and AMBER on both LLaVA-1.5-7B and Qwen2.5-VL-7B.  $\uparrow/\downarrow$  indicate that higher/lower is better; **bold/underline** mark the best/second-best within each backbone; gray highlights TLVS, with  $\uparrow/\downarrow$  showing absolute gains over the vanilla baseline.

Model	POPE				HallusionBench				CHAIR (COCO)				AMBER					
	Adv. $\uparrow$	Pop. $\uparrow$	Rand. $\uparrow$	Avg. $\uparrow$	qAcc $\uparrow$	fAcc $\uparrow$	aAcc $\uparrow$	CHAIR $\downarrow$	CHAIR $\downarrow$	Recall $\uparrow$	CHAIR $\downarrow$	Cover $\uparrow$	Hal $\downarrow$	Cog $\downarrow$	F1-Et	F1-A $\uparrow$	F1-R $\uparrow$	F1 $\uparrow$
LLaVA-1.5-7B	81.80	84.36	89.12	85.09	13.19	16.18	45.70	46.0	12.8	78.9	12.0	50.3	36.4	4.6	83.2	65.6	62.4	74.7
+VCD	81.33	85.06	87.16	84.52	12.75	19.08	46.41	47.8	14.1	<b>82.7</b>	10.1	51.2	43.6	4.3	84.1	63.6	66.4	74.8
+SPARC	81.03	86.05	89.00	85.36	13.19	17.34	35.25	43.4	16.8	61.1	14.5	44.6	<u>23.1</u>	2.4	88.5	66.7	65.5	77.9
+ASD	—	—	—	<u>87.87</u>	—	—	—	40.0	11.3	82.0	—	—	—	—	—	—	—	—
+SHARP	79.00	83.20	85.30	82.50	—	—	—	<u>34.8</u>	<u>10.6</u>	59.7	8.5	<b>52.1</b>	39.2	4.8	—	—	—	—
+VTI	80.43	85.53	88.50	84.82	13.19	15.90	44.11	35.8	11.1	76.8	4.7	47.2	27.0	1.8	85.8	67.5	<b>68.3</b>	77.1
+VISTA	80.98	88.87	89.25	86.37	14.29	19.36	45.82	44.6	13.4	77.0	6.0	39.9	27.1	2.4	87.6	67.3	55.9	77.6
+DPO (RLHF-V)	<u>84.03</u>	85.94	89.12	86.36	<b>16.70</b>	<b>20.81</b>	<b>51.31</b>	43.2	11.9	76.4	5.7	49.7	27.3	2.6	90.7	72.6	64.6	80.9
+TLVS (ours)	<b>87.85</b>	<b>89.42</b>	<b>92.08</b>	<b>89.78</b>	<b>15.83</b>	<b>19.95</b>	<b>46.50</b>	<b>31.2</b>	<b>9.7</b>	<b>82.2</b>	<b>3.8</b>	<b>51.8</b>	<b>19.7</b>	<b>1.5</b>	<b>91.5</b>	<b>73.0</b>	<b>69.7</b>	<b>82.2</b>
Qwen2.5-VL-7B	84.92	86.78	89.01	86.90	23.74	31.21	55.54	24.0	7.0	<b>86.7</b>	5.4	<b>51.6</b>	29.0	1.9	97.3	<b>83.8</b>	72.6	89.4
+VCD	83.92	85.80	88.41	86.04	31.43	35.26	61.82	29.8	9.7	58.0	6.2	50.3	36.5	2.1	91.6	78.6	68.0	83.8
+DPO (RLHF-V)	<u>84.99</u>	<u>86.81</u>	<u>89.72</u>	<u>87.17</u>	<b>34.95</b>	<b>39.60</b>	<b>62.62</b>	<b>13.4</b>	<b>3.5</b>	<b>62.7</b>	<b>3.1</b>	<b>50.7</b>	<b>21.9</b>	<b>1.1</b>	<b>98.5</b>	82.5	65.9	88.3
+TLVS (ours)	<b>85.10</b>	<b>87.01</b>	<b>89.88</b>	<b>87.33</b>	<b>33.08</b>	<b>37.19</b>	<b>63.23</b>	<b>12.0</b>	<b>5.2</b>	<b>83.9</b>	<b>3.8</b>	<b>51.2</b>	<b>16.7</b>	<b>1.0</b>	<b>97.3</b>	<b>73.0</b>	<b>69.7</b>	<b>89.8</b>

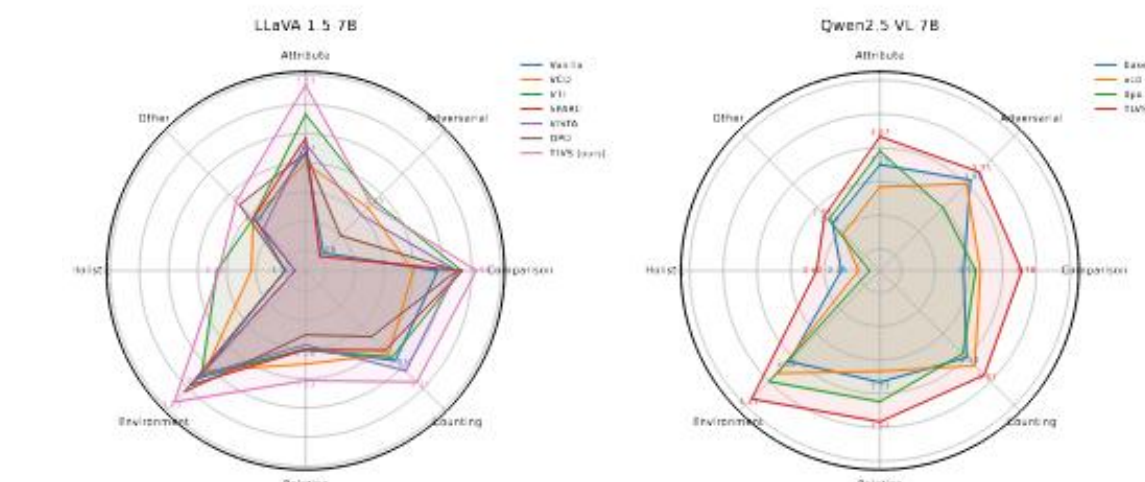


Figure 7. MMHal-Bench performance across question categories (attributes, adversarial objects, comparisons, counting, spatial, environmental, holistic, and others).

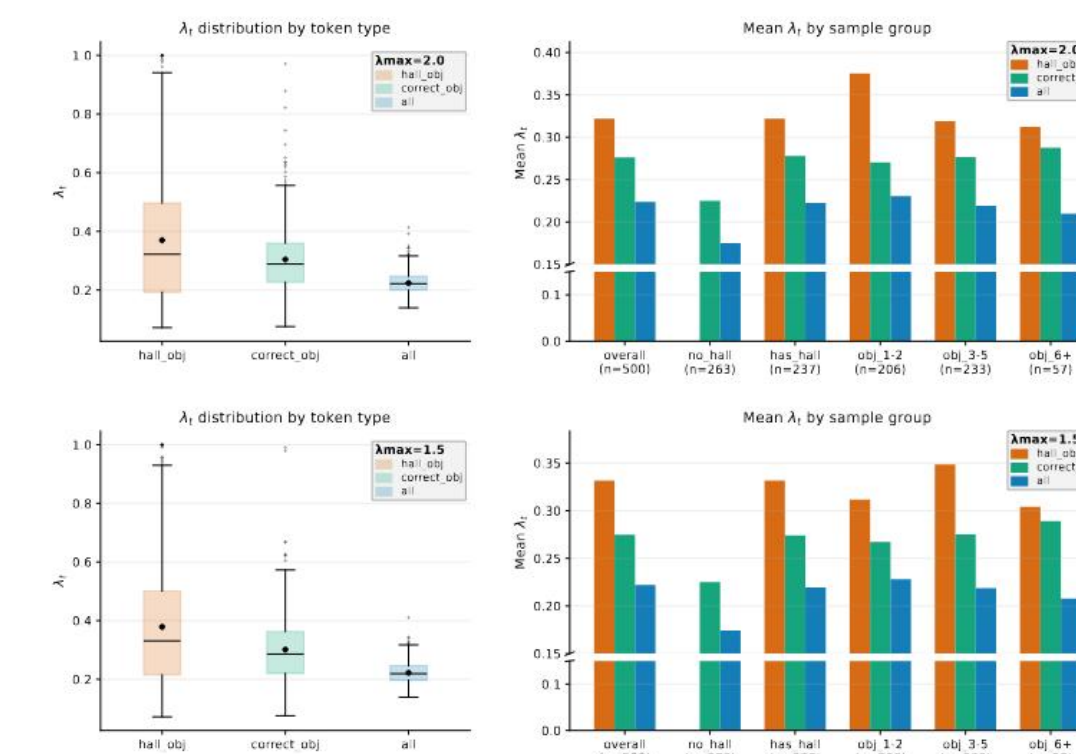


Figure 8. Adaptive steering focuses on critical tokens.  $\lambda_t$  distribution (left) and mean  $\lambda_t$  by token group (right).

Table 2. Ablation on refinement and token-level adaptive steering.  $\uparrow/\downarrow$  indicate higher/lower is better; green/red arrows show absolute gains over the vanilla baseline within each backbone.

Method	Components		POPE				AMBER			
	Refine	Adapt.	Avg $\uparrow$	CHAIR $\downarrow$	Cover $\uparrow$	Hal $\downarrow$	F1 $\uparrow$			
<b>LLaVA-1.5-7B</b>										
Vanilla	—	—	85.09	12.0	50.3	36.4	74.7			
TLVS-Init	×	×	86.83 $\uparrow$ 1.74	7.4 $\downarrow$ 4.6	45.4 $\downarrow$ 4.9	35.9 $\downarrow$ 0.5	75.0 $\uparrow$ 0.3			
TLVS-Refine	✓	×	88.71 $\uparrow$ 3.62	4.1 $\downarrow$ 7.9	44.5 $\downarrow$ 5.8	28.1 $\downarrow$ 8.3	80.8 $\uparrow$ 6.1			
TLVS	✓	✓	89.78 $\uparrow$ 4.69	3.8 $\downarrow$ 8.2	51.8 $\uparrow$ 15	19.7 $\downarrow$ 16.7	82.2 $\uparrow$ 7.5			
<b>Qwen2.5-VL-7B</b>										
Vanilla	—	—	86.90	5.4	51.6	29.0	89.4			
TLVS-Init	×	×	86.83 $\downarrow$ 0.07	5.0 $\downarrow$ 0.4	50.8 $\downarrow$ 0.8	26.8 $\downarrow$ 2.2	87.6 $\downarrow$ 1.8			
TLVS-Refine	✓	×	87.00 $\uparrow$ 0.10	4.2 $\downarrow$ 1.2	48.9 $\downarrow$ 2.7	20.0 $\downarrow$ 9.0	88.8 $\downarrow$ 0.6			
TLVS	✓	✓	87.33 $\uparrow$ 0.43	3.8 $\downarrow$ 1.6	51.2 $\downarrow$ 0.4	16.7 $\downarrow$ 12.9	89.8 $\uparrow$ 0.4			

Table 3. Efficiency and Overhead. (i) Fwd/Token, the average number of forward passes per generated token; (ii) ms/Token, wall-clock latency per token; (iii) Sec/Sample, end-to-end generation time per sample; and (iv) Peak VRAM, the maximum GPU memory footprint during inference. Lower is better for all four metrics ( $\downarrow$ ).

Method	Fwd/Token $\downarrow$	ms/Token $\downarrow$	Sec/Sample $\downarrow$	Peak VRAM (GB) $\downarrow$
Base	1.00	31.25	3.34	13.95
TLVS	2.00	64.01	6.79	14.03
Overhead (Ours/Base)	2.00×	2.05×	2.04×	1.01×

E-mail: 202421044797@mail.scut.edu.cn  
 Paper: <https://arxiv.org/abs/2606.07647>  
<https://zrpchuang.github.io/>  
 Looking for PhD opportunities for Fall 2027!