

# Time-R1: Post-Training Large Vision Language Model for Temporal Video Grounding

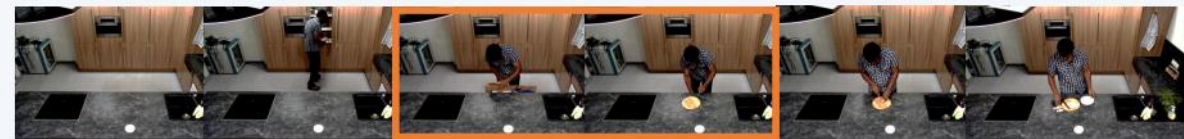
- AIM3 Lab, Renmin University of China
- MiLM Plus, Xiaomi Inc

# Introduction



- **Challenge:** Large Vision-Language Models (LVLMs) underperform on Temporal Video Grounding (TVG) because Supervised Fine-Tuning (SFT) over-penalizes reasonable timestamp deviations via autoregressive loss, leading to poor generalization.
- **Solution:** We propose Time-R1, a Reinforcement Learning (RL) based post-training framework that directly optimizes the task metric (IoU) to overcome SFT's limitations.
- **Contributions:**
  - **Time-R1:** An RL framework to enhance TVG capabilities.
  - **TimeRFT:** A data-efficient RL fine-tuning strategy and dataset.
  - **TVGBench:** A balanced benchmark designed for evaluating LVLMs on TVG.

## Improved Temporal Video Grounding (video length: 148s)



**Question:** Where is the event "He took out cutting board".

Ground Truth	13.5s	18.1s
TimeSuite		20.0s   23.0s
Gemini-2.5-Pro	15.0s	16.0s
VideoChat-Flash	16.4s	19.2s
Trace	15.0s	20.0s
Time-R1 (ours)	13.4s	18.2s

**Time-R1 thinking (ours):** The event "He took out cutting board" occurs between 13.40 and 18.20 seconds. This is when the person is seen placing the cutting board on the counter.

## Improved Short Video QA (video length: 25s)



**Question:** What letter is likely to be shown next?

(A) Y (B) H (C) Z

### Base model thinking:

... Given this sequence, the next logical letter would be the next letter in the alphabet after 'F', which is 'G'. The letters 'D', 'E', 'F', and 'G' are now on the table, forming the sequence 'DEFG'. Therefore, the next letter to be shown is likely to be 'G'.

### Time-R1 thinking (ours):

The sequence of letters shown in the video is "D", "E", "F", and "G". This suggests a pattern of showing the next letter in the alphabet. The letters are being placed in a specific order, and the next logical letter after "G" would be "H".

## Improved Long Video QA (video length: 59m)



**Question:** What do the expanding red lines on the map in the first few minutes of the video stand for?

(A) The Yellow River (B) The Silk Road (C) Du Fu's route to Xi'an (D) The Yangtze River

### Base model thinking:

... The map is used to illustrate the historical journey of Du Fu, a famous Chinese poet, as indicated by the text "Du Fu's route to Xi'an" and the historical context provided by the narrator. Given this analysis, the expanding red lines on the map in the first few minutes of the video stand for Du Fu's route to Xi'an....

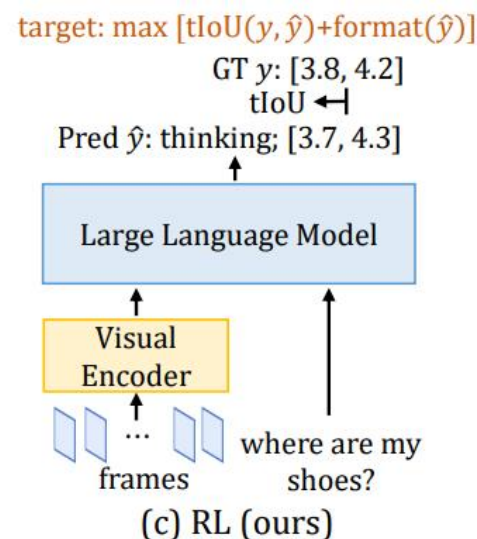
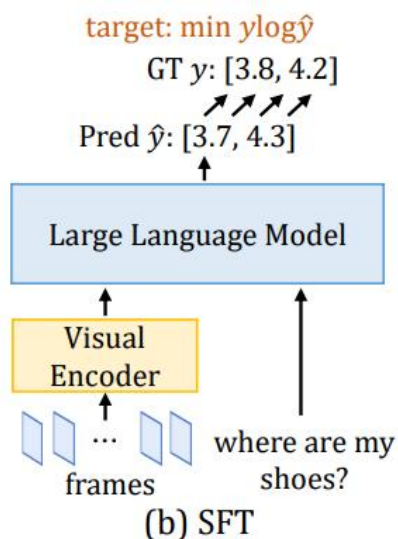
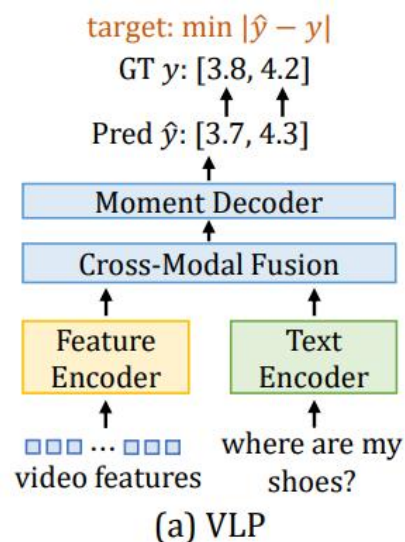
### Time-R1 thinking (ours):

... 2. The map is shown with red lines expanding from a central point, which is likely Xi'an, a major city on the Silk Road. 3. The map is labeled with "Silk Road" and "The Great Silk Road," confirming that the red lines represent the Silk Road. ...

# Time-R1: Reinforcement Learning Framework



- **Core Idea:** Shift from token-matching SFT to outcome-driven RL.
- **Process:** Video/Text Input -> LVLM -> <think>Reasoning</think><answer>Timestamp</answer>
- **Reward Function:** Total Reward =  $r_{\text{tIoU}}$  (Timestamp-aware IoU) +  $r_{\text{form}}$  (Format Reward)
  - **$r_{\text{tIoU}}$ :** Standard IoU plus a penalty for timestamp center deviation, ensuring more precise localization.
  - **$r_{\text{form}}$ :** Encourages the structured "think-then-answer" output format.



$$\text{IoU} = \frac{[t_s, t_e] \cap [t'_s, t'_e]}{[t_s, t_e] \cup [t'_s, t'_e]}$$

$$r_{\text{tIoU}}(o) = \text{IoU} \cdot \left(1 - \frac{|t_s - t'_s|}{t}\right) \cdot \left(1 - \frac{|t_e - t'_e|}{t}\right)$$

# TimeRFT: Efficient RL Fine-Tuning



- **Data Efficiency:** Fine-tuned on only 2.5K samples selected from 339K, focusing on medium-difficulty instances ( $\text{IoU} \approx 0.3$ ).
- **Dynamic Hard Sampling:** A multi-epoch strategy that filters out easy samples ( $\text{IoU} > 0.7$ ) after each epoch to focus on challenging cases.
- **Cold Start:** An initial SFT phase with a few CoT examples to suppress hallucinations, stabilize training, and reduce reasoning length.

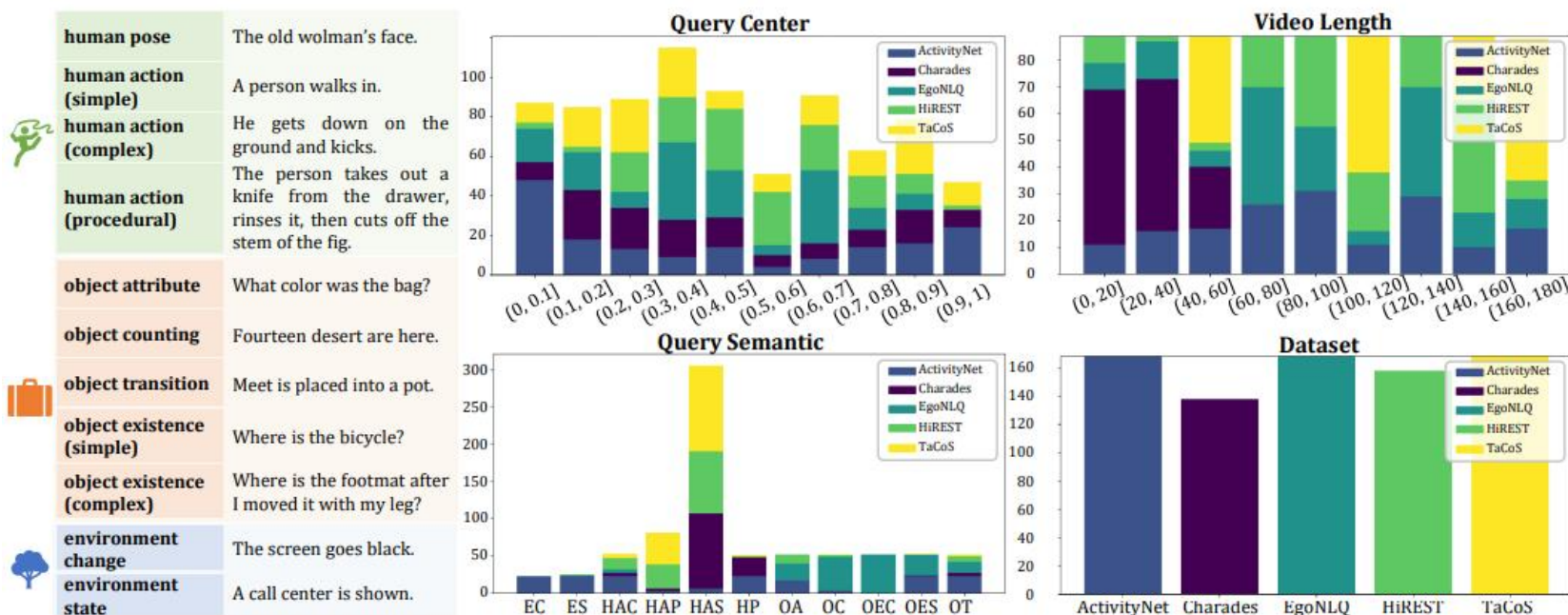




# TVGBench: A Balanced Benchmark



- **Features:** Lightweight (800 samples), with balanced distributions of data sources, video durations, and query temporal locations.
- **Innovation:** First to introduce 11 query semantic categories for fine-grained model analysis.

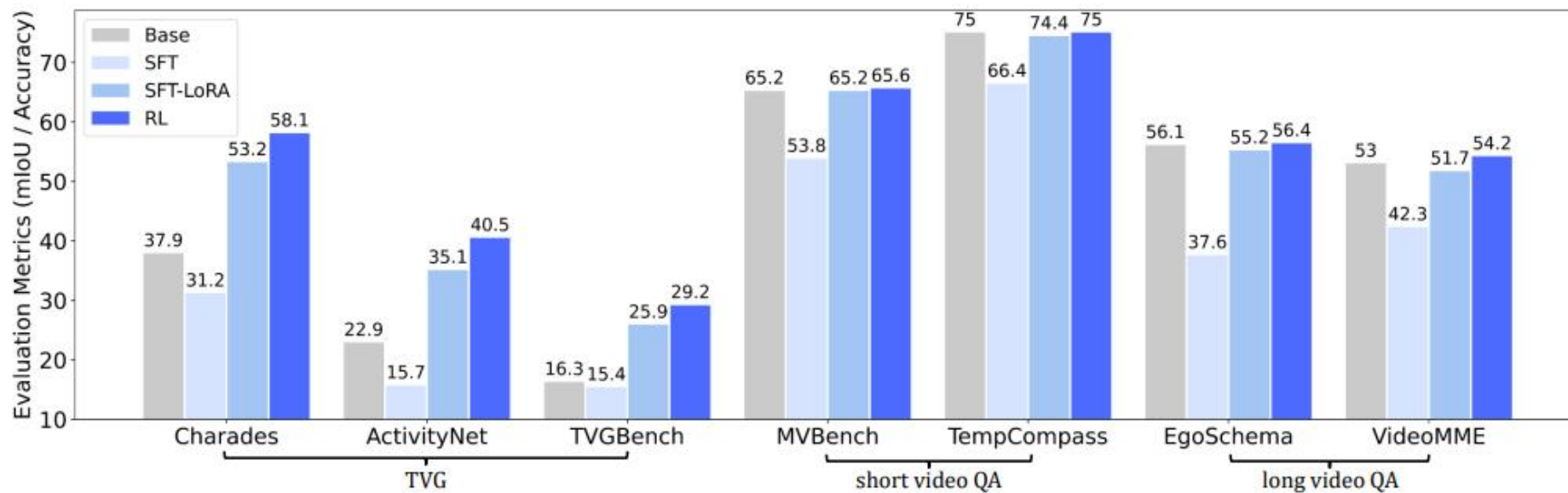


## State-of-the-Art Performance

- **Zero-Shot:** Time-R1 (trained on 2.5K data) outperforms all SFT-based models and surpasses Gemini-2.5-Pro on TVGBench (41.8 vs. 39.1, R1@0.3).
- **Fine-Tuned:** Achieves top performance on Charades-STA and ActivityNet, exceeding all prior LVLMs and most classic VLP methods.

Type	Method	Charades-STA			ActivityNet			TVGBench		
		R1@0.3	R1@0.5	R1@0.7	R1@0.3	R1@0.5	R1@0.7	R1@0.3	R1@0.5	R1@0.7
VLP	2D-TAN* [68]	57.3	45.8	27.9	60.4	43.4	25.0	-	-	-
	UniVTG* [31]	72.6	60.2	38.6	56.1	43.4	24.3	-	-	-
	SSRN* [71]	-	65.5	42.6	-	54.5	33.2	-	-	-
	SnAG* [39]	-	64.6	46.2	-	48.6	30.6	-	-	-
	EaTR* [23]	-	68.4	44.9	-	58.2	37.6	-	-	-
	Gemini-2.5-Pro [10]	-	-	-	-	-	-	39.1	24.4	12.8
SFT	Momentor [43]	42.6	26.6	11.6	42.9	23.0	12.4	-	-	-
	ChatVTG [44]	52.7	33.0	15.9	40.7	22.5	9.4	-	-	-
	TimeChat [47]	-	32.2	13.4	36.2	20.2	9.5	22.4	11.9	5.3
	VTG-LLM [18]	-	33.8	15.7	-	-	-	-	-	-
	HawkEye [53]	50.6	31.4	14.5	49.1	29.3	10.7	-	-	-
	VTimeLLM [22]	51.0	27.5	11.4	44.0	27.8	14.3	-	-	-
	VideoChat-TPO [57]	58.3	40.2	18.4	-	-	-	-	-	-
	VideoExpert [70]	61.5	40.3	20.9	-	-	-	-	-	-
	TimeSuite [65]	69.9	48.7	24.0	-	-	-	31.1	18.0	8.9
	VideoMind [32]	73.5	59.1	31.2	48.4	30.3	15.7	-	-	-
	VideoChat-Flash [28]	74.5	53.1	27.6	-	-	-	32.8	19.8	10.4
	TRACE [19]	-	40.3	19.4	-	37.7	24.0	37.0	25.5	14.6
	HawkEye* [53]	72.5	58.3	28.8	55.9	34.7	17.9	-	-	-
	TimeSuite* [65]	79.4	67.1	43.0	-	-	-	-	-	-
	VideoChat-TPO* [57]	77.0	65.0	40.7	-	-	-	-	-	-
	VideoExpert* [70]	74.3	60.8	36.5	-	-	-	-	-	-
RL	Time-R1 (ours)	<b>78.1</b>	<b>60.8</b>	<b>35.3</b>	<b>58.6</b>	<b>39.0</b>	<b>21.4</b>	<b>41.8</b>	<b>29.4</b>	<b>16.4</b>
	Time-R1 (ours)*	82.8	72.2	50.1	73.3	55.6	34.0	-	-	-

- **Generalization:** RL preserves and improves general video QA capabilities, whereas SFT causes catastrophic forgetting.
- **Data Efficiency:** 2.5K (RL) > 339K (SFT-LoRA).



- **Framework Generality:** Consistently effective across various model architectures and sizes (Qwen-VL, MiMo-VL, InternVL).
- **Component Effectiveness:** Ablation studies confirm the critical roles of the TimeRFT strategy, rtIoU reward, and cold start in boosting performance.

Model	Type	R1@0.3	R1@0.5	R1@0.7
Qwen-2.5-VL-3B [4]	Base	11.5	6.5	3.8
	Time-R1	<b>33.5</b>	<b>21.0</b>	<b>10.5</b>
Qwen-2.5-VL-7B [4]	Base	24.9	16.0	8.0
	Time-R1	<b>41.6</b>	<b>28.5</b>	<b>15.6</b>
MiMo-VL-7B [56]	Base	22.4	12.6	6.6
	Time-R1	<b>41.2</b>	<b>27.8</b>	<b>15.1</b>
InternVL3-2B [72]	Base	16.3	6.3	2.3
	Time-R1	<b>21.8</b>	<b>9.5</b>	<b>4.1</b>
InternVL3-8B [72]	Base	17.4	8.3	3.4
	Time-R1	<b>38.0</b>	<b>22.5</b>	<b>9.2</b>

	tIoU	GF	ME	SF	TVGBench		
					R1@0.3	R1@0.5	R1@0.7
1	✗	✗	✗	✗	38.0	24.8	13.2
2	✓	✗	✗	✗	36.0	23.6	12.9
3	✗	✓	✗	✗	37.2	25.0	13.4
4	✗	✗	✓	✗	39.9	26.0	14.2
5	✓	✓	✗	✗	38.4	25.6	14.1
6	✓	✗	✓	✗	39.4	26.5	16.4
7	✓	✓	✓	✗	41.6	28.5	15.6
8	✓	✓	✓	✓	41.8	29.4	16.4

Reward Design	R1@0.3	R1@0.5	R1@0.7
$r_{\text{tIoU}} + r_{\text{format}}$ (Ours)	<b>41.8</b>	<b>29.4</b>	<b>16.4</b>
$r_{\text{format}}$ only	27.1	18.0	10.1
$r_{\text{tIoU}}$ (w/o format)	40.5	27.6	15.4
$r_{\text{IoU}} + r_{\text{format}}$	41.4	28.0	15.8
$r_{\text{em}} + r_{\text{format}}$	26.5	16.8	9.1
$r_{\text{abs}} + r_{\text{format}}$	39.1	27.8	14.8
$r_{\text{rmse}} + r_{\text{format}}$	38.9	27.0	15.8
$r_{\text{center}} + r_{\text{format}}$	37.6	25.9	15.0

