

# MUSEG: REINFORCING VIDEO TEMPORAL UNDERSTANDING VIA TIMESTAMP-AWARE MULTI-SEGMENT GROUNDING

Anonymous authors

Paper under double-blind review

## ABSTRACT

Video temporal understanding is crucial for multimodal large language models (MLLMs) to reason over events in videos. Despite recent advances in general video understanding, current MLLMs still struggle with fine-grained temporal reasoning. While reinforcement learning (RL) has been explored to address this issue recently, existing RL approaches remain limited in performance on time-sensitive tasks. In this work, we propose **MUSEG**, a novel RL-based method that enhances temporal understanding by introducing timestamp-aware multi-segment grounding. MUSEG enables MLLMs to align queries with multiple relevant video segments, promoting more comprehensive temporal reasoning. To facilitate effective learning, we design a customized RL training recipe with phased rewards that progressively guides the model toward temporally grounded reasoning. Extensive experiments on temporal grounding and time-sensitive video question answering (QA) tasks demonstrate that MUSEG significantly outperforms existing methods and generalizes well across diverse temporal understanding scenarios.

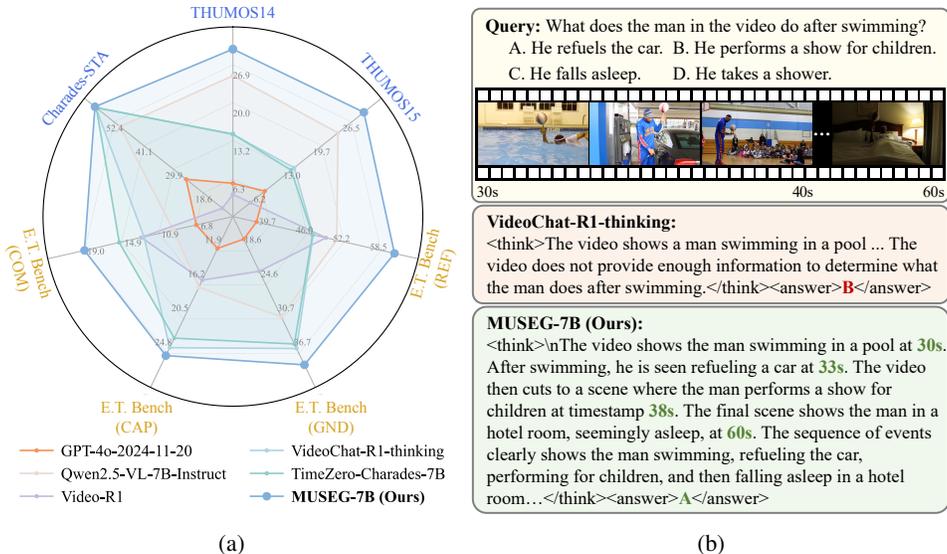


Figure 1: (a) Performance of our MUSEG-7B on various temporal grounding (Charades-STA, THUMOS14 and THUMOS15) and broader time-sensitive video understanding (E.T. Bench Subset) tasks. (b) An example comparing our MUSEG-7B with previous models. MUSEG-7B performs more precise, timestamp-aware reasoning by effectively using multiple key temporal cues to derive the correct answer.

## 1 INTRODUCTION

Video temporal understanding (Liu et al., 2024a; Chen et al., 2024; Cheng et al., 2025b) refers to tasks of comprehending events based on temporal dynamics such as temporal grounding (Gao et al., 2017), dense video captioning (Wang et al., 2024), grounded video question answering (Xiao et al.,

054 2024), etc. This capability is vital for multimodal large language models (MLLMs) (Hurst et al.,  
055 2024; Team et al., 2023; Bai et al., 2025) in understanding complex temporal structures in videos  
056 and making accurate, context-aware predictions or decisions based on when and how events unfold.

057 Despite rapid progress and impressive results in general video understanding, current MLLMs still  
058 show significant limitations in temporal understanding (Liu et al., 2024b; Li et al., 2025c). Early  
059 efforts to address this are mainly based on supervised fine-tuning (SFT) to improve temporal com-  
060 prehension (Bai et al., 2025; Liu et al., 2024a; Li et al., 2025a). As reinforcement learning (RL)  
061 has been shown to significantly improve complex reasoning and comprehension in large language  
062 models (LLMs) (Guo et al., 2025), recent studies have extended RL techniques to the video do-  
063 main (Feng et al., 2025; Li et al., 2025b; Wang et al., 2025b; Zhang et al., 2025b), encouraging  
064 models to “reason before answering”. This typically involves designing a format reward to ensure a  
065 structured reasoning process and an answer reward such as Intersection over Union (IoU) to measure  
066 the correctness of the predictions.

067 However, directly applying RL to video temporal understanding tasks has not achieved the same  
068 level of performance improvement as in textual domains (Feng et al., 2025; Li et al., 2025b). We  
069 attribute this limitation to two key challenges. First, most existing methods (Li et al., 2025b; Wang  
070 et al., 2025b) rely solely on single-segment temporal grounding, where each input query corresponds  
071 to only one video segment. This limits the ability to capture fine-grained, multi-segment temporal  
072 information, which is essential for complex video understanding tasks. Second, although temporal  
073 understanding depends fundamentally on reasoning over temporal cues, current RL approaches of-  
074 ten fail to model them effectively. Reasoning process of current models typically consists of brief  
075 descriptions of video content, lacking detailed temporal analysis of key events, as illustrated in Fig-  
076 ure 1 (b). Therefore, we argue that advancing MLLMs in video temporal understanding requires  
077 rethinking both the *training task design* and the *RL training recipe*.

078 In this paper, we propose timestamp-aware **M**ulti-**S**egment **G**rounding (MUSEG), an RL-based  
079 method designed to enhance the temporal understanding and reasoning capabilities of MLLMs. On  
080 the task side, we incorporate *multi-segment grounding* into the training process, enabling models to  
081 learn from queries that align with multiple relevant video segments. This promotes stronger temporal  
082 understanding and better generalization to a wide range of time-sensitive tasks. On the training side,  
083 we introduce a customized RL training recipe with phased rewards, which progressively encourage  
084 the model to establish temporally grounded reasoning processes. Our recipe features a dedicated  
085 segment matching reward and a timestamp reward, encouraging models to perform fine-grained  
086 temporal reasoning over multiple segments as shown in Figure 1 (b). Additionally, we employ a  
087 multi-phase training strategy that balances guided learning and exploration. As illustrated in Figure 1  
088 (a), MUSEG achieves significant improvements on temporal grounding benchmarks and generalizes  
089 effectively to other time-sensitive video understanding tasks. Our contributions can be summarized  
as follows:

- 090 • We propose MUSEG, a novel RL-based method for video temporal understanding, which en-  
091 ables MLLMs to reason over multiple temporally distributed events by incorporating multi-  
092 segment grounding into training.
- 093 • We design a tailored RL training recipe featuring novel reward functions and a multi-phase  
094 training strategy, effectively promoting fine-grained and temporally grounded reasoning.
- 095 • We conduct extensive experiments and analyses, showing that MUSEG consistently outper-  
096 forms existing methods on video temporal understanding benchmarks, and validating the ef-  
097 fectiveness of our task and training designs.

## 099 2 RELATED WORK

101 **Video Temporal Understanding.** Previous research on video temporal understanding mainly fo-  
102 cuses on cross-references and alignments between videos and texts (Arnab et al., 2021; Luo et al.,  
103 2021; Liu et al., 2021; Xu et al., 2021; Wang et al., 2021). Recent advances in video temporal un-  
104 derstanding have moved from these cross-modal attention-based local feature matching approaches  
105 to broader time-sensitive tasks, such as temporal grounding (Gao et al., 2017), dense video cap-  
106 tioning (Wang et al., 2024), and grounded video question answering (Xiao et al., 2024). These  
107 methods attempt to fuse video temporal features and text features with LLMs to enhance model  
performance (Liu et al., 2024a; Li et al., 2025c; Yan et al., 2025).

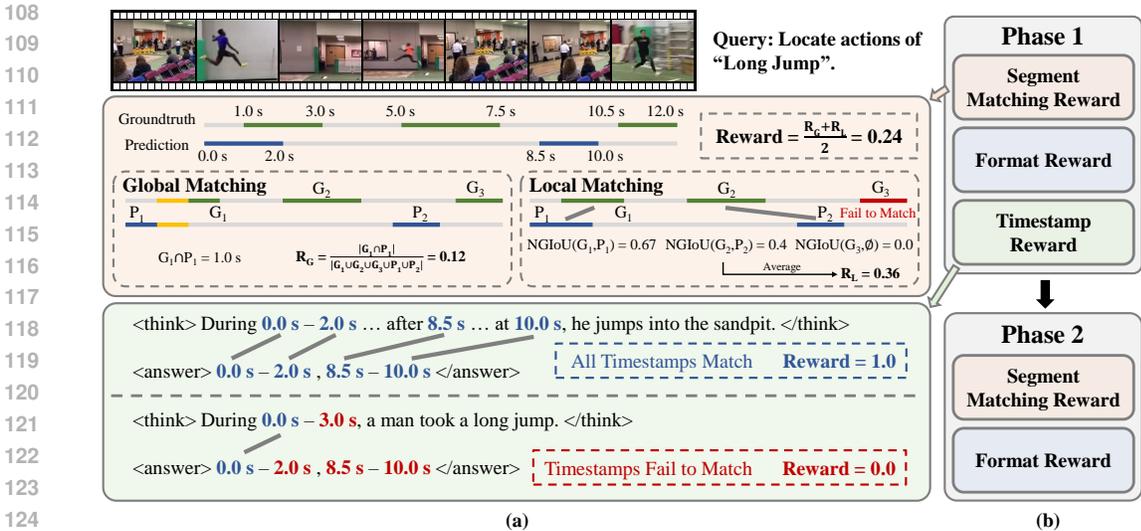


Figure 2: Overview of MUSEG. (a) Our proposed segment matching reward (up) and timestamp reward (down). (b) RL-based training process with phased rewards of MUSEG.

However, these methods involve automatic construction of Chain-of-Thoughts (CoTs) (Wei et al., 2022) in training data. Thus, they are limited by fixed reasoning patterns and potential issues with data quality, and model performance remains suboptimal on temporal understanding tasks, and struggle to generalize to complex scenarios (Liu et al., 2024a; Chen et al., 2024; Cai et al., 2024; Huang et al., 2024). Instead, we empower models to develop temporal-aware reasoning patterns autonomously, rather than constraining them to learn predefined patterns.

**RL for Video Understanding.** RL has been widely adopted in various textual tasks (Shao et al., 2024; Ouyang et al., 2022; Schulman et al., 2017). Recent works apply RL to general video question answering (Feng et al., 2025; Chen et al., 2025b; Dang et al., 2025; Wang et al., 2025a; Chen et al., 2025a; Zhang et al., 2025a) and temporal grounding tasks (Li et al., 2025b; Cheng et al., 2025a). However, these methods primarily provide rewards based on correctness of answers. The effective utilization of temporal information in the reasoning processes is not accounted for in the rewards, which may result in less effective CoTs. Models still struggle on complex temporal grounding tasks, and there is still room for improvement in generalizing to broader temporal understanding scenarios. In contrast, our approach imposes supervision on both reasoning processes and final answers.

### 3 PRELIMINARIES: REWARD DESIGN IN GRPO

Group Relative Policy Optimization (GRPO) (Shao et al., 2024) is a reinforcement learning (RL) training algorithm that has been widely adopted to enhance the reasoning abilities of large language models (LLMs). A key factor in GRPO is the design of the reward function. Prior work has demonstrated that verifiable rewards, though simple, can be highly effective in improving LLM reasoning (Guo et al., 2025; Yang et al., 2025). For instance, DeepSeek-R1 (Guo et al., 2025) incorporates two rule-based rewards:

- **Accuracy Rewards:** Assess whether the model produces correct outputs. For math tasks, correctness is verified by rule-based checkers, while for coding tasks, correctness is validated by compilers combined with pre-defined test cases.
- **Format Rewards:** Ensure that model outputs follow the required structure, namely the “<think>...</think> <answer>...</answer>” format.

While such rewards have proven effective in textual domains, we will show that directly applying them to video temporal understanding tasks leads to suboptimal performance.

## 4 METHOD

In this work, we propose a novel RL-based approach for video temporal understanding. The key components include multi-segment grounding as the training task (Section 4.1), a carefully designed

reward function with timestamp awareness (Section 4.2), and a new training paradigm with phased rewards (Section 4.3).

#### 4.1 MULTI-SEGMENT GROUNDING TASK

Temporal grounding is the task that requires models to match text queries with corresponding video segments, which is capable of improving temporal understanding abilities of MLLMs (Liu et al., 2024a; Bai et al., 2025). Broadly, temporal grounding queries can be categorized into two types: *single-segment grounding*, where each query is associated with a single video segment, and *multi-segment grounding*, where a query may correspond to multiple video segments.

Previous RL-based temporal grounding approaches (Li et al., 2025b; Wang et al., 2025b) typically adopt single-segment grounding as the training task. However, our preliminary empirical study reveals that a notable portion of single-segment grounding questions can be solved through unintended shortcuts, for example, by identifying key objects rather than reasoning about the temporal structure of events. As shown in Table 1, this shortcut accounts for about 30% of cases, based on a manual check of 50 single-segment grounding questions sampled from E.T. Bench (Liu et al., 2024a). These findings suggest that relying solely on single-segment grounding tasks is insufficient for enhancing the temporal understanding abilities of MLLMs.

In contrast, multi-segment grounding queries are difficult to be answered by shortcuts, as shown in Table 1, so we add them to our training process. We ensure the number of single-segment grounding and multi-segment grounding queries are balanced, and our selected data are diverse in scenarios. Experiments demonstrate that incorporating multi-segment grounding significantly enhances model performance (see Section 5.3).

Table 1: Results of preliminary empirical study. We sample single-segment grounding and multi-segment grounding queries from E.T. Bench (Liu et al., 2024a), and examine whether they can be answered by shortcut of recognizing key objects.

Query Type	w/ Shortcut	Total
Single-Segment	15	50
Multi-Segment	4	50

#### 4.2 REWARD DESIGN

We utilize segment matching reward (Section 4.2.1) to evaluate outputs for multi-segment grounding tasks, timestamp reward (Section 4.2.2) to stimulate model ability for temporal-aware reasoning, and format reward (Section 4.2.3) to encourage models to think before providing answers.

##### 4.2.1 SEGMENT MATCHING REWARD

Segment matching reward is designed to align model outputs with ground truths. It consists of two parts, global matching and local matching, to enhance model abilities of understanding overall video contents, and grasping detailed events, respectively.

*Global matching* is shown in upper left area of Figure 2 (a). We measure the overlap ratio among all the ground truth segments  $\{G_i\}$  and predicted segments  $\{P_j\}$ :

$$r_G = \frac{\sum_{i,j} |G_i \cap P_j|}{|(\cup_i G_i) \cup (\cup_j P_j)|}, \text{ where } G_i \text{ and } P_j \text{ are represented as intervals.} \quad (1)$$

In the *local matching* process, we pair ground truths and predictions one-to-one as  $\{(G_n, P_n)\}_{n=1}^N$ , where  $N = \max(|\{G_i\}|, |\{P_j\}|)$ . As shown in upper right area of Figure 2 (a), we sort  $\{G_i\}$  and  $\{P_j\}$  according to their start timestamps, and match  $G_k$  with  $P_k$ , where  $k$  is the new index after sorting and  $1 \leq k \leq \min(|\{G_i\}|, |\{P_j\}|)$ . For the rest of ground truths or predictions, we match them with empty segments  $\phi$ . We also explore other matching strategies in Section 6.1. After matching, we leverage the normalized version of GIoU (Rezatofighi et al., 2019), denoted as NGIoU, as the metric to assess the overlap between paired ground truth  $G_n$  and prediction  $P_n$ , where  $1 \leq n \leq N$ . The value of NGIoU ranges from 0 to 1 and is defined as follows:

$$\text{NGIoU}(G_n, P_n) = \frac{1}{2} \left( 1 + \frac{|G_n \cap P_n|}{|G_n \cup P_n|} - \frac{|\mathcal{C} \setminus (G_n \cup P_n)|}{|\mathcal{C}|} \right), \quad (2)$$

where  $\mathcal{C}$  is the shortest video segment covering  $G_n$  and  $P_n$ . We use GIoU instead of IoU because it better guides model optimization when there is no overlap between the predicted video segment

and the ground truth. Specifically, to encourage the model outputs to align more closely with the ground truth, we impose a penalty when the number of predicted segments differs from that of the ground truth. This is implemented by setting the NGIoU score to 0 when paired with  $\phi$ , i.e.,  $\text{NGIoU}(\cdot, \phi) = 0$ , and  $\text{NGIoU}(\phi, \cdot) = 0$ . Finally, the average NGIoU of all pairs is calculated:

$$r_L = \frac{\sum_{n=1}^N \text{NGIoU}(G_n, P_n)}{N} \quad (3)$$

And the final segment matching reward is defined as

$$r_M = \frac{r_G + r_L}{2} \quad (4)$$

#### 4.2.2 TIMESTAMP REWARD

Explicitly incorporating temporal information into reasoning process during video comprehension helps models better understand complex temporal structures and events in videos, while neglecting temporal information may lead to misconceptions about video content (see the example in Figure 1 (b)). Previous works (Feng et al., 2025; Yu et al., 2025) reveal the importance of model ability of temporal-aware reasoning. Unfortunately, it remains a challenging problem to stimulate this ability.

To tackle this problem, we design the timestamp reward  $r_T$  to enforce models to include timestamps which occur in the final answers in their reasoning processes. Suppose  $\{T_A^i\}$  and  $\{T_R^i\}$  are timestamps occurring in the answer and reasoning process of a model output, then

$$r_T = I_{\{T_R^i\} \subset \{T_A^i\}} \quad (5)$$

where  $I$  is indicator function. As shown in lower part of Figure 2 (a), when all the timestamps occurring in the answer are found in thinking process, models get the reward. If some timestamps fail to match, the reward is set zero. With the timestamp reward, we encourage models to focus on temporal details during reasoning instead of thinking purely based on overall video contents.

#### 4.2.3 FORMAT REWARD

Our format reward follows DeepSeek-R1 (Guo et al., 2025), enforcing models to output their thinking processes and final answers in format “<think>...</think><answer>...</answer>”:

$$r_F = \begin{cases} 1, & \text{if } o_i \text{ has right format} \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

### 4.3 TRAINING RECIPE WITH PHASED REWARDS

Our preliminary experiments indicate that explicitly encouraging models to output timestamps assists models in establishing a timestamp-aware reasoning strategy in early stages, but later leads to performance drop once the model becomes more capable. A possible reason is that enforcing such behavior restricts models to explore and develop more flexible reasoning strategies. To address this, we propose a training recipe with phased rewards, as illustrated in Figure 2(b).

In the early stage, all three rewards, segment matching reward  $r_F$ , timestamp reward  $r_T$ , and format reward  $r_M$ , are combined to form the final reward:

$$r_1 = [\beta r_T + (1 - \beta) r_F] + \alpha r_M. \quad (7)$$

In the later stage, we remove the timestamp reward to allow for more flexible reasoning patterns, yielding the following final reward:

$$r_2 = r_F + \alpha r_M. \quad (8)$$

Training with phased rewards leads to greater performance improvements than using either  $r_1$  or  $r_2$  alone throughout the entire process. Further analysis is provided in Section 6.2. Search of hyperparameters  $\alpha$  and  $\beta$  is introduced in Appendix C.

Table 2: Results of MLLMs on in-domain and out-of-domain tasks. \*Results are copied from original paper. Detailed model versions and introduction of other baselines can be found in Appendix D.

Model	In-Domain				Out-of-Domain									
	Charades-STA (Single-Seg)	THUMOS14 (Multi-Seg)	THUMOS15 (Multi-Seg)	Perception Test (Multi-Seg)	E.T. Bench					E.T. Bench (Subset)				
					REF	GND	CAP	COM	AVG	REF	GND	CAP	COM	AVG
<b>API-based Models</b>														
GPT-4o	25.1	5.5	6.7	-	-	-	-	-	-	37.4	16.5	11.6	6.8	18.1
<b>Open-source ~ 7B Models</b>														
Qwen2.5-VL-7B	50.2	24.9	23.4	25.3	53.1	30.7	16.2	11.3	27.8	<u>51.0</u>	30.3	16.5	9.3	26.8
+vanilla SFT	28.1	15.5	15.6	20.3	24.3	11.3	15.3	6.6	14.4	27.8	12.6	15.0	8.7	16.0
+vanilla GRPO	53.9	<u>25.8</u>	<u>25.6</u>	<u>30.0</u>	54.4	<u>37.6</u>	<u>23.5</u>	20.6	<u>34.0</u>	50.9	<u>36.6</u>	<u>23.7</u>	<u>17.8</u>	<u>32.3</u>
E.T. Chat	45.6	23.7	24.9	9.2	38.4*	<b>38.0*</b>	16.7*	13.5*	26.7	31.8*	33.8*	17.1*	11.1*	23.5
TRACE-7B	29.9*	7.6	7.8	14.0	33.6*	33.8*	20.3*	<b>25.8*</b>	28.4	25.2	17.2	14.7	5.3	15.6
Video-R1	11.3	3.5	3.4	5.7	50.3	25.3	15.6	12.4	25.9	49.2	22.2	15.6	12.8	25.0
VideoChat-R1	<u>59.4</u>	14.3	13.4	27.1	55.8	35.6	22.1	19.5	33.3	47.0	35.9	<u>24.1</u>	12.5	29.9
TimeZero	59.2	14.4	12.7	26.8	<u>55.9</u>	35.8	21.4	17.1	32.6	46.9	35.1	22.9	15.2	30.0
MUSEG-7B (Ours)	<b>59.7</b>	<b>29.7</b>	<b>29.3</b>	<b>31.7</b>	<b>61.9</b>	37.5	<b>23.7</b>	<u>24.0</u>	<b>36.8</b>	<b>60.8</b>	<b>38.8</b>	<b>25.1</b>	<b>19.0</b>	<b>35.9</b>
<b>Open-source ~ 3B Models</b>														
Qwen2.5-VL-3B	41.4	<u>12.6</u>	<u>12.8</u>	19.4	<u>51.7</u>	20.4	13.6	8.0	23.4	52.9	20.4	12.7	<u>7.6</u>	23.4
TEMPURA	<u>44.5</u>	8.7	12.1	<u>20.7</u>	46.3	<u>26.1</u>	<u>14.4</u>	<b>10.2</b>	<u>24.3</u>	<b>56.4</b>	<u>22.8</u>	<u>13.3</u>	3.5	<u>24.0</u>
MUSEG-3B (Ours)	<b>53.7</b>	<b>21.0</b>	<b>20.3</b>	<b>29.1</b>	<b>53.9</b>	<b>30.0</b>	<b>18.7</b>	<u>8.8</u>	<b>27.9</b>	<u>54.3</u>	<b>28.7</b>	<b>18.3</b>	<b>11.8</b>	<b>28.3</b>

## 5 EXPERIMENTS

### 5.1 IMPLEMENTATIONS

Our training dataset is constructed from E.T. Instruct 164k (Liu et al., 2024a) and Charades-STA (Gao et al., 2017). For E.T. Instruct 164k, we only sample data from temporal video grounding (TVG) and temporal action localization (TAL) tasks. Our final training dataset consists of 12.6k samples, including 6,967 with a single segment and 5,633 with multiple segments as ground truths.

We train MUSEG-7B and MUSEG-3B based on Qwen2.5-VL-7B-Instruct and Qwen2.5-VL-3B-Instruct (Bai et al., 2025) (abbreviated as Qwen2.5-VL-7B and Qwen2.5-VL-3B in Table 2), respectively. They are trained with timestamp reward ( $r_1$ ) for 400 steps and without timestamp reward ( $r_2$ ) for another 500 steps. For comparison, we also conduct SFT and naive RL experiments on Qwen2.5-VL-7B-Instruct with our constructed dataset as baselines. For naive RL experiment, we remove timestamp reward, and retain only global matching part of segment matching reward as well as format reward. All other settings remain consistent with those used for training MUSEG-7B. Training details can be found in Appendix B.

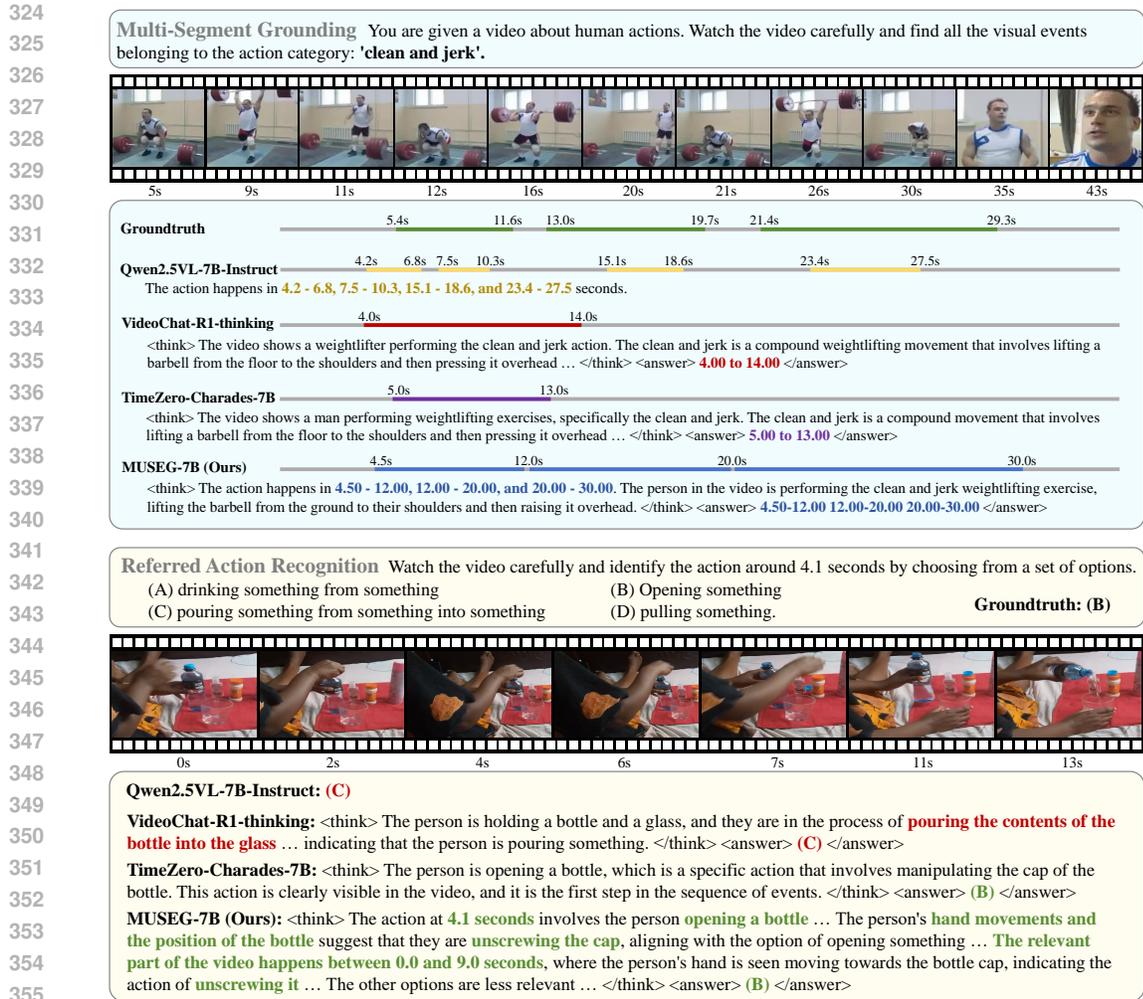
### 5.2 BENCHMARKS AND EVALUATION METRICS

We evaluate MUSEG-7B and MUSEG-3B on both temporal grounding tasks (in-domain) and broader time-related tasks (out-of-domain). For in-domain evaluation, we use the Charades-STA (Gao et al., 2017) test set for single-segment grounding and report performance using mIoU. For multi-segment grounding, we adopt the validation sets of THUMOS14, THUMOS15 (Idrees et al., 2017), and Perception Test (Patraucean et al., 2023), and measure F1 scores averaged over four IoU thresholds (0.1, 0.3, 0.5, and 0.7), following previous work (Liu et al., 2024a).

For out-of-domain evaluation, we assess model generalization on a variety of time-related tasks in E.T. Bench (Liu et al., 2024a), including referring (REF), grounding (GND), dense captioning (CAP), and complex understanding (COM). We adopt the metrics from the original benchmark: accuracy for referring, F1 score for grounding, sentence similarity for dense captioning, and recall for complex understanding.

### 5.3 MAIN RESULTS

As shown in Table 2, MUSEG-7B and MUSEG-3B outperform other SFT- or RL-based methods on most in-domain and out-of-domain tasks among all ~ 7B and ~ 3B models, and even surpass GPT-4o. Our method shows a significant advantage over base models. MUSEG-7B achieves more than 10% performance enhancement on all the tasks compared to its base model Qwen2.5-VL-7B-Instruct. Also, it is worth noting that our model gets doubled performance on complex understanding task, showing strong ability of generalization.



356 Figure 3: Cases of our MUSEG-7B and baselines on multi-segment grounding (in-domain) and  
 357 referred action recognition (out-of-domain) tasks.

358 Furthermore, we compare MUSEG-7B with vanilla GRPO baseline (“+vanilla GRPO” in Table 2),  
 359 which is trained on the same data as MUSEG-7B, and with VideoChat-R1 (“VideoChat-R1” in Ta-  
 360 ble 2), which is trained exclusively on single-segment grounding using vanilla GRPO with Charades-  
 361 STA as the training dataset. As depicted in Table 2, VideoChat-R1 experiences a performance de-  
 362 cline on multi-segment grounding tasks, whereas the vanilla GRPO baseline with multi-segment  
 363 grounding task exhibits limited performance enhancement. This suggests that simply introducing  
 364 multi-segment grounding during training is insufficient to yield substantial gains. In contrast, our  
 365 MUSEG-7B delivers significantly larger gains across both grounding tasks and a broader set of time-  
 366 sensitive tasks, while maintaining comparable general video QA performance to its base model (see  
 367 Appendix E). These results further demonstrate that the effectiveness of MUSEG-7B stems from its  
 368 innovative design of training tasks and reward recipes.

369 Qualitative case studies presented in Figure 3 provide additional evidence of the effectiveness of our  
 370 proposed model. The first case is a multi-segment grounding task (in-domain) with the query “clean  
 371 and jerk”. VideoChat-R1-thinking and TimeZero-Charades-7B only recognize the video segment  
 372 corresponding to the first attempt, consistent with the fact that they are trained only with single-  
 373 segment grounding tasks. In contrast, MUSEG-7B accurately localizes all three weight-lifting at-  
 374 tempts. The performance gap highlights effectiveness of multi-segment grounding training tasks.

375 The second case involves referred action recognition (out-of-domain) query about event happening  
 376 around 4.1 seconds. Seen from the video, the person first opens the bottle, and then pour water out  
 377 from it. VideoChat-R1 incorrectly aligns the event of pouring water from the bottle (occurring at  
 11 seconds) with a 4.1-second timestamp, demonstrating a temporal misalignment in its reasoning.

Table 3: Results with different matching strategies. For all the experiments, we train Qwen2.5-VL-7B for 900 steps same as the training process of MUSEG-7B.

Local Matching Strategy	Charades-STA	THUMOS14	THUMOS15	E.T. Bench (Subset)				
				REF	GND	CAP	COM	AVG
w/o Local Matching	59.3	26.0	25.7	46.2	37.0	22.5	15.3	30.3
w/ Local Matching (Maximum)	58.2	28.6	28.2	56.2	31.5	24.7	17.4	32.5
w/ Local Matching (Sequential)	<b>59.7</b>	<b>29.7</b>	<b>29.3</b>	<b>60.8</b>	<b>38.8</b>	<b>25.1</b>	<b>19.0</b>	<b>35.9</b>

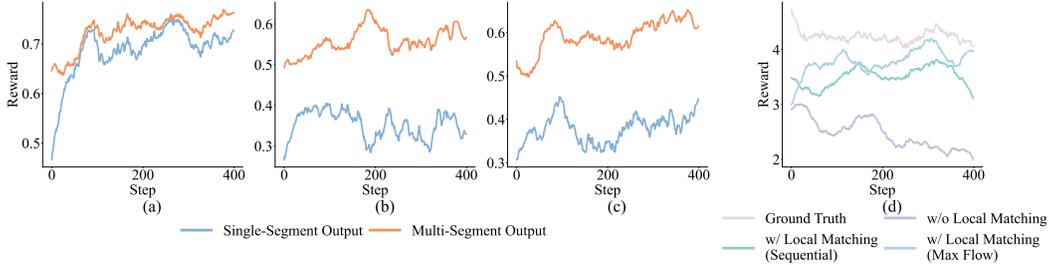


Figure 4: Segment matching reward (a) w/o local matching, (b) w/ local matching (sequential), and (c) w/ local matching (maximum). (d) Evolution of numbers of predicted segments during training process. For all the plots, we only consider queries whose ground truths are more than one segments.

TimeZero-Charades-7B provides the correct answer but lacks precise timestamp references in its explanation. In contrast, MUSEG-7B exhibits superior temporal reasoning capability: it not only identifies the bottle-opening action around 4.1 seconds but also accurately localizes the corresponding video segment.

## 6 ANALYSES

### 6.1 LOCAL MATCHING STRATEGIES

We delve deeper to verify effectiveness of local matching in segment matching reward. We conduct experiments of removing local matching, and only keeping global matching (“w/o Local Matching” in Table 3). Additionally, we explore another design, which involves matching ground truths and predictions to maximize average overlap (“w/ Local Matching (Maximum)” in Table 3). We do this by calculating maximum weighted matching in bipartite graph. For ground truth segments  $\{G_i\}$  and predicted segments  $\{P_j\}$ , we construct a complete bipartite graph  $\mathcal{G}$  as

$$\mathcal{G} = \{(G_i, P_j, W_{i,j})\}, \tag{9}$$

where  $(G_i, P_j, W_j)$  denotes an edge connecting  $G_i$  and  $P_j$  with weight  $W_{i,j} = \text{NGIoU}(G_i, P_j)$ , then we calculate  $r_L$  as follows:

$$r_L = \frac{\text{Matching}(\mathcal{G})}{\max(|\{G_i\}|, |\{P_j\}|)} \tag{10}$$

where  $\text{Matching}(\cdot)$  is the maximum weighted matching function. Table 3 shows that including local matching boost overall model performance compared to only keeping global matching. Additionally, sequential matching reaches better performance than maximum matching. Therefore, we finally adopt sequential matching in MUSEG.

We also notice that drops of model performance on multi-segment grounding are much larger than single-segment grounding when local matching is removed. To better understand its reason, we examine differences in rewards model would get when it produces a single segment or at least two segments for a query whose ground truth consists of more than one segments. As shown in Figure 4 (a), (b), and (c), local matching strategies impose significant penalties on segment matching rewards when model output only contains a single segment, but the penalties imposed by global matching are relatively weak. We further report evolution of numbers of predicted segments during training process in Figure 4 (d). When we remove local matching, numbers of predicted segments significantly drop and their gaps from ground truths become larger. This indicates that local matching can help better align numbers of predicted segments with ground truths.

Table 4: Results with different training recipes.

Training Paradigms	Charades-STA	THUMOS14	THUMOS15	E.T. Bench (Subset)				
				REF	GND	CAP	COM	AVG
w/o Timestamp Reward	56.9	28.4	28.3	55.1	37.6	22.3	13.2	32.1
w/ Timestamp Reward	57.3	26.1	24.6	57.3	28.9	22.0	16.1	31.1
w/ Timestamp Reward for 400 Steps	<b>59.7</b>	<b>29.7</b>	<b>29.3</b>	<b>60.8</b>	<b>38.8</b>	<b>25.1</b>	<b>19.0</b>	<b>35.9</b>

## 6.2 DESIGN OF PHASED REWARDS

In this section, we explore the effectiveness of our proposed training recipe with phased rewards. We compare it against training model with or without timestamp reward during the whole training process in Table 4. From the table we can see that our proposed recipe of training the model with timestamp reward for 400 steps and without timestamp reward for another 500 steps reaches the highest performance. We further change the total training steps and report the results in Figure 5 (a). We can see that our proposed recipe consistently outperforms other training strategies, showing effectiveness over different data scales. We also explore model performance when we vary number of steps of keeping timestamp reward. Figure 5 (b) shows that when the model is trained with timestamp reward for 400 steps, its performance reaches the peak.

To understand the underlying reasons, we examine values of segment matching reward, which reflects accuracy of model output, throughout the training process. As illustrated in Figure 6, when the model is trained either without the timestamp reward or with the timestamp reward applied throughout the entire training process, there is minimal improvement in performance after 400 steps as reasoning forms of models stabilize. In contrast, if the timestamp reward, which is designed to guide the model in referencing specific timestamps during reasoning, is removed after 400 steps, the model can continue to freely explore more effective reasoning strategies, leading to continuous enhancement in its segment matching reward in subsequent steps.

## 7 CONCLUSION

In this work, we introduce MUSEG, a RL-based method to improve video temporal understanding abilities of MLLMs. Experiments demonstrate effectiveness of our method on improving model performance on single-segment and multi-segment grounding tasks, as well as broader time-sensitive scenarios. We hope our proposed method will inspire future research on enhancing temporal understanding abilities of MLLMs.

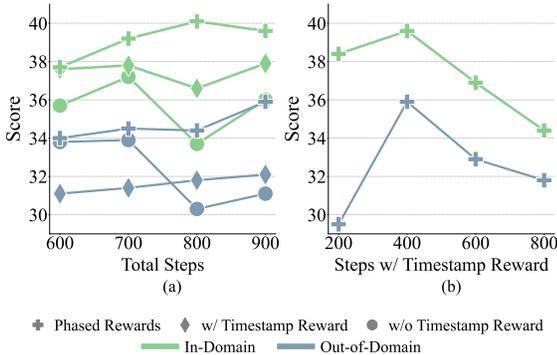


Figure 5: (a) Model performance with different training recipes. For setting of phased rewards, we train models with timestamp reward for 300 steps when total steps are 600 and 700, for 400 steps when total steps are 800 and 900. (b) Model performance when we vary number of steps with timestamp reward, keeping total steps to be 900. For all the experiments, we report average score of Charades-STA, THUMOS14 and THUMOS15 as in-domain score, and average score of E.T. Bench (Subset) as out-of-domain score.

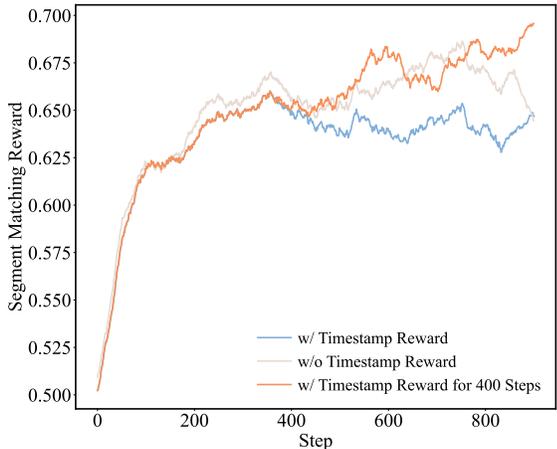


Figure 6: Rewards with different training recipes. We also report timestamp reward during training.

486 REPRODUCIBILITY STATEMENT  
487

488 We train MUSEG-7B and MUSEG-3B using publicly available models and datasets. Implementa-  
489 tion details can be found in Section 5.1 and Appendix B. Additionally, we include our training and  
490 inference code in supplementary material to enhance reproducibility.  
491

492 REFERENCES  
493

- 494 Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid.  
495 ViViT: A Video Vision Transformer. *2021 IEEE/CVF International Conference on Computer  
496 Vision (ICCV)*, pp. 6816–6826, 2021.
- 497 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang,  
498 Shijie Wang, Jun Tang, et al. Qwen2.5-VL Technical Report. *arXiv preprint arXiv:2502.13923*,  
499 2025.
- 500 Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai Zhang, Feng Yao, Fangrui Zhu, Jing Gu,  
501 Yiwu Zhong, Yuzhang Shang, et al. Temporalbench: Benchmarking fine-grained temporal under-  
502 standing for multimodal video models. *arXiv preprint arXiv:2410.10818*, 2024.
- 503 Guo Chen, Yicheng Liu, Yifei Huang, Yuping He, Baoqi Pei, Jilan Xu, Yali Wang, Tong Lu, and  
504 Limin Wang. CG-Bench: Clue-grounded Question Answering Benchmark for Long Video Un-  
505 derstanding. *arXiv preprint arXiv:2412.12075*, 2024.
- 506 Xinlong Chen, Yuanxing Zhang, Yushuo Guan, Bohan Zeng, Yang Shi, Sihan Yang, Pengfei Wan,  
507 Qiang Liu, Liang Wang, and Tieniu Tan. VersaVid-R1: A Versatile Video Understanding and Rea-  
508 soning Model from Question Answering to Captioning Tasks. *arXiv preprint arXiv:2506.09079*,  
509 2025a.
- 510 Yi Chen, Yuying Ge, Rui Wang, Yixiao Ge, Lu Qiu, Ying Shan, and Xihui Liu. Exploring the Ef-  
511 fect of Reinforcement Learning on Video Understanding: Insights from SEED-Bench-R1. *arXiv  
512 preprint arXiv:2503.24376*, 2025b.
- 513 Jen-Hao Cheng, Vivian Wang, Huayu Wang, Huapeng Zhou, Yi-Hao Peng, Hou-I Liu, Hsiang-Wei  
514 Huang, Kuang-Ming Chen, Cheng-Yen Yang, Wenhao Chai, et al. TEMPURA: Temporal Event  
515 Masked Prediction and Understanding for Reasoning in Action. *arXiv preprint arXiv:2505.01583*,  
516 2025a.
- 517 Zixu Cheng, Jian Hu, Ziquan Liu, Chenyang Si, Wei Li, and Shaogang Gong. V-STaR: Bench-  
518 marking Video-LLMs on Video Spatio-Temporal Reasoning. *arXiv preprint arXiv:2503.11495*,  
519 2025b.
- 520 Jisheng Dang, Jingze Wu, Teng Wang, Xuanhui Lin, Nannan Zhu, Hongbo Chen, Wei-Shi Zheng,  
521 Meng Wang, and Tat-Seng Chua. Reinforcing Video Reasoning with Focused Thinking. *arXiv  
522 preprint arXiv:2505.24718*, 2025.
- 523 Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Benyou  
524 Wang, and Xiangyu Yue. Video-R1: Reinforcing Video Reasoning in MLLMs. *arXiv preprint  
525 arXiv:2503.21776*, 2025.
- 526 Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
527 Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-MME: The First-Ever Comprehensive Eval-  
528 uation Benchmark of Multi-modal LLMs in Video Analysis. In *Proceedings of the Computer  
529 Vision and Pattern Recognition Conference*, pp. 24108–24118, 2025.
- 530 Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. TALL: Temporal Activity Localization  
531 via Language Query. In *Proceedings of the IEEE international conference on computer vision*,  
532 pp. 5267–5275, 2017.
- 533 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
534 Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing Reasoning Capability in  
535 LLMs via Reinforcement Learning. *arXiv preprint arXiv:2501.12948*, 2025.

- 540 Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xiangyu Zeng, Cheng Liang, Tao Wu, Xi Chen,  
541 Liang Li, and Limin Wang. Online Video Understanding: A Comprehensive Benchmark and  
542 Memory-Augmented Method. *arXiv preprint arXiv:2501.00584*, 2024.
- 543
- 544 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
545 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. GPT-4o System Card. *arXiv preprint*  
546 *arXiv:2410.21276*, 2024.
- 547 H. Idrees, A. R. Zamir, Y. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah. The THUMOS  
548 challenge on action recognition for videos “in the wild”. *Computer Vision and Image Under-*  
549 *standing*, 155:1–23, 2017.
- 550
- 551 Hongyu Li, Jinyu Chen, Ziyu Wei, Shaofei Huang, Tianrui Hui, Jialin Gao, Xiaoming Wei, and  
552 Si Liu. LLaVA-ST: A Multimodal Large Language Model for Fine-Grained Spatial-Temporal  
553 Understanding. *arXiv preprint arXiv:2501.08282*, 2025a.
- 554 Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,  
555 Ping Luo, et al. MVBench: A Comprehensive Multi-modal Video Understanding Benchmark.  
556 In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 22195–22206,  
557 2024.
- 558
- 559 Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu Zeng, Yanan He, Yali Wang, Yu Qiao,  
560 Yi Wang, and Limin Wang. VideoChat-R1: Enhancing Spatio-Temporal Perception via Rein-  
561 forcement Fine-Tuning. *arXiv preprint arXiv:2504.06958*, 2025b.
- 562 Yun Li, Zhe Liu, Yajing Kong, Guangrui Li, Jiyuan Zhang, Chao Bian, Feng Liu, Lina Yao, and  
563 Zhenbang Sun. Exploring the Role of Explicit Temporal Modeling in Multimodal Large Language  
564 Models for Video Understanding. *arXiv preprint arXiv:2501.16786*, 2025c.
- 565
- 566 Ye Liu, Zongyang Ma, Zhongang Qi, Yang Wu, Ying Shan, and Chang W Chen. E.T. Bench:  
567 Towards Open-Ended Event-Level Video-Language Understanding. *Advances in Neural Infor-*  
568 *mation Processing Systems*, 37:32076–32110, 2024a.
- 569 Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and  
570 Lu Hou. TempCompass: Do Video LLMs Really Understand Videos? *Findings of the Association*  
571 *for Computational Linguistics: ACL 2024*, pp. 8731–8772, 2024b.
- 572
- 573 Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video Swin Trans-  
574 former. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
575 *(CVPR)*, pp. 3192–3201, 2021.
- 576 Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. CLIP4Clip: An  
577 Empirical Study of CLIP for End to End Video Clip Retrieval. *Neurocomputing*, 508:293–304,  
578 2021.
- 579
- 580 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
581 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
582 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
583 27730–27744, 2022.
- 584 Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens, Larisa Markeeva, Dylan Ba-  
585 narse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, et al. Perception Test: A  
586 Diagnostic Benchmark for Multimodal Video Models. *Advances in Neural Information Process-*  
587 *ing Systems*, 36:42748–42761, 2023.
- 588
- 589 Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese.  
590 Generalized Intersection Over Union: A Metric and a Loss for Bounding Box Regression. In  
591 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
592 June 2019.
- 593 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy  
Optimization Algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- 594 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,  
595 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. DeepSeekMath: Pushing the Limits of  
596 Mathematical Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300*, 2024.  
597
- 598 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,  
599 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: A Family of  
600 Highly Capable Multimodal Models. *arXiv preprint arXiv:2312.11805*, 2023.
- 601 Jiawei Wang, Liping Yuan, Yuchen Zhang, and Haomiao Sun. Tarsier: Recipes for Training and  
602 Evaluating Large Video Description Models. *arXiv preprint arXiv:2407.00634*, 2024.  
603
- 604 Ning Wang, Guangming Zhu, Liang Zhang, Peiyi Shen, Hongsheng Li, and Cong Hua. Spatio-  
605 temporal interaction graph parsing networks for human-object interaction recognition. In *Pro-  
606 ceedings of the 29th ACM international conference on multimedia*, pp. 4985–4993, 2021.
- 607 Qi Wang, Yanrui Yu, Ye Yuan, Rui Mao, and Tianfei Zhou. VideoRFT: Incentivizing Video Reason-  
608 ing Capability in MLLMs via Reinforced Fine-Tuning. *arXiv preprint arXiv:2505.12434*, 2025a.  
609
- 610 Ye Wang, Boshen Xu, Zihao Yue, Zihan Xiao, Ziheng Wang, Liang Zhang, Dingyi Yang, Wenxuan  
611 Wang, and Qin Jin. TimeZero: Temporal Video Grounding with Reasoning-Guided LVLM. *arXiv  
612 preprint arXiv:2503.13377*, 2025b.
- 613 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
614 Zhou, et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *Advances  
615 in neural information processing systems*, 35:24824–24837, 2022.
- 616 Junbin Xiao, Angela Yao, Yicong Li, and Tat-Seng Chua. Can I Trust Your Answer? Visually  
617 Grounded Video Question Answering. *Proceedings of the IEEE/CVF Conference on Computer  
618 Vision and Pattern Recognition (CVPR)*, pp. 13204–13214, 2024.  
619
- 620 Hu Xu, Gargi Ghosh, Po-Yao (Bernie) Huang, Dmytro Okhonko, Armen Aghajanyan, and Flo-  
621 rian Metzger Luke Zettlemoyer Christoph Feichtenhofer. VideoCLIP: Contrastive Pre-training for  
622 Zero-shot Video-Text Understanding. In *Conference on Empirical Methods in Natural Language  
623 Processing*, 2021.
- 624 Yibin Yan, Jilan Xu, Shangzhe Di, Yikun Liu, Yudi Shi, Qirui Chen, Zeqian Li, Yifei Huang, and  
625 Weidi Xie. Learning Streaming Video Representation via Multitask Training. *arXiv preprint  
626 arXiv:2504.20041*, 2025.  
627
- 628 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
629 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 Technical Report. *arXiv preprint  
630 arXiv:2505.09388*, 2025.
- 631 En Yu, Kangheng Lin, Liang Zhao, Yana Wei, Zining Zhu, Haoran Wei, Jianjian Sun, Zheng Ge, Xi-  
632 angyu Zhang, Jingyu Wang, et al. Unhackable Temporal Rewarding for Scalable Video MLLMs.  
633 *arXiv preprint arXiv:2502.12081*, 2025.
- 634 Haoji Zhang, Xin Gu, Jiawen Li, Chixiang Ma, Sule Bai, Chubin Zhang, Bowen Zhang, Zhichao  
635 Zhou, Dongliang He, and Yansong Tang. Thinking With Videos: Multimodal Tool-Augmented  
636 Reinforcement Learning for Long Video Reasoning. *arXiv preprint arXiv:2508.04416*, 2025a.  
637
- 638 Xingjian Zhang, Siwei Wen, Wenjun Wu, and Lei Huang. TinyLLaVA-Video-R1: Towards Smaller  
639 LMMs for Video Reasoning. *arXiv preprint arXiv:2504.09641*, 2025b.  
640  
641  
642  
643  
644  
645  
646  
647

## 648 A LLM USAGE STATEMENT

649  
650 Throughout the completion of this work, LLMs are solely used for the purpose of spelling checking  
651 and grammatical error detection for the manuscript. They are not employed for any other purposes,  
652 such as deriving idea of this paper, or validating the methods or results.  
653

## 654 B IMPLEMENTATION DETAILS

655  
656 We leverage 7B and 3B models of Qwen2.5-VL (Bai et al., 2025) series as our base models. They  
657 are trained on large scale image and video data and demonstrate strong instruction following and  
658 reasoning abilities. Additionally, there are special designs in Qwen2.5-VL to enable models to  
659 process absolute timestamps and dynamic resolutions of video frames. During training and inference  
660 of MUSEG-7B and MUSEG-3B, we set maximum total video tokens to be 3584 and maximum  
661 number of frames to be 448.

662 We train MUSEG-7B and MUSEG-3B for 900 steps in total, including 400 steps with timestamp  
663 reward and another 500 steps without timestamp reward. We set `batch_size = 14` and `learning_rate =`  
664 `1e-5`. We set  $\alpha = 2$  in phase 1 and phase 2 reward, and  $\beta = 0.4$  in phase 1 reward (see experiments  
665 of hyperparameters search in Appendix C). Considering that base models have been trained on  
666 temporal-related data and already have strong abilities of instruction-following, we do not include  
667 SFT stage in our experiments as DeepSeek-R1 (Guo et al., 2025).

668 It takes about 22 hours for MUSEG-7B phase 1 training, 27 hours for MUSEG-7B phase 2 training,  
669 16 hours for MUSEG-3B phase 1 training and 20 hours for MUSEG-3B phase 2 training on 8 A100-  
670 80G GPUs.  
671

## 672 C HYPERPARAMETERS SEARCH

673  
674 For  $\alpha$  and  $\beta$ , we conduct experiments using a small dataset comprising 1400 samples for training,  
675 along with a randomly selected subset of 200 samples from THUMOS14 and THUMOS15 for  
676 evaluation, aiming at identifying the optimal combination of hyperparameters. Based on results  
677 presented in Table 5 and Table 6, we select the best combination of  $\alpha = 2$  and  $\beta = 0.4$ .  
678

679 Table 5: Model performance with different settings of  $\alpha$ .

Hyperparameters	THUMOS (subset)
$\alpha = 1, \beta = 0.4$	23.4
$\alpha = 1.5, \beta = 0.4$	27.6
$\alpha = 2, \beta = 0.4$	<b>28.8</b>
$\alpha = 3, \beta = 0.4$	19.3
$\alpha = 4, \beta = 0.4$	21.5

687 Table 6: Model performance with different settings of  $\beta$ .

Hyperparameters	THUMOS (subset)
$\alpha = 2, \beta = 0.2$	28.4
$\alpha = 2, \beta = 0.4$	<b>28.8</b>
$\alpha = 2, \beta = 0.6$	28.7

## 694 D INTRODUCTION OF BASELINES

695  
696 Our baselines can be categorized into SFT-based methods and RL-based methods. We introduce  
697 SFT-based models first:  
698

699  
700 **E.T. Chat (~7B)**: It compresses video frames into single tokens using a Q-Former-based compres-  
701 sor with cross-attention, and generates timestamps with special tokens. It is trained on E.T. Instruct  
164k, a dataset covering 9 tasks across 14 sources.

**TRACE ( ~ 7B):** It is trained with a causal event modeling framework, integrating timestamp, salient score, and textual caption prediction tasks. Its training data include 1.9M samples from Valley, TextVR, ShareGPT4Video, and 0.9M samples from ActivityNet Captions and InternVid.

**TEMPURA ( ~ 3B):** It is trained with masked event prediction reasoning, event segmentation and dense captioning tasks. Its training data consist of 500k samples.

Then we introduce RL-based models:

**Video-R1 ( ~ 7B):** It is trained by SFT with 165k samples and RL with 260k samples. Its training data consist of various general image question answering and video question answering tasks.

**VideoChat-R1-thinking ( ~ 7B, abbreviated as VideoChat-R1 in Table 2):** It is trained with temporal grounding, object tracking, video captioning and grounded video question answering tasks, with a total data scale of 18.0k samples.

**TimeZero-Charades-7B ( ~ 7B, abbreviated as TimeZero in Table 2):** It is trained towards temporal grounding tasks. A version of its models is trained with Charades-STA (Gao et al., 2017).

Additionally, we report performance of GPT-4o-2024-11-20 (Hurst et al., 2024) (abbreviated as GPT-4o in Table 2) for reference. In consideration of inference costs, we do not report results of GPT-4o on Perception Test and the whole set of E.T. Bench. Only results on a subset of 470 samples of E.T. Bench, specified by the original paper, are reported.

## E MODEL PERFORMANCE ON GENERAL VIDEO QA TASKS

We present model performance on general video QA benchmarks, MVBench (Li et al., 2024) and Video-MME (Fu et al., 2025), in Table 7. Results indicate that our overall performance on MVBench and Video-MME is comparable to that of our base model. Our approach does not negatively impact model performance on general video QA tasks.

Table 7: Model performance on MVBench and Video-MME.

Model	MVBench	Video-MME			
		SHORT	MEDIUM	LONG	AVG
Qwen2.5-VL-7B-Instruct	65.7	71.8	62.7	52.6	62.4
MUSEG-7B	<b>67.4</b>	<b>71.9</b>	<b>63.7</b>	<b>53.6</b>	<b>63.1</b>

## F TRAINING PROMPTS

We prompt models to include timestamps in their reasoning processes and ensure that the timestamps are consistent with those in answers using instruction in Table 8.

Table 8: Training prompts.

Training prompts
{QUESTION} First, output reasoning process in <think> </think> tags. The reasoning process must REFER TO SPECIFIC TIMESTAMPS TO TELL WHERE YOU GET THE INFORMATION FROM THE VIDEO. Then summarize your reasoning process above and output selected segments like <answer>X.XX-X.XX</answer>, where X denotes arabic numbers. If there are multiple segments, separate them with spaces like <answer>X.XX-X.XX X.XX-X.XX</answer>. Your output format should be like <think>...</think><answer>...</answer>.