Figure 6. **Up:** Spatial Temporal CoT SFT + GRPO; **Down:** Regular SFT + GRPO

# A. Appendix

## A.1. Future Work

Several promising directions remain for advancing multimodal spatial-temporal reasoning:

1. **Expanding Modalities and Data Diversity:** Incorporate additional modalities and more diverse datasets to improve models' spatial-temporal inference abilities.
2. **Implicit Reconstruction Objectives:** Encourage models to learn deeper representations of spatial and temporal dynamics through implicit reconstruction tasks.
3. **Generalized Video Reasoning Models:** Develop efficient and robust video reasoning models with strong generalization in spatial-temporal reasoning.
4. **Efficient Inference and Scalability:** Advance computationally efficient inference methods to enable large-scale, practical deployment.

Addressing these challenges will be critical to furthering the state-of-the-art in multimodal spatial-temporal reasoning.

## A.2. Inference Details

In Ego-bench, all closed-source and open-source models sample 8 frames for inference. The closed source models have a single frame size of $224 \times 224$, while the open source models follow the default optimal size for each model. Note: Due to large fluctuations in the performance of closed-source models over time, we have deliberately indicated that all closed-source models here were tested on 7 March 2025.
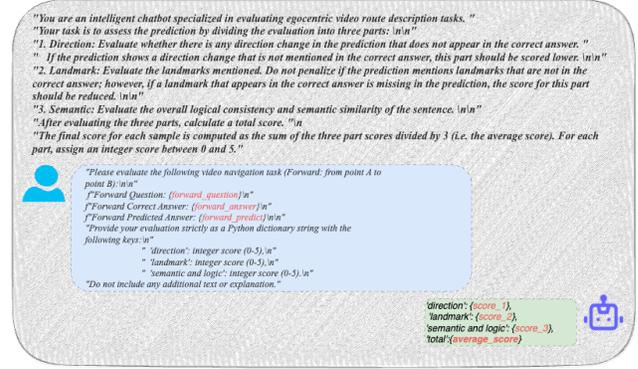


Figure 7. **The Prompt for the evaluation.** this prompt is scored in three parts, *Direction:* focuses on MLLM's ability to perceive spatial and temporal changes. *Landmarks:* focuses on MLLM's ability to perceive spatial changes and self-orientation. *Logical Semantics:* focuses on MLLM's ability to organize logical language.

## A.3. ST-R1 Training Details

As shown in Table 5 and Table 6, we list the key training parameters in the Lora Supervised finetuning and GRPO, and all experiments were run on $4 \times$ H100 GPUs.

| Frames | 8 |
| --- | --- |
| max prompt length | 4096 |
| base learning rate | 1e-5 |
| gradient accumulation steps | 1 |
| $\beta$ | 0.1 |
| Group | 8 |
| base batch size(per bs) | 1 |

Table 5. GRPO Training Detail.

| Frames | 8 |
| --- | --- |
| Quantization | 4-bit |
| Lora rank | 8 |
| Lora alpha | 8 |
| Lora dropout | 0.1 |
| optimizer | AdamW |
| learning rate schedule | cosine decay |
| base learning rate | 8e-5 |
| base batch size(per bs) | 1 |
| total epochs | 1 |
| warmup ratio | 0.1 |
| max gradient norm | 0.3 |

Table 6. Supervised finetuning Training Detail.

As shown in Figure 8, the loss curve and acc curve of CoT sft are illustrated in the figure. As shown in Figure 9,
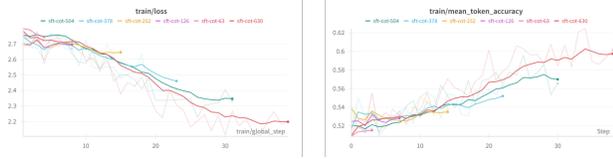
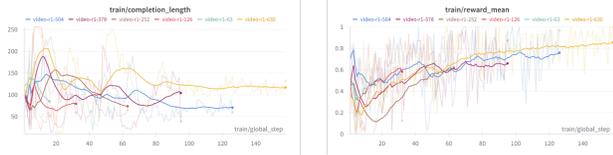Figure 8. **Left:** Loss curve. **Right:** Mean Token Accuracy



Figure 9. **Left:** Completion Length. **Right:** Reward Mean Accuracy

the graph illustrates some of the metrics for the final reinforcement learning phase of R1.