

LongVideoAgent: Multi-Agent Reasoning with Long Videos

Runtao Liu* Ziyi Liu* Jiaqi Tang Yue Ma Renjie Pi Jipeng Zhang Qifeng Chen†

Hong Kong University of Science and Technology

rliuay@connect.ust.hk, ziyiliu0811@outlook.com, cqf@ust.hk

<https://longvideoagent.github.io/>

Abstract

Recent advances in multimodal LLMs and systems that *use tools* for long-video QA point to the promise of reasoning over hour-long episodes. However, many methods still compress content into lossy summaries or rely on limited toolsets, weakening temporal grounding and missing fine-grained cues. We propose a multi-agent framework in which a master LLM coordinates a grounding agent to localize question-relevant segments and a vision agent to extract targeted textual observations. The master agent plans with a step limit, and is trained with reinforcement learning to encourage concise, correct, and efficient multi-agent cooperation. This design helps the master agent focus on relevant clips via grounding, complements subtitles with visual detail, and yields interpretable trajectories. On our proposed *LongTVQA* and *LongTVQA+* which are episode-level datasets aggregated from TVQA/TVQA+, our multi-agent system significantly outperforms strong non-agent baselines. Experiments also show reinforcement learning further strengthens reasoning and planning for the trained agent.

1 Introduction

Multimodal large language models (MLLMs) extend LLMs beyond text to perceive and reason over multimodal signals, such as visual frames, audio, and subtitles. A key emerging challenge is robust *long video* understanding, where information is sparsely distributed across hours of content and multiple modalities (e.g., frames, and dialogue cues). Early instruction-tuned systems such as Video-LLaMA (Zhang et al., 2023; Lin et al., 2024) demonstrated that LLMs can be adapted to jointly process sampled video frames, marking an initial step toward multimodal video reasoning. However, current models remain limited to

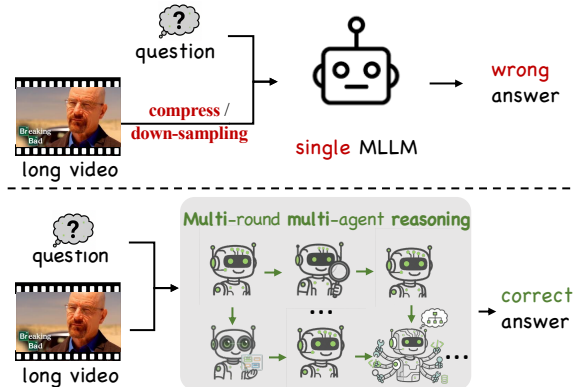


Figure 1: Traditional single-pass MLLMs ingest entire long videos in a single context, typically through heavy downsampling and compression, which often causes them to miss crucial evidence and produce wrong answers. In contrast, LONGVIDEOAGENT conducts *multi-agent*, *multi-round*, and *multimodal* reasoning to extract sparse, task-relevant cues and answer correctly.

short clips or coarse summaries and struggle with fine-grained, temporally extended queries. Crucially, most prior systems are *non-agentic* models: they process a static, pre-encoded or down-sampled video. Converting the full visual stream into compressed representations in the LLM’s textual space shifts the burden of temporal reasoning to this early stage—often lossy and irreversible, making it difficult to recover fine-grained evidence. These limitations motivate an *agentic*, tool-augmented paradigm that can actively decide what to observe next, when to query external visual or other tools, and when enough grounded evidence has been gathered to respond. Despite recent advances, the field still lacks a solution that jointly achieves efficiency, multimodal completeness, and fine-grained temporal reasoning in long videos.

Recent works have begun to frame long video understanding as an *agent-driven process*, rather than a passive encoding task. Notably, VideoAgent (Fan et al., 2024; Wang et al., 2024b) introduced an agent-based framework where a central

* Equal Contribution. † Corresponding author.

LLM actively conducts video analysis. In this paradigm, the LLM agent iteratively queries external vision models (tools) to retrieve and interpret video frames, progressively compiling the information needed to answer a given query. This interactive strategy mirrors human cognitive behavior and has demonstrated promising effectiveness. These findings highlight the potential of *tool-augmented LLM agents* in achieving both efficiency and accuracy. However, the initial incarnation of VideoAgent relies on a less powerful toolset, primarily generic vision-language foundation models for captioning and image retrieval. Such tools are often insufficient for capturing fine-grained semantics, precise object references, or subtle temporal cues. This restricts the agent’s ability to understand complex scenes and reason over long temporal spans. Moreover, current frameworks underutilize the LLM’s inherent reasoning abilities and lack mechanisms for multi-step decision making or reinforcement-based planning.

In this paper, shown as Figure 1 we address these challenges by proposing a new *multi-agent-based framework* for long video understanding that strategically incorporates agents. Our system adopts a multi-agent architecture, where a central MASTERAGENT is responsible for reasoning and answering, while coordinating with other specialized agents. Specifically, a GROUNDINGAGENT locates video segments relevant to the question, and a VISIONAGENT extracts detailed visual information from the selected clips (e.g., objects, faces, actions). The master agent gathers these outputs to iteratively reason over the accumulated evidence. To guide the reasoning process, we design a reward-driven training strategy that encourages the master agent to conduct structured, multi-step reasoning. In each iteration, the master agent generates sub-queries, invokes either the grounding or vision agent as needed, and integrates the returned information before deciding on the next step. When it determines that enough evidence has been collected, it produces a final answer. By designing a reward function that penalizes irrelevant tool use and incoherent reasoning, we guide the agent to “think” in a proper format, effectively learning when to explore the video with tools and when it has gathered sufficient evidence to answer the question. Furthermore, to evaluate long-form video reasoning in a realistic setting, we construct a new benchmark dataset *LongTVQA*

and *LongTVQA+*. This dataset extends the well-known TVQA video question answering task to much longer video durations, providing a rigorous testbed for our agent.

Our *Agent-with-Tools* approach demonstrates superior performance on the LongTVQA benchmark, outperforming all existing baselines by a significant margin. Through ablation studies, we show that both the multi-agent architecture and the reward-guided training contribute critically to the agent’s gains. Our system not only achieves higher accuracy, but also exhibits interpretable decision-making, coordinating sub-agents to select relevant video segments and extract fine-grained visual information essential for reasoning. These results underscore the benefit of an agentic framework for long video understanding.

Our contributions are threefold: (i) a modular *multi-agent* architecture in which a master LLM coordinates grounding and vision specialists; (ii) a *reward-driven* agentic reinforcement learning training scheme that promotes concise, step-wise reasoning; and (iii) episode-level long video datasets LongTVQA and LongTVQA+ are proposed under which our system achieves state-of-the-art results.

2 Related Work

2.1 Video Question Answering

Early work focused on memory and attention mechanisms over appearance–motion features (Gao et al., 2018). This evolved into multimodal transformers designed for efficient frame sampling (Lei et al., 2021). Recent trends emphasize retrieval-aware reasoning and efficient tokenization for long videos, as well as integrating LLM-based reasoning with video encoders (Zhang et al., 2023) and employing agentic planners that iteratively gather evidence (Wang et al., 2024b). Long-form systems further explore sparse memory and temporal grounding techniques to handle hour-scale inputs (Song et al., 2024). These developments motivate long-form VideoQA systems that selectively retrieve segments under a limited context budget.

2.2 LLM Agents

LLM agents couple chain-of-thought with *actions*: planning, tool calls, and iterative evidence gathering. Foundational agent ideas include ReAct, Self-Ask, and WebGPT (Yao et al., 2022; Press et al.,

2022; Nakano et al., 2021). Toolformer shows self-supervised API-calling, while orchestration frameworks (HuggingGPT/Gorilla-style) route subtasks to expert models (Schick et al., 2023; Shen et al., 2023). In multimodal settings, MM-ReAct wires LLMs to vision experts via prompting, and program-of-thought systems like ViperGPT compose perception modules through executable code for transparent, verifiable reasoning (Yang et al., 2023; Surís et al., 2023). For long videos, agentic designs such as VideoAgent/VideoAgent-style frameworks use memory, targeted retrieval, and temporal grounding to operate under strict context budgets while improving faithfulness (Wang et al., 2024b). Beyond planning, video-RAG pipelines extract ASR/OCR/objects and retrieve evidence to augment LVLMs for factual responses (Luo et al., 2024). In addition, long-horizon multimodal agents with persistent memory and structural planning further enhance reliability for extended videos, e.g., Long-Seeing, VideoTree, and Koala (Long et al., 2025; Wang et al., 2025b; Tan et al., 2024); and general reasoning paradigms such as Chain-of-Thought, Least-to-Most, Tree-of-Thoughts, and Generative Agents provide foundations for decomposition and memory (Wei et al., 2022; Zhou et al., 2022; Yao et al., 2023; Park et al., 2023). Retrieval-first paradigms like Retrieving-to-Answer complement agent pipelines with a retrieve-then-reason template (Pan et al., 2023). (We also include the alternate ReAct entry for key consistency (Yao et al., 2022).)

2.3 Multi-Modal LLMs

Modern MLLMs combine strong vision encoders with instruction-tuned LLMs. CLIP pretraining provides broad visual–text transfer (Radford et al., 2021). Flamingo introduces a perceiver-style resampler for few-shot multimodal learning (Alayrac et al., 2022); BLIP-2/InstructBLIP bridge frozen encoders and LLMs (Li et al., 2023; Dai et al., 2023). Recent visually instruction-tuned MLLMs (Tang et al., 2025; Pi et al., 2024, 2025), such as LLaVA (Liu et al., 2023), scale visual instruction tuning using open components, while LLaVA-OneVision (Li et al., 2024a) unifies high-resolution perception with token-efficient processing for both images and videos. Recent video-tuned variants (e.g., VideoLLaVA) and training-free token schedulers (e.g., SlowFast-LLaVA) further improve temporal cov-

erage and efficiency (Lin et al., 2024; Xu et al., 2024b). Proprietary MLLMs (GPT-4/4o; Gemini 1.5) show long-context multimodal reasoning (Achiam et al., 2023; Gemini Team, 2024), while open models (Qwen2-VL, InternVL) narrow the gap via dynamic resolution, OCR, and video pipelines (Wang et al., 2024a; Chen et al., 2024). Complementary advances focus on unifying image–video tokens with few, informative representations (e.g., MiniGPT4-Video, VideoChatGPT, Video-LaVIT, LLaMA-VID, LongVU, PLLaVA, LLaVA-Video, Chat-UniVi) (Ataallah et al., 2024; Maaz et al., 2024; Jin et al., 2024b; Li et al., 2024b; Shen et al., 2024; Xu et al., 2024a; Zhang et al., 2024c; Jin et al., 2024a), and on long-context optimization or adaptive input selection (e.g., InternVideo2.5, LongVLM, Long Context training, self-adaptive sampling, simple-but-effective alignment, and question-instructed tuning) (Wang et al., 2025a; Weng et al., 2024; Zhang et al., 2024b; Han et al., 2023; Zhang et al., 2024a; Romero and Solorio, 2024). Comprehensive analyses of video understanding in large multimodal models (e.g., Apollo) situate these models within broader capabilities and evaluation protocols (Zohar et al., 2025). For key harmonization with the bibliography, we also include the alternate VideoLLaMA entry (Zhang et al., 2023). However, most models still face long-video constraints (context length, retrieval). This motivates combining video-native encoders, instruction tuning, retrieval, and tool use for scalable long-form VideoQA.

3 Method

As shown in Figure 2, we cast long-video QA as *multi-agent reasoning*, where a master agent LLM coordinates a grounding agent to temporally localize question-relevant segments and a vision agent to extract targeted observations from those segments. The system proceeds iteratively, maintaining a running context that accumulates subtitles, relevant segment tags, and vision observations, and it produces an answer once the master agent judges that sufficient evidence has been gathered. For open-source LLMs serving as the master agent, we apply reinforcement learning to encourage accurate, concise, and cooperation-efficient behavior while keeping the other agents frozen. At inference, the process yields clear, step-by-step traces aimed at solving the question at hand.

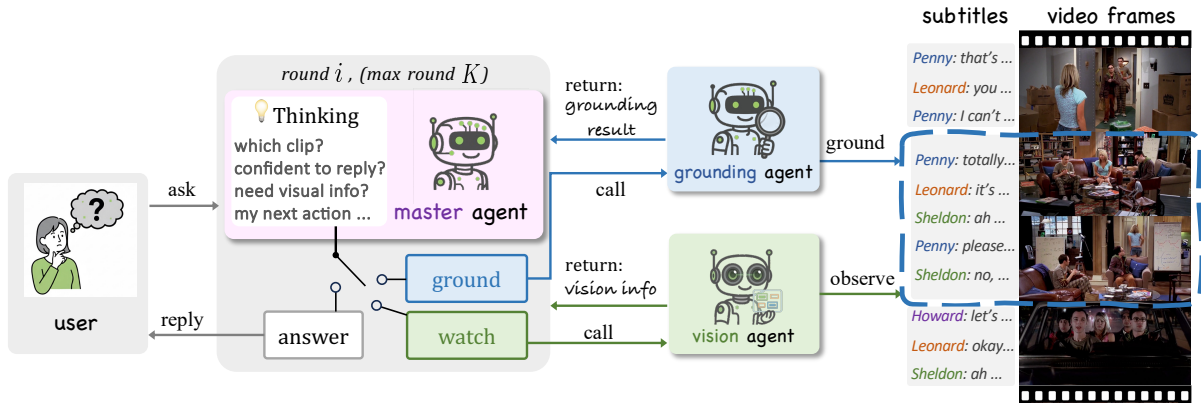


Figure 2: Architecture of LONGVIDEOAGENT. A MASTERAGENT runs for up to K rounds, collaborating with a GROUNDINGAGENT to localize relevant clips from videos and a VISIONAGENT to read fine-grained cues from the localized frames. Evidence accumulates until the MASTERAGENT feels confident to answer the user.

Table 1: System prompt for LONGVIDEOAGENT.

System Prompt — LONGVIDEOAGENT

You are an agent that answers questions about a long video episode. You may use two tools: a *grounding agent* to localize relevant segments and a *vision agent* to extract visual facts from the localized segment. Produce *concise, direct* answers.

Context you may receive. All subtitles and the user question q . When a segment has been localized, you will also have a tag $\langle \text{clip}_X \rangle$ (e.g., $\langle \text{clip}_2 \rangle$). When the vision agent has been called, you will see its textual response.

Available actions (choose exactly one per turn).

A — Visual query: If current visual information is insufficient, or you need visual details conditioned on the subtitles for the current $\langle \text{clip}_X \rangle$, call the vision engine with $\langle \text{visual_query} \rangle \text{query} \langle / \text{visual_query} \rangle$.

B — (Re)Grounding: If the current text/visual evidence conflicts with the question, or the current location cannot support a confident answer, call the grounding agent with $\langle \text{request_grounding} \rangle \text{query} \langle / \text{request_grounding} \rangle$.

C — Answer: If evidence is sufficient, return the final answer with $\langle \text{answer} \rangle \dots \langle / \text{answer} \rangle$. The answer must be concise and direct.

Guidelines. (1) Be conservative with tool calls; answer when sufficient. (2) Do not hallucinate visual details; only use the vision agent for facts not inferable from subtitles. (3) Each turn targets the current $\langle \text{clip}_X \rangle$ (if any); if none exists, prefer (re)grounding before visual query.

3.1 Multi-agent System Framework

Master agent behavior and training. Specifically, the master agent follows the instruction schema in the *System Prompt* (Table 1) and the multi-turn policy in Algorithm 1 that coordinates two other specialist agents: a grounding agent and a vision agent. Given an episode with its full subtitles and a question, the master runs a bounded loop (at most K steps). At each turn it emits exactly one structured action token, $\langle \text{visual_query} \rangle$ for a visual read, $\langle \text{request_grounding} \rangle$ for (re)localization, or $\langle \text{answer} \rangle$ to terminate. After the corresponding agent is invoked, its textual output is appended to the context of the master agent. For open-source masters, we optimize the policy with GRPO while keeping the grounding and vision agents fixed. The rollouts terminated by action tokens in Algorithm 1 provide the trajectories for training and evaluation.

Grounding agent. Given the question and subtitles, the grounding agent proposes a temporal segment and returns a symbolic tag $\langle \text{clip}_X \rangle$ mark-

ing the relevant portion of the episode. By default the window context is 1; when larger, the agent outputs a short run of consecutive tags. The master may re-query grounding to refine or validate the segment as reasoning progresses.

Vision agent. Conditioned on $\langle \text{clip}_X \rangle$ and an on-demand prompt that specifies the current visual need, the vision agent extracts textual observations from frames within the localized segment (e.g., objects/entities, attributes, actions, OCR/on-screen text, scene cues). These observations are appended to the context and guide the next decision; the loop terminates when the master judges the accumulated visual evidence sufficient to answer.

3.2 Reinforcement Learning for LONGVIDEOAGENT

For open-source LLMs serving as the master agent, we fine-tune the master with GRPO while keeping the grounding and vision agents frozen. Long-video QA is cast as a finite-horizon deci-

Algorithm 1 LONGVIDEOAGENT with Multi-Turn Reasoning

Require: Subtitles \mathcal{S} ; question q ; video V ; MASTERAGENT parameters π_θ ; maximum steps K ; GROUNDINGAGENT; VISIONAGENT.

Ensure: Final answer \hat{y} .

```
1: Initialize rollout sequence  $y \leftarrow \emptyset$ 
2: Initialize step count  $t \leftarrow 0$ 
3: while  $t < K$  do
4:   Initialize current action LLM rollout sequence  $y_t \leftarrow \emptyset$ 
5:   while True do
6:     Generate thinking token  $y_i \sim \pi_\theta(\cdot \mid \mathcal{S}, q, V, y + y_t)$ 
7:     Append  $y_i$  to rollout sequence  $y_t \leftarrow y_t + y_i$ 
8:     if  $y_i$  in [ $\langle \text{/visual\_query} \rangle$ ,  $\langle \text{/request\_grounding} \rangle$ ,  $\langle \text{/answer} \rangle$ ,  $\langle \text{eos} \rangle$ ] then
9:       break
10:    end if
11:   end while
12:    $y \leftarrow y + y_t$ 
13:   if  $\langle \text{visual\_query} \rangle$  detected in  $y_t$  then
14:     Extract visual query  $q_{vis} \leftarrow \text{PARSE}(y_t, \langle \text{visual\_query} \rangle, \langle \text{/visual\_query} \rangle)$ 
15:     Retrieve vision results  $d = \text{VISIONAGENT}(q_{vis}, V)$ 
16:     Insert visual results into rollout  $y \leftarrow y + d$ 
17:   else if  $\langle \text{request\_grounding} \rangle$  detected in  $y_t$  then
18:     Retrieve grounding results  $\text{clipTag} = \text{GROUNDINGAGENT}(q, \mathcal{S})$ 
19:     Insert clip tag into rollout  $y \leftarrow y + \text{clipTag} + \mathcal{S}(\text{clipTag})$ 
20:   else if  $\langle \text{answer} \rangle$  detected in  $y_t$  then
21:     Extract predicted answer  $\hat{y} \leftarrow \text{PARSE}(y_t, \langle \text{answer} \rangle, \langle \text{/answer} \rangle)$ ; Normalize  $\hat{y}$  (trim spaces/punctuation)
22:     Insert final answer into rollout  $y \leftarrow y + \text{"The answer is: " + } \hat{y}$ 
23:     return final answer  $\hat{y}$ 
24:   else
25:     Ask for rethink  $y \leftarrow y + \text{"The action is not correct. Only } \langle \text{visual\_query} \rangle, \langle \text{request\_grounding} \rangle, \text{ or } \langle \text{answer} \rangle.$ "
26:   end if
27:   Increment step count  $t \leftarrow t + 1$ 
28: end while
29: return final generated response  $y$  for  $q$ 
```

sion process: at each action step after reasoning the policy emits exactly one structured action token ($\langle \text{visual_query} \rangle$, $\langle \text{request_grounding} \rangle$, or $\langle \text{answer} \rangle$).

Trajectory. A full response terminates upon emitting $\langle \text{answer} \rangle \dots \langle \text{/answer} \rangle$ or reaching K steps. We index decision steps by $t \in \{0, 1, \dots, T\}$ with $T \leq K$. At step t , the policy π_θ first plans and then emits a contiguous action string a_t ending with exactly one closing tag from $\{\langle \text{/visual_query} \rangle, \langle \text{/request_grounding} \rangle, \langle \text{/answer} \rangle\}$. If not terminating, the system appends feedback from the invoked agent o_t (e.g., a vision observation or a clip tag) to the context for the next step.

Rewards. We use two simple, rule-based rewards as supervision for reinforcement learning: (i) *Structural validity* $r_t^{\text{fmt}} \in \{0, 1\}$ grants 1 if the action string contains exactly one top-level tag with proper closure and no extraneous text; otherwise 0. (ii) *Answer correctness* $r^{\text{ans}} \in [0, 1]$ is awarded at termination via exact match on the multiple-choice answer; if no valid $\langle \text{answer} \rangle$ appears, $r^{\text{ans}} = 0$.

Objective and optimization. We seek a policy that produces well-formed actions at every step and a correct final answer. To balance these goals, the trajectory reward return is $R(\tau) = \alpha \sum_{t=0}^T r_t^{\text{fmt}} + r^{\text{ans}}$ where $\alpha > 0$ weights the per-step structural shaping and r^{ans} supplies the terminal task reward. r_t^{fmt} encourages the master to emit exactly one correct action tag at each decision, while r^{ans} evaluates only the final $\langle \text{answer} \rangle$. If no valid and correct answer is produced, $r^{\text{ans}} = 0$.

We optimize the master agent with GRPO on sampled rollouts: for each episode, the policy generates an action sequence, receives structural rewards at action boundaries and a terminal answer reward, and we compute sequence-level advantages with a learned value baseline. Policy updates follow the GRPO objective with standard clipping and entropy regularization, while the grounding and vision agents remain frozen. This minimal, two-signal objective provides sufficient guidance to learn structured, multi-turn coordination without additional dense rewards.

Table 2: Performance on *LongTVQA* and *LongTVQA+*. The left block lists model attributes (*Agentic, Input, RL fine-tune*); the right block reports validation accuracy (%). GPT-4o and Gemini-2.5 Pro are *multimodal* baselines that process and accept the full long video directly. Methods labeled *Agentic* indicate the model operates as the MASTERAGENT; methods labeled *AgenticRL* additionally denote RL fine-tuning. Parenthesized **green** numbers denote absolute gains over the immediately preceding (non-agentic or non-RL) setting. We observe that: (i) our multi-agent framework, LONGVIDEOAGENT, consistently outperforms the non-agentic counterparts; (ii) agentic RL yields additional gains, especially for smaller open-source models; (iii) using frames provides visual evidence beyond subtitles, and generally outperforms subtitle-only inputs; (iv) closed-source models remain strong, but the gap narrows much when open-source models adopt agentic designs and agentic RL.

Method	Multi-agent	Input	RL Finetune	Accuracy (%)	
				LongTVQA	LongTVQA+
<i>Closed-source (M)LLMs</i>					
GPT-4o	✗	Subtitle+Frame	✗	70.78	78.32
Gemini-2.5 Pro	✗	Subtitle+Frame	✗	<u>78.90</u>	<u>81.28</u>
GPT5-mini	✗	Subtitle	✗	<u>62.40</u>	<u>66.70</u>
Agentic-GPT5-mini	✓	Subtitle+Frame	✗	71.11(+8.71)	78.90(+12.20)
Grok	✗	Subtitle	✗	76.90	81.80
Agentic-Grok	✓	Subtitle+Frame	✗	82.65(+5.75)	85.60(+3.80)
<i>Open-source LLMs</i>					
DeepSeek-R1(671B)	✗	Subtitle	✗	68.99	75.04
Agentic-DeepSeek-R1(671B)	✓	Subtitle+Frame	✗	70.30(+1.31)	79.70(+4.66)
Agentic-Qwen2.5(3B)	✓	Subtitle+Frame	✗	<u>23.50</u>	<u>27.70</u>
AgenticRL-Qwen2.5(3B)	✓	Subtitle+Frame	✓	47.40(+23.90)	50.10(+22.40)
Agentic-Qwen2.5(7B)	✓	Subtitle+Frame	✗	46.10	60.30
AgenticRL-Qwen2.5(7B)	✓	Subtitle+Frame	✓	60.20(+14.10)	70.80(+10.50)

4 Experiments

4.1 Datasets

We build *LongTVQA* and *LongTVQA+* on top of TVQA and TVQA+. TVQA spans six TV shows with 152.5K multiple-choice QAs over 21.8K clips (60–90s) with subtitles and moment annotations; questions require joint dialogue–visual reasoning (Lei et al., 2018). TVQA+ refines a subset with spatio-temporal grounding—adding precise timestamps and 310.8K frame-level boxes for referenced entities (29.4K QAs from 4,198 clips, mainly TBBT)—supporting joint QA and temporal/spatial localization (Lei et al., 2020).

To obtain *LongTVQA* and *LongTVQA+*, we aggregate all clips from the same TV episode into a single *episode-level* (hour-scale) sequence. For each episode, we merge the visual stream, subtitles, and all associated questions; clip timestamps are re-indexed into the episode timeline, and TVQA+ bounding boxes are preserved at their corresponding frames. Unless otherwise noted, we report results on the original validation splits after this episode-level aggregation.

4.2 Baselines

We include both open-source and closed-source models (see Table 2), including representative open-source LLMs such as *DeepSeek-R1* (Guo

et al., 2025) and *Qwen2.5-3B/7B* (Qwen et al., 2025), and closed-source models such as *Grok*, *GPT5-mini* (OpenAI, 2025), *GPT-4o* (OpenAI) and *Gemini 2.5 Pro* (Comanici et al., 2025). We adopt unified experimental settings to ensure comparability across backbones. We evaluate *base LLMs in a non-agent mode* and *our agent system*, with the base LLM serving as the master agent in both cases. In the non-agent mode, the model consumes the full subtitles and does not invoke grounding or vision agents. We then compare our agent system against the corresponding non-agent runs on the *same* backbones, so that observed gains can be attributed to agentic behavior rather than backbone differences. For open-source backbones, we also report results *with vs. without* reinforcement learning under the identical evaluation protocol described in §4.3. Closed-source models are evaluated as released, without additional training.

4.3 Experimental Setup

By default we use *Grok-4-fast-reasoning* for temporal localization and *GPT-4o* as the vision agent. The window context is set to 1, meaning the agent conditions on a single localized clip (no adjacent clips), and the maximum execution steps are $K=5$. All methods read the full episode subtitles. In the non-agent setting, no external modules are in-

Table 3: A case of LONGVIDEOAGENT showing the collaboration among multiple agents.

Question & Choices
Question: Where is Sheldon sitting when he is accompanied by a man? Choices: a0: The Mall a1: A Movie Theatre a2: A Park a3: A Bus Stop ✓ a4: The Store
LONGVIDEOAGENT Execution
MASTER AGENT: <think> Subtitles alone do not specify the location; I will first localize with grounding, then request a visual read. </think> <request_grounding> ⇒ GROUNDING AGENT
GROUNDING AGENT: Grounding information. Subtitles: “Sheldon: ... apparently we’re just two peas in a pod. ... Mrs Cooper: Here. Thank you.” Clip tag: <s05e06_seg02_clip_15>
MASTER AGENT: <think> The grounded subtitles still don’t reveal the setting; I’ll request a visual description of the localized clip. </think> <visual_query> visual description of scene in <i>The Big Bang Theory</i> S05E06 segment 02 clip 15 where Sheldon is with a man </visual_query> ⇒ VISION AGENT
VISION AGENT: Visual description. Sheldon is sitting on a bench at night in an urban sidewalk setting; a waste container and building windows are visible nearby.
MASTER AGENT: <think> Bench + sidewalk + trash can + windows strongly indicate a bus stop rather than a mall, theatre, store, or park. No further tools needed. </think> <answer> a3: A Bus Stop ✓ </answer> ⇒ USER

voked. In the agent setting, the master agent receives symbolic temporal tag(s) <clip_X>, which is produced by the grounding agent, that marks the grounded clip(s) on the episode timeline, and it may request on-demand *textual* observations from the vision agent via prompting for selected frames within that segment (e.g., objects and attributes, OCR/text, brief scene cues). The master agent only consumes text (subtitles, the <clip_X> tag, and optional visual observations); no raw images are passed to the master agent. We report results on the *validation* splits of LongTVQA and LongTVQA+, using *Answer Accuracy (Acc)* as the primary metric (the questions are multiple choice) and additionally *Grounding Accuracy* for experiments that involve clip grounding. For reinforcement learning, we use GRPO with a learning rate of 5×10^{-6} , up to 2,000 optimization steps, a KL coefficient of 10^{-3} , batch size 4, rollout count $N=4$, and temperature 1.0. Training Qwen2.5-7B took 12 hours on 4x NVIDIA H800 GPUs, while the 3B variant took 6 hours under the same setup.

4.4 Performance

Table 2 presents overall validation accuracy. Moving from the non-agent setting to our multi-agent

Table 4: Ablations and analysis.

(a) **Comparison of non-agent vs. multi-agent performance.** Agentic components progressively improve performance: adding grounding outperforms the non-agent baseline, and adding vision agent yields the best results.

Setting	Accuracy (%)
Non-agent (Text-only)	64.3
Multi-Agent (Grounding)	69.0
Multi-Agent (Grounding + Vision)	74.8

(b) **Effect of max steps K .** Increasing the MASTERAGENT step budget generally raises both grounding and overall accuracy until reaching a saturation position.

K	Grounding Accuracy (%)	Accuracy (%)
2	67.00	68.30
5	71.00	73.67
10	72.00	73.67

(c) **Effect of evidence window size.** Larger temporal windows supply richer context for grounding and vision.

Window	Grounding Accuracy (%)	Accuracy (%)
1	71.67	70.33
2	78.67	75.00
3	81.94	77.26

(d) **Vision model ablation.** Stronger VISIONAGENT shows higher overall accuracy, reflecting higher quality extraction of visual information from frames.

Vision model	Grounding Acc. (%)	Acc. (%)
Qwen3-VL-235B	71.00	73.67
GPT-4o	73.30	78.00

framework yields significant gains. This provides direct evidence for the effectiveness of a *multi-agentic* pipeline that can localize the relevant clips and performs targeted visual inspection. In addition, for several open-source LLMs (as master agent), reinforcement learning consistently improves over their inference-only counterparts under identical prompts and evaluation; notably, the Qwen2.5-7B model with RL attains accuracy comparable to *GPT-5-mini* (closed-source) on our protocol. Illustrative examples in Table 3 and Table 5 demonstrate the efficacy of our approach, with additional cross-domain and open-ended results provided in the appendix.

4.5 Ablation Studies and Analysis

Execution step limit K . Table 4b varies the *upper bound* on agent actions per question. Increasing K from 2 to 5 raises temporal localization accuracy from 67.0 to 71.0 and answer accuracy from 68.30 to 73.67. As increasing the limit to $K = 10$ brings no further improvement, we set $K = 5$ as the default.

Table 5: Example showing LONGVIDEOAGENT reasoning and repeatedly invoking the GROUNDING and VISION agents to gather evidence and answer correctly.

Question & Choices
Question: What side of the bed is Sheldon when he is closer to the window? Choices: a0: the left side ✓ a1: the right side a2: the middle a3: not on the bed a4: the chair
LONGVIDEOAGENT Execution
MASTER AGENT: <think> Subtitles alone cannot disambiguate left vs. right; I will first ground the segment, then request visual evidence. </think> <request_grounding> localize the referenced segment </request_grounding> ⇒ GROUNDING AGENT
GROUNDING AGENT: Grounding information. Subtitles (localized): <s10e04_seg02_clip_09> "Sheldon: Okay. Now, on this side, ... Bernadette: Raj knows,"
MASTER AGENT: <think> Grounded subtitles still don't reveal which side is nearer the window; I'll request a visual read. </think> <visual_query> Big Bang Theory season 10 episode 4 — bedroom scene layout (Sheldon & Amy); identify bed/window configuration </visual_query> ⇒ VISION AGENT
VISION AGENT: Visual description (1). Bedroom scene; bed centered against the wall; Sheldon sits on the bed; general layout and decor described, but the side nearest the window is not specified.
MASTER AGENT: <think> The first visual read lacks the window-bed side relation; I need a more precise visual query focused on the window's position relative to Sheldon. </think> <visual_query> The Big Bang Theory S10E04 seg02 clip 09 — which side of the bed is next to the window, and where is Sheldon relative to it? </visual_query> ⇒ VISION AGENT
VISION AGENT: Visual description (2). A window with patterned curtains is on the left side of the bed; Sheldon is seated against the headboard near that window.
MASTER AGENT: <think> Window is on the left side of the bed and Sheldon is closer to it; therefore he is on the left side. No further tools needed. </think> <answer> a0: the left side ✓ </answer> ⇒ USER

Window context. With GPT-4o as the vision agent, expanding the evidence window from 1 to 2 adjacent clips improves localization from 71.67 to 78.67 (+7.00) and answer accuracy from 70.33 to 75.00 (+4.67); using 3 clips yields further gains to 81.94 and 77.26 (additional +3.27 and +2.26), as shown in Table 4c. The adjacent context helps disambiguate references across shots, but larger windows require more visual queries and latency while delivering smaller accuracy gains; hence we keep *Window*= 1 for the default and study larger windows in ablations.

Vision model. Table 4d compares perception backbones. GPT-4o attains 73.30 localiza-

tion and 78.00 answer accuracy, outperforming Qwen3-VL-235B-a22b at 71.00 and 73.67 by +2.30 and +4.33, respectively. The gap indicates that stronger visual recognition (small objects, OCR, fine attributes) translates into better end-task accuracy in long-form QA, so we adopt GPT-4o as the default vision agent.

Contribution of agentic components. Table 4a decomposes the gains when moving from a single LLM to a multi-agent, multimodal system. Adding temporal grounding to the same backbone increases answer accuracy from 64.3 to 69.0 (+4.7), showing that identifying the relevant clip filters distractors and focuses reasoning. Enabling vision after grounding further lifts accuracy to 74.8 (+5.8 over grounding; +10.5 overall): targeted visual inspection complements subtitles with concrete object/text cues and can validate or refine grounding through repeated calls when uncertain. Because backbones and prompts are held fixed, these improvements are attributable to the agentic procedure. We suggest grounding narrows the context length for reasoning and guides the master agent's attention, while vision supplies the missing fine-grained evidence.

5 Conclusion

We presented a multi-agent framework, LONGVIDEOAGENT, for long-form video question answering in which a MASTER agent coordinates a GROUNDINGAGENT for temporal localization and a VISIONAGENT for targeted perception. The framework is model-agnostic: we evaluate it with both closed- and open-source LLMs; for open-source masters, we fine-tune with GRPO to encourage accurate, concise, and cooperation-efficient behavior while keeping the other agents frozen. Equipped with a unified context and GRPO training that combines structural and answer rewards, the system where open-source LLMs act as the master agent yields transparent, step-by-step traces and achieves strong gains on *LongTVQA* / *LongTVQA+* over non-agent baselines. Ablations show that grounding+vision is essential, modest step limits suffice, adjacent-window context helps, and stronger perception yields higher accuracy, validating the effectiveness of the framework. Future work includes richer modalities (like audio track and knowledge background), finer grounding and larger-scale RL training.

Limitations

Our work has several practical limitations. First, based on TVQA and TVQA+, we rely on provided subtitles as the primary textual channel and do not process raw audio; in future work we plan to integrate an audio-to-subtitles (ASR) module to capture raw speech. Second, the vision and grounding modules are kept fixed during RL. Jointly optimizing them could further improve robustness and accuracy. Lastly, the reward is intentionally simple (format + answer correctness), which may still have room for improvements.

Acknowledgments

The work described in this paper was supported by a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China (Project Reference Number: AoE/E-601/24-N).

References

- Josh Achiam and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Jean-Baptiste Alayrac and 1 others. 2022. Flamingo: A visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*.
- Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and Mohamed Elhoseiny. 2024. Minigt4-video: Advancing multimodal llms for video understanding with interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, Hewei Guo, Conghui He, Botian Shi, Zhenjiang Jin, Chao Xu, and 16 others. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Preprint, arXiv:2404.16821*.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.
- Wenliang Dai, Junnan Li, Dongxu Li, and 1 others. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*.
- Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. 2024. VideoAgent: A memory-augmented multimodal agent for video understanding. *arXiv preprint arXiv:2403.11481*.
- Chaoyou Fu and 1 others. 2025. Video-mme: The first-ever comprehensive evaluation benchmark of multimodal large language models in video analysis.
- Jingjing Gao, Runzhou Ge, Kai Chen, and Ram Nevatia. 2018. Motion-appearance co-memory networks for video question answering. In *CVPR*. Dual attention over appearance and motion.
- Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Wei Han, Hui Chen, Min-Yen Kan, and Soujanya Poria. 2023. Self-adaptive sampling for efficient video question-answering on image-text models. *arXiv preprint arXiv:2307.04192*.
- Kairui Hu and 1 others. 2025. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional videos. *arXiv preprint arXiv:2501.13826*, arXiv:2501.13826.
- Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. 2024a. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13700–13710.
- Yang Jin, Zhicheng Sun, Kun Xu, Liwei Chen, Hao Jiang, Quzhe Huang, Chengru Song, Yuliang Liu, Di Zhang, Yang Song, and 1 others. 2024b. Video-lavit: Unified video-language pre-training with decoupled visual-motional tokenization. *arXiv preprint arXiv:2402.03161*.
- Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. 2021. Less is more: Clipbert for video-and-language learning via sparse sampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7331–7341.
- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. 2018. Tvqa: Localized, compositional video question answering. In *EMNLP*.
- Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. 2020. Tvqa+: Spatio-temporal grounding for video question answering. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 8211–8225.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and 1 others. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*.
- Yanwei Li, Chengyao Wang, and Jiaya Jia. 2024b. Llama-vid: An image is worth 2 tokens in large language models. In *European Conference on Computer Vision*, pages 323–340. Springer.
- Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. 2024. Video-llava: Learning united visual representation by alignment before projection. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pages 5971–5984.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*. LLaVA.
- Lin Long, Yichen He, Wentao Ye, Yiyuan Pan, Yuan Lin, Hang Li, Junbo Zhao, and Wei Li. 2025. Seeing, listening, remembering, and reasoning: A multimodal agent with long-term memory. *arXiv preprint arXiv:2508.09736*.
- Yongdong Luo, Xiawu Zheng, Guilin Li, Shukang Yin, Haojia Lin, Chaoyou Fu, Jinfa Huang, Jiayi Ji, Fei Chao, Jiebo Luo, and 1 others. 2024. Video-rag: Visually-aligned retrieval-augmented long video comprehension. *arXiv preprint arXiv:2411.13093*.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Khan. 2024. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12585–12602.
- Reiichiro Nakano and 1 others. 2021. Webgpt: Browser-assisted question-answering with human feedback. In *NeurIPS*.
- OpenAI. [Hello GPT-4o](#).
- OpenAI. 2025. [Introducing GPT-5](#).
- Junting Pan, Ziyi Lin, Yuying Ge, Xiatian Zhu, Renrui Zhang, Yi Wang, Yu Qiao, and Hongsheng Li. 2023. Retrieving-to-answer: Zero-shot video question answering with frozen large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 272–283.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- Renjie Pi, Tianyang Han, Wei Xiong, Jipeng Zhang, Runtao Liu, Rui Pan, and Tong Zhang. 2024. Strengthening multimodal large language model with bootstrapped preference optimization. In *European Conference on Computer Vision*, pages 382–398. Springer.
- Renjie Pi, Kehao Miao, Li Peihang, Runtao Liu, Jiahui Gao, Jipeng Zhang, and Xiaofang Zhou. 2025. Pointing to a llama and call it a camel: On the sycophancy of multimodal large language models. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 20177–20191.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Jacob Andreas. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*. Introduces Self-Ask and search integration.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Alec Radford, Jong Wook Kim, Chris Hallacy, and 1 others. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.
- David Romero and Tamar Solorio. 2024. Question-instructed visual descriptions for zero-shot video question answering. *arXiv preprint arXiv:2402.10698*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, Zhuang Liu, Hu Xu, Hyunwoo J. Kim, Bilge Soran, Raghuraman Krishnamoorthi, Mohamed Elhoseiny, and Vikas Chandra. 2024. LongVU: Spatiotemporal adaptive compression for long video-language understanding. *arXiv preprint arXiv:2410.17434*.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, and 1 others. 2023. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. In *NeurIPS*.
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, and 1 others.

2024. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In *ICCV*.
- Reuben Tan, Ximeng Sun, Ping Hu, Jui-hsien Wang, Hanieh Deilamsalehy, Bryan A Plummer, Bryan Russell, and Kate Saenko. 2024. Koala: Key frame-conditioned long video-llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13581–13591.
- Jiaqi Tang, Jianmin Chen, Wei Wei, Xiaogang Xu, Runtao Liu, Xiangyu Wu, Qipeng Xie, Jiafei Wu, Lei Zhang, and Qifeng Chen. 2025. [Robust-r1: Degradation-aware reasoning for robust visual understanding](#). *Preprint*, arXiv:2512.17532.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Ke-Yang Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *ArXiv*, abs/2409.12191.
- Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. 2024b. VideoAgent: Long-form video understanding with large language model as agent. *arXiv preprint arXiv:2403.10517*.
- Yi Wang, Xinhao Li, Ziang Yan, Yinan He, Jiashuo Yu, Xiangyu Zeng, Chenting Wang, Changlian Ma, Haiyan Huang, Jianfei Gao, and 1 others. 2025a. Internvideo2. 5: Empowering video mllms with long and rich context modeling. *arXiv preprint arXiv:2501.12386*.
- Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal. 2025b. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3272–3283.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. 2024. Longvlm: Efficient long video understanding via large language models. In *European Conference on Computer Vision*, pages 453–470. Springer.
- Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. 2024a. Pllava: Parameter-free llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*.
- Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin Dehghan. 2024b. Slowfast-llava: A strong training-free baseline for video large language models. *arXiv preprint arXiv:2407.15841*.
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, and 1 others. 2023. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.
- Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. 2024a. A simple llm framework for long-range video question-answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21715–21737.
- Hang Zhang, Xin Li, and Lidong Bing. 2023. Videollama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Hao-ran Tan, Chunyuan Li, and Ziwei Liu. 2024b. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024c. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and 1 others. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.
- Orr Zohar, Xiaohan Wang, Yann Dubois, Nikhil Mehta, Tong Xiao, Philippe Hansen-Estruch, Licheng Yu, Xiaofang Wang, Felix Juefei-Xu, Ning Zhang, and 1 others. 2025. Apollo: An exploration of video understanding in large multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18891–18901.

A Appendix

A.1 LLM Usage Disclosure

We used OpenAI ChatGPT only for minor language edits (grammar, spelling, clarity). It did not contribute to the ideas, methods, experiments, analyses, or technical content, and all edits were reviewed and verified by the authors.

A.2 Cross-Domain Generalization Beyond TV Shows

To evaluate generalization beyond TV shows and dialogue-driven reasoning, we additionally report results on two long-video benchmarks: VideoMMU and VideoMME (Hu et al., 2025; Fu et al., 2025). VideoMMU comprises expert-level educational videos across six disciplines and 30 subjects, while VideoMME covers six visual domains with multiple question types and videos up to 60 minutes. In subtitle-sparse or subtitle-absent settings, we also evaluate using ASR-generated transcripts so the framework can still receive language context for retrieval and reasoning. As summarized in Table 6, the multi-agent variants improve across both benchmarks in this heterogeneous setting, providing appendix-only evidence that the framework can transfer beyond TV-show-heavy dialogue data.

A.3 Open-Ended QA Extension

To address whether the framework generalizes beyond multiple-choice supervision, we evaluate a small open-ended extension on LongTVQA+. In this setting, we remove answer options and train the master to generate a free-form textual answer directly. This setting also removes the rule-based answer reward used in the main paper. Instead, during training we use an external LLM judge (Grok-4) that scores each model prediction against the gold answer on a 0–2 scale. A score of 0, 1, or 2 indicates an incorrect, partially correct, or correct answer, respectively. Table 7 reports preliminary results. Compared with the no-RL baseline, the Qwen2.5-3B model with RL improves from 0.53 to 0.84, suggesting that the framework may also benefit open-ended learning. We present this as preliminary evidence that the method is not restricted to multiple-choice formulation.

A.4 More cases of LONGVIDEOAGENT

Tables 8, 9, and 10 provide additional cases of LONGVIDEOAGENT.

Table 6: Cross-domain long-video evaluation beyond TV shows using VideoMMU and VideoMME. Scores are reported as accuracy (%). The three agentic settings use GPT-5[M] + Grok[G] + Grok[V], Grok4[M] + Grok[G] + Grok[V], and DeepSeekR1[M] + Qwen[G] + Qwen-VL[V], respectively, where [M], [G], and [V] denote the master, grounding, and vision models. For subtitle-sparse inputs, transcript-like text is obtained from ASR when subtitles are unavailable.

Model	Agentic	Open-source	VideoMMU	VideoMME
Gemini-2.5 Pro	No	No	62.67	72.56
GPT-5[M]	Yes	No	78.28	74.88
Grok4[M]	Yes	No	80.63	71.02
DeepSeekR1[M]	Yes	Yes	77.59	68.73

Table 7: Open-ended QA extension on LongTVQA+ evaluated with an external LLM judge (Grok-4) on a 0–2 scale. This appendix-only study removes answer choices and does not use the rule-based answer reward.

Method	Size	Agentic	RL	LLM-judge
Grok-4	-	Yes	No	1.10
DeepSeek	671B	Yes	No	1.02
Qwen2.5-72B	72B	Yes	No	0.88
Qwen2.5-3B (no RL)	3B	Yes	No	0.53
Qwen2.5-3B (with RL)	3B	Yes	Yes	0.84

Table 8: Additional case study 1.

Question & Choices
Question: What color shirt was Sheldon wearing when he was knocking on Penny’s door ? Choices: a0: Blue a1: Black a2: Yellow ✓ a3: Red a4: Purple
LONGVIDEOAGENT Execution
MASTER AGENT: <code><think></code> ... To determine the color of Sheldon’s shirt, visual information from this specific clip is required. I lack that knowledge, so I need to search for visual details of the scene. <code></think></code> <code><visual_query></code> visual description of scene in <i>The Big Bang Theory</i> S05E06 segment 02 clip 15 where Sheldon is with a man <code></visual_query></code> ⇒ VISION AGENT VISION AGENT: <i>Visual description.</i> He’s wearing a yellow short-sleeve T-shirt with a repeating multicolored (red/green/yellow) Rubik’s-cube print, layered over a light gray long-sleeve shirt.
MASTER AGENT: <code><think></code> The Vision Agent described Sheldon’s shirt in <i>The Big Bang Theory</i> S05E06, segment 02, clip 15 as a yellow T-shirt with a Rubik’s cube print. The question asks for the shirt color when Sheldon knocks on Penny’s door, and option a2 (Yellow) is correct. Thus, the answer is a2. <code></think></code> <code><answer></code> a2 ✓ <code></answer></code> ⇒ USER

Table 9: Additional case study 2.

Question & Choices
<p>Question: What was on the counter in front of Leonard when he was saying he had some friends who were not genius ?</p> <p>Choices: a0: A bowl of salad . a1: A pot roast . a2: A plate of steaks . a3: Several bottles . ✓ a4: A computer .</p>
LONGVIDEOAGENT Execution
<p>MASTER AGENT: <code><think></code> The provided subtitles describe the dialogue but contain no visual details about the counter in front of Leonard during the key line ("some of my best friends aren't geniuses"). To determine what was on the counter, visual information from the clip is necessary, as the text alone is insufficient. <code></think></code> <code><visual_query></code> Big Bang Theory S03E06 seg02 clip03 visual description kitchen counter Leonard Penny friends conversation <code></visual_query></code> ⇒ VISION AGENT</p> <p>VISION AGENT: <i>Visual description.</i> When Leonard delivers that line the kitchen/living-room counters and table are cluttered. Key visible items: - A cluster of brown beer bottles (several) on the coffee/table area in front of the couch. - Additional clear glass bottles and smaller glass containers mixed in with the beer bottles. - An aluminum soda/beer can (blue) on the counter near the bottles. - A plastic water bottle (clear) on the counter. - A yellow plastic bowl (being held/placed on the table). - A pizza box on the dining counter. - Snack/chip bags near the kitchen area. - Blue disposable plates and a red plastic cup on the table. - Various condiment/olive-oil-type bottles and small jars near the kitchen prep area. - Refrigerator with photos/magnets in the background. Overall it's a party/after-party spread: multiple empty/half bottles, a pizza box, snack bags, disposable plates/cups and a yellow serving bowl.</p> <p>MASTER AGENT: <code><think></code> The subtitles pinpoint the key line to the end of the clip, during the conversation in the kitchen area. The visual description confirms the setting is the cluttered kitchen counter with party remnants, including multiple beer bottles, clear glass bottles, a water bottle, a can, plates, bowls, pizza box, etc. The question asks what was on the counter in front of Leonard, and option a3 (Several bottles .) matches the description of several bottles. Thus, the answer is a3. <code></think></code> <code><answer></code> a3 ✓ <code></answer></code> ⇒ USER</p>

Table 10: Additional case study 3.

Question & Choices
<p>Question: What is behind Penny when she asks how the potato clock works ?</p> <p>Choices: a0: Curtains a1: Window a2: Door a3: Books ✓ a4: Table</p>
LONGVIDEOAGENT Execution
<p>MASTER AGENT: <code><think></code> The provided subtitles from clip discuss Bernadette, Howard, Sheldon, and Penny waiting for someone (likely Professor Proton), but make no mention of a potato clock or Penny asking about it. There are significant discrepancies between the provided information and the question, and no visual details are given. I should request grounding to relocate to the correct clip containing the potato clock scene with Penny. <code></think></code> <code><request_grounding></code> What is behind Penny when she asks how the potato clock works in The Big Bang Theory? <code></request_grounding></code> ⇒ GROUNDING AGENT</p> <p>GROUNDING AGENT: <i>New clip.</i> <code><clip_12></code> UNKNAME : I'm sorry to hear about your troubles . Uh , thanks . <code><eos></code> Penny : But if you do n't mind me asking , uh , the potato clock , how does it work ? (...) <code><eos></code> UNKNAME : No , I get that . <code></clip_12></code></p> <p>MASTER AGENT: <code><think></code> The initial clip was irrelevant, and this is the correct relocated scene. However, to determine what's behind Penny (e.g., curtains, window, door, books, table), I need visual details of the scene, which are not available in the text. <code></think></code> <code><visual_query></code> What is behind Penny in The Big Bang Theory season 6 episode 22 scene where she asks Professor Proton about the potato clock? <code></visual_query></code> ⇒ VISION AGENT</p> <p>VISION AGENT: <i>Visual description.</i> A bookshelf filled with books.</p> <p>MASTER AGENT: <code><think></code> The search for visual details returns a description of a bookshelf filled with books behind her, which aligns with option a3 (Books). This is sufficient to answer without discrepancies or need for further actions. <code></think></code> <code><answer></code> a3 ✓ <code></answer></code> ⇒ USER</p>