

TreeReasoner: Reinforcing Tool-Augmented Tree-of-Videos Reasoning

Anonymous CVPR submission

Paper ID 3

Abstract

001 We present *TreeReasoner*, a tool-augmented, tree-structured
 002 reasoning framework that recasts long-video understanding
 003 as an active hypothesis-verification problem over a vast vi-
 004 sual search space. By maintaining multiple parallel reason-
 005 ing paths, the model systematically explores the tempo-
 006 ral dimension and, guided by intermediate hypotheses,
 007 invokes frame-level tools such as temporal zooming, tem-
 008 poral jumping, and sliding to incrementally search a min-
 009 imal yet sufficient chain of evidence. The entire frame-
 010 work is trained end-to-end with *Tree-of-Tool Relative Policy*
 011 *Optimization (ToT-RPO)* following a supervised fine-tuning
 012 warmup, achieving superior video-understanding accuracy
 013 while decoding far fewer frames than existing methods and
 014 exhibiting interpretable temporal localization and causal-
 015 verification behaviors. Experiments on six long-video reason-
 016 ing benchmarks show that *TreeReasoner* consistently
 017 outperforms both standard IO and naive tool-calling base-
 018 lines. Transferability experiments on hallucination further
 019 confirm its generalization and reduced hallucination ten-
 020 dencies. These experiments validate the stability and effi-
 021 ciency of *TreeReasoner* in complex temporal scenarios.

022 1. Introduction

023 Previous progress in multimodal large language models
 024 (MLLMs) [1, 13, 14, 30, 36, 37, 42, 44, 46, 57] has driven
 025 significant advancements in the ability of models to under-
 026 stand videos. However, long-video understanding remains
 027 an unresolved and open challenge. Unlike image under-
 028 standing or short video understanding, long-video under-
 029 standing requires processing visual contextual information
 030 spanning several hours or even longer. Due to constraints
 031 in computation and memory, it is impractical to conduct a
 032 comprehensive analysis of the entire visual content.

033 Mainstream end-to-end long-video understanding solu-
 034 tions typically employ uniform frame sampling at fixed in-
 035 tervals or with a fixed number of frames. They alleviate
 036 computational burdens through visual compression [9, 19,
 037 21, 31] or by expanding the context length of MLLMs [55,

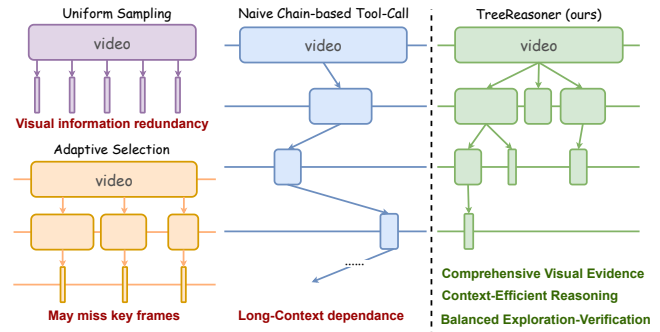


Figure 1. **Schematic illustrating different approaches for long-video understanding.** Left: Traditional sampling methods. *Uniform Sampling* (top) suffers from visual information redundancy, while *Adaptive Selection* (bottom) may miss key frames. Middle: A *Naive Chain-based Tool-Call* approach, which relies on a single reasoning path and can suffer from long-context dependence. Right (Ours): The *TreeReasoner* framework, which uses a parallel, tree-structured search to achieve comprehensive visual evidence, context-efficient reasoning, and balanced exploration-verification.

64], aiming to incorporate as many frames as possible for
 understanding. For most long-video understanding tasks,
 completing the task does not actually require accessing in-
 formation from the entire video. Consequently, a growing
 body of research has begun exploring more efficient
 video sampling strategies [2, 43, 47, 49, 60, 61]. These ap-
 proaches leverage vision-language models to select frames
 that are most relevant to the specific task, thereby signifi-
 cantly reducing the computational overhead of long-video
 understanding (Fig. 1). Nevertheless, since frame sam-
 pling and the model’s frame understanding process often
 cannot be optimized in an end-to-end manner, these meth-
 ods, which are analogous to agent workflows, have clearly
 imposed limitations on performance of long-video under-
 standing systems (Fig. 2).

Recently, OpenAI o3 [26] has presented an alternative
 perspective on visual understanding and reasoning: visual
 perception can be treated as a tool for invocation, enabling
 the model to autonomously learn how to identify critical
 visual cues. This insight inspires us to develop an end-to-
 end optimization framework, which empowers the model

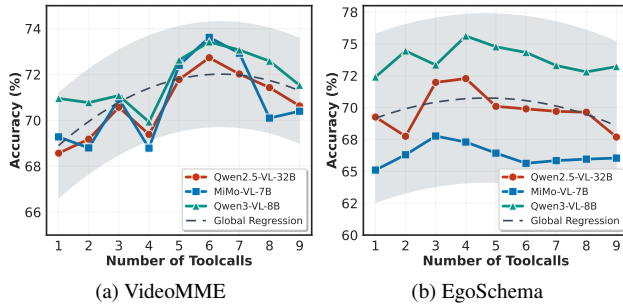


Figure 2. **Tool-call Scaling.** Empirical analysis on VideoMME (a) and EgoSchema (b) showing how performance varies under different fixed tool-call times.

to strategically explore the vast visual search space of long videos. The goal is to identify the minimal chain of evidence required to answer specific questions, thereby achieving efficient resolution of long-video understanding tasks.

In this work, we introduce *TreeReasoner*, a novel framework that reformulates video understanding as an active search problem through tool-augmented tree-of-thought reasoning. Our key insight is that effective video comprehension mirrors human cognitive processes: we don’t passively observe entire videos but actively seek specific temporal evidence through hypothesis-driven exploration. *TreeReasoner* achieves this by maintaining multiple reasoning trajectories in a tree structure, where each path represents a distinct hypothesis about the temporal location and relationship of relevant events. Crucially, this tree-based exploration is intrinsically coupled with tool utilization—the model employs frame extraction, temporal navigation, and region focusing tools to gather evidence, with each tool-call guided by the current reasoning state and contributing to trajectory expansion.

The synergy between tree-structured reasoning and tool augmentation addresses the fundamental challenges of video understanding. The tree structure enables systematic exploration of multiple temporal hypotheses in parallel, preventing premature commitment to potentially incorrect temporal interpretations. Meanwhile, tool augmentation transforms each reasoning step from passive prediction to active verification, allowing the model to adaptively zoom into specific temporal regions, extract key frames based on intermediate hypotheses, and incrementally build minimal evidence chains. This bidirectional relationship—where tree exploration guides tool selection and tool outputs inform trajectory prioritization—creates an efficient search mechanism through the vast temporal space.

We train *TreeReasoner* using Tree-of-Tool Relative Policy Optimization (*ToT-RPO*), which learns to balance exploration of diverse reasoning paths with exploitation of promising trajectories. Through reinforcement learning, the model develops sophisticated temporal search strategies, learning when to extract frames, how to navigate tempo-

ral neighborhoods, and which visual details require verification. Our experiments demonstrate that *TreeReasoner* significantly outperforms existing methods on challenging video understanding benchmarks, achieving superior accuracy while requiring substantially fewer processed frames. The emergent search behaviors reveal interpretable reasoning patterns, including temporal bracketing, causal chain verification, and adaptive temporal resolution adjustment.

2. Related Work

2.1. Visual Agentic Reasoning

Reinforcement Learning from Verifiable Rewards (RLVR) has achieved remarkable success in the domain of Large Language Models [8, 35, 56]. Models trained with algorithms such as GRPO/PPO demonstrate strong performance across a wide range of complex reasoning and agentic tasks. Inspired by these advances, recent efforts have extended RLVR-inspired paradigms to Vision-Language Models (VLMs), yielding promising results. One line of research enhances native visual reasoning through supervised fine-tuning (SFT) or reinforcement learning (RL), encouraging models to reason truly over visual inputs rather than relying solely on the text-based reasoning capabilities inherited from their LLM backbones [16, 18, 27, 32, 41, 48, 51, 53, 54, 59, 63]. Another complementary approach equips VLMs with external tools (e.g., zoom-in/search/crop/coding api, etc) and employs end-to-end RL to further strengthen their practical agentic capabilities [17, 22, 24, 38, 39, 65, 67]. In contrast to these works, this paper focuses on naive long-video agentic reasoning and proposes a tree-based reasoning framework that generalizes beyond simple chain-of-thought (CoT) processes.

2.2. Long-Video Understanding and Reasoning

Existing end-to-end long-video Multimodal Large Language Models (MLLMs) primarily fall into two categories. First, given the visual redundancy inherent in long videos, numerous studies [9, 19, 21, 31, 46, 57] have attempted to design visual compression algorithms to reduce the number of video tokens, thereby enabling long-video understanding with acceptable computational overhead. The second research direction [13, 55, 64] draws on techniques from long-context large language models, leveraging context extension to increase the input sequence length. Additionally, a substantial body of work has focused on developing complex Agent systems for long-video understanding [2, 43, 47, 49, 60, 61], their core design involves algorithms that select key frames from long videos for understanding and reasoning. However, due to the difficulty of implementing end-to-end optimization, such methods struggle to further push the boundaries of performance. Recently, there have been many efforts [3, 5, 7, 10, 12, 20, 40] to enhance video reasoning capabilities using RL, yet these

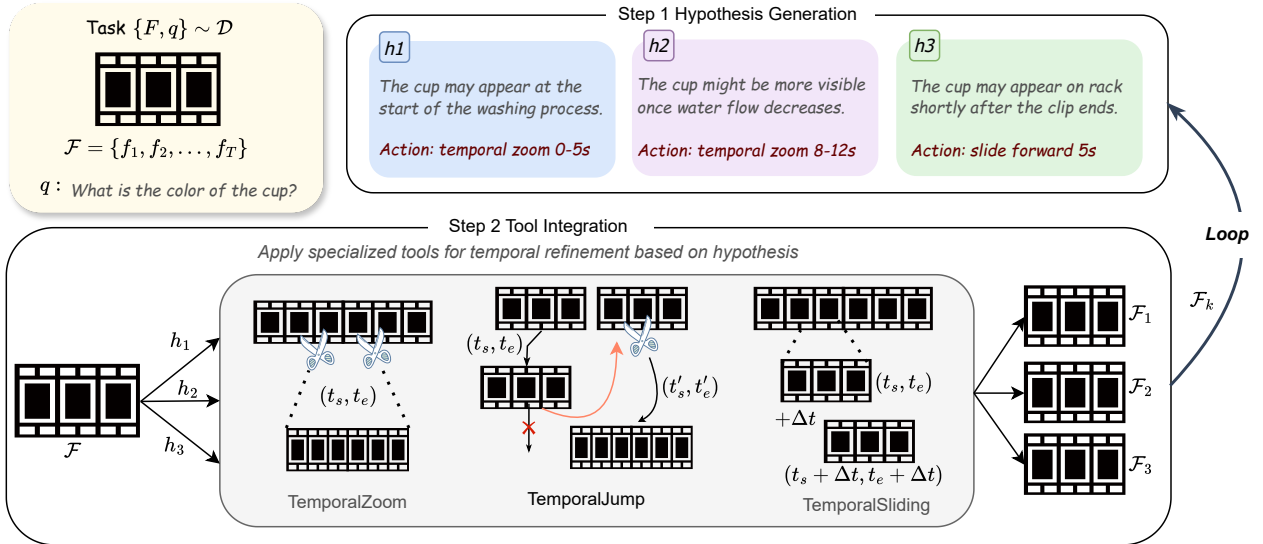


Figure 3. **Overview of TreeReasoner.** Given a question and a long video, the model first generates n parallel hypotheses about where and when relevant evidence might occur. For each hypothesis it invokes a specialized, frame-level tool—zoom, jump, or slide—to extract a short, task-specific clip. Every tool-call spawns a new child node, and the process is repeated breadth-first until an answer is reached or the budget is exhausted, yielding a minimal yet sufficient chain of evidence.

150 still fail to address the issue of excessive overhead dur- 178
 151 ing long-video understanding and reasoning. Unlike pre- 179
 152 vious approaches, this paper proposes a solution that sup- 180
 153 ports end-to-end optimization and enables efficient long- 181
 154 video understanding through temporal processing. 182

155 3. Methodology 183

156 In this section, we introduce our tool-augmented video re- 185
 157asoning as an active search and critic over temporal hypo- 186
 158theses under partial observability with tree-based expansion 187
 159policy. Specifically, in section 3.1, we formulate our theo- 188
 160retical framework on learning objectives and constraints. 189
 161Then in section 3.2, we dive into details of our reasoning 190
 162process end-to-end. Finally, in section 3.3 we introduce 191
 163the reward formulation and training strategy based on our 192
 164proposed Tree-of-Tool Relative Policy Optimization (ToT- 193
 165RPO) algorithm for tree-based video reasoning.

166 3.1. Problem Formulation 194

167 We formulate video reasoning process as an active search 194
 168and criticism over temporal hypotheses under partial ob- 195
 169servability via the guidance of a tree-based policy. Given 196
 170a video sequence $\mathcal{F} = \{f_1, f_2, \dots, f_T\}$ with T total frames 197
 171and a language query q , our reasoning objective is to dis- 198
 172cover or search an optimal evidence set $\mathcal{E}^* \subseteq \mathcal{F}$ that maxi- 199
 173mizes accuracy while minimizing computational cost. This 200
 174constitutes a multi-objective optimization problem: 201

$$175 \mathcal{E}^*, \tau^* = \arg \max_{\mathcal{E}, \tau} \mathbb{P}(y^* | q, \mathcal{E}, \tau) \quad \text{s.t.} \quad C(\mathcal{E}, \tau) \leq C_{\text{budget}}. \quad 202$$

176 Here y^* denotes the ground truth answer, τ represents the 203
 177 reasoning trajectory, $C(\mathcal{E}, \tau)$ measures computational cost, 204
 205

178 which contains both the number of reasoning trajectories 179
 179 and the computational complexity of each trajectory. C_{budget} 180
 180 means the controlled reasoning budget. The fundamen- 181
 181 tal challenge emerges from the exponential search space 182
 182 $|\mathcal{P}(\mathcal{F})| = 2^T$ over the frame sequences and the partial ob- 183
 183 servability constraint that limits frame access at any given 184
 184 time step. We model partial observability through a state- 185
 185 dependent visibility function $\mathcal{V}_t : \mathcal{S} \times \mathcal{A} \rightarrow 2^T$ that deter- 186
 186 mines which video segments become observable after tak- 187
 187 ing action a in state s . This formulation captures the realis- 188
 188 tic constraint that video understanding models cannot pro- 189
 189 cess entire video frames simultaneously due to memory and 190
 190 computational limitations. Additionally, we impose tempo- 191
 191 ral coherence constraints on searched evidence frame chains 192
 192 through conditional dependencies: 193

$$193 \forall \mathcal{E} = \{f_{i_1}, \dots, f_{i_k}\}, i_1 < \dots < i_k : \quad (1)$$

$$194 \mathbb{P}(f_{i_{j+1}} | f_{i_1}, \dots, f_{i_j}, q) > \mathbb{P}(f_{i_{j+1}} | q), \quad 195$$

194 thus ensuring that evidence frames form meaningful tempo- 194
 195 ral narratives rather than disconnected visual elements. 195

196 3.2. Tree-of-video Reasoning 196

197 Our approach models the reasoning process as a directed 197
 198 tree $\mathcal{T} = (\mathcal{N}, \mathcal{E}_{\text{tree}})$ where nodes represent reasoning 198
 199 states and edges encode tool-augmented state transitions. 199
 200 Given the challenge of searching through massive frame se- 200
 201 quences in long-video understanding, chain-like reasoning 201
 202 methods (e.g., Chain-of-Thought) may fail to accurately lo- 202
 203 calize relevant clips from the outset. Their depth-first ex- 203
 204 pansion tends to trap the policy model in irrelevant frame 204
 205 regions, with limited ability to recover or escape. Moreover, 205

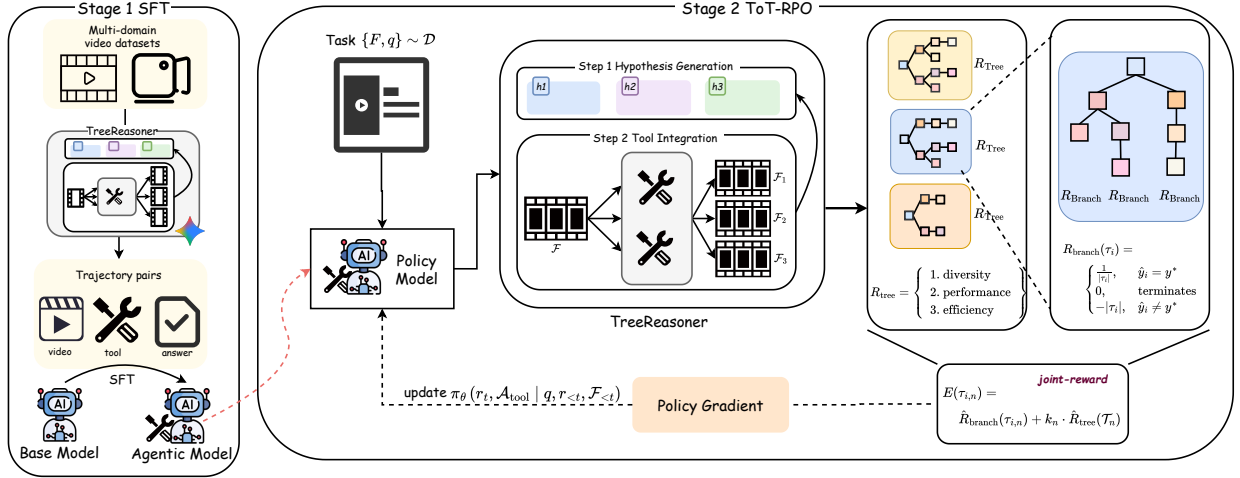


Figure 4. **Training paradigm of TreeReasoner.** *Stage 1:* the agent is warmed up with SFT on multi-turn, tool-augmented reasoning trajectories distilled from a teacher model. *Stage 2:* we continue with *ToT-RPO*, sampling entire reasoning trees per video-question pair and optimizing a composite reward that balances answer accuracy, search efficiency, tool utility, and inter-branch diversity. The whole pipeline is end-to-end and produces a single policy that jointly learns when to reason, which tool to call, and where to look next.

current RL algorithms face significant convergence difficulties in multi-turn tool-use settings. To address these issues, we adopt a breadth-first tree search strategy, enabling parallel exploration of multiple key frame intervals while progressively refining clip localization via tool interaction. This leads to earlier correct answer retrieval or timely early stopping, substantially improving reasoning efficiency.

3.2.1. Reasoning Process Representation

In our work, the reasoning state $s_t \in \mathcal{S}$ of each node $n_t \in \mathcal{N}$ is comprehensively characterized by a multi-modal tuple:

$$s_t = \langle q, \mathcal{H}_t, \mathcal{O}_t, \mathcal{R}_t, \mathcal{A}_t \rangle, \quad (2)$$

where q represents the input query. $\mathcal{H}_t = \{h_1, \dots, h_t\}$ denotes the sequence of hidden representations encoding temporal hypotheses. $\mathcal{O}_t = \{(o_i, t_i, b_i)\}$ contains observed visual elements o_i with their temporal locations t_i and spatial bounding boxes b_i . \mathcal{R}_t maintains the textual reasoning trace. $\mathcal{A}_t = \{a_1, \dots, a_{t-1}\}$ represents the history of actions taken up to the current state.

The state transition function $\delta : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ governs how states evolve through actions. For tool-based actions, the transition incorporates environmental feedback:

$$s_{t+1} = \delta(s_t, a_t) = \begin{cases} \delta_{\text{tool}}(s_t, a_t, \text{Env}(a_t, s_t)) & \text{if } a_t \in \mathcal{A}_{\text{tool}} \\ \delta_{\text{answer}}(s_t, a_t) & \text{if } a_t \in \mathcal{A}_{\text{answer}} \\ \delta_{\text{terminal}}(s_t, a_t) & \text{if } a_t \in \mathcal{A}_{\text{terminal}} \end{cases} \quad (3)$$

where $\text{Env}(a_t, s_t)$ executes the tool action with parameters θ_a in the video environment and returns observations.

3.2.2. Tool Integration

The action space $\mathcal{A} = \mathcal{A}_{\text{tool}} \cup \mathcal{A}_{\text{answer}} \cup \mathcal{A}_{\text{terminal}}$ encompasses three distinct categories of operations. Answer actions $\mathcal{A}_{\text{answer}}$ involve generating the final response directly

to the query based on the reasoning trace and collected visual evidence. Tool actions $\mathcal{A}_{\text{tool}}$ enable active video exploration with potential environmental interaction through parameterized functions. Inspired by natural human behaviors when watching videos, we design the following concrete tools:

- **TemporalZoom**(t_s, t_e) performs temporal zooming within the current clip, extracting frames at higher temporal resolution r from time interval $[t_s, t_e]$
- **TemporalJump**(t'_s, t'_e) jumps out of the current clip to extract frames from a different temporal interval $[t'_s, t'_e]$ in the global video, avoiding trapping in a local loop.
- **Sliding**($t_s, t_e, \Delta t$) slides a temporal window starting from interval $[t_s, t_e]$ with stride Δt , progressively exploring adjacent temporal segments.

These tool actions are practically useful as they mirror natural human video-watching behaviors: carefully examining details within interesting segments (zooming), jumping to other parts of the video to seek relevant evidence (jumping), and smoothly browsing continuous video content (sliding). Each tool action is parameterized by a continuous parameter vector $\theta_a \in \mathbb{R}^{d_a}$ learned through the policy network. This parameterization allows the policy model to learn optimal tool usage strategies through RL rather than relying on hand-crafted heuristics.

3.2.3. Hierarchical Tree Expansion

Tree expansion follows a principled hierarchical search strategy that balances exploration breadth with computational efficiency. At each state s_t , the policy network π_θ generates a probability distribution $\pi_\theta(a_t|s_t)$ over actions conditioned on the current state, where the policy network π_θ employs a transformer architecture with specialized attention mechanisms for processing temporal hypotheses, vi-

267 sual observations, and textual reasoning traces. First from a
268 high-level perspective, the search process maintains k par-
269 allel trajectories with diversity regularization. For each tra-
270 jectory i , we sample actions at timestep t according to the
271 constrained policy distribution:

$$a_t^{(i)} \sim \pi_\theta(\cdot | s_t^{(i)}) \quad \text{s.t.} \quad \sum_{j=1}^k \text{KL}[\pi(\cdot | s_t^{(i)}) || \pi(\cdot | s_t^{(j)})] \geq \beta_{\text{div}}, \quad (4)$$

272 where the diversity constraint β_{div} prevents trajectory col-
273 lapse and encourages exploration of alternative reasoning
274 paths. This constraint is enforced through a diversity-
275 augmented sampling procedure that rejects actions leading
276 to excessive similarity with existing trajectories.

277 Then, dive into the internal trajectory, the hierarchical
278 expansion strategy operates across multiple temporal scales
279 through a coarse-to-fine refinement process. Initially, the
280 search explores broad temporal regions using low tempo-
281 ral resolution sampling: $\mathcal{L}_{\text{coarse}}(t_s, t_e) = \{t_s + k \cdot \Delta_{\text{coarse}} : k \in \mathbb{Z}, t_s \leq t_s + k \cdot \Delta_{\text{coarse}} \leq t_e\}$ where Δ_{coarse} represents
282 the coarse temporal resolution. Promising regions identified
283 during coarse search are subsequently refined using higher
284 temporal resolution: $\mathcal{L}_{\text{fine}}(t_s, t_e) = \{t_s + k \cdot \Delta_{\text{fine}} : k \in \mathbb{Z}, t_s \leq t_s + k \cdot \Delta_{\text{fine}} \leq t_e\}$ with $\Delta_{\text{fine}} \ll \Delta_{\text{coarse}}$. Note that
285 this tree-structured, frame-granular enhancement expansion
286 strategy can enable the policy model to concurrently focus
287 on key frame intervals, allowing it to rapidly identify impor-
288 tant clues or terminate search early—avoiding wasted com-
289 putation on irrelevant frames or iterative trial-and-error over
290 incorrect frame parts (as in deep tool-call chain), thereby
291 significantly reducing inference time.

295 3.3. Learn to Tree-of-video Reasoning

296 3.3.1. Warmup SFT

297 To enable model to learn how to reason with tool use based
298 on external execution feedback, we first generate multi-
299 round multimodal reasoning data using Gemini2.5-Pro API
300 within our agent data workflow. For an input video query q ,
301 the Gemini2.5 model generates l rounds of reasoning output
302 R_{query} , we filter the output with wrong answers and finally
303 construct a SFT dataset \mathcal{D} . This dataset is used to train the
304 model to predict correct tool-use action $(r_t, \mathcal{A}_{\text{tool}})$ on input
305 query q , previously selected video frame sequence $\mathcal{F}_{<t}$ and
306 reasoning result $r_{<t}$. The SFT objective is to unlock the
307 policy’s reasoning ability with specific tool use. Formally,
308 we minimize following negative log-likelihood loss:

$$\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(q, r, \mathcal{A}_{\text{tool}}, \mathcal{F}) \sim \mathcal{D}} \left[\sum_{t=1}^T \log \pi_\theta(r_t, \mathcal{A}_{\text{tool}} | q, r_{<t}, \mathcal{F}_{<t}) \right]. \quad (5)$$

310 3.3.2. Reinforcement Learning Via ToT-RPO

311 Although SFT enables the model to imitate tool-augmented
312 reasoning trajectories, it cannot optimize long-term deci-

sion outcomes or balance exploration and efficiency. To ad-
dress these limitations, we further refine TreeReasoner with
Tree-of-Tool Relative Policy Optimization (ToT-RPO) [8],
which extends GRPO to structured reasoning trees for end-
to-end optimization over multi-branch trajectories. Af-
ter the SFT stage, the model is trained via *ToT-RPO*
by sampling a group of tree-based trajectories for each
video–query pair (Sec 3.2) and optimizing composite re-
wards that integrate accuracy, reasoning efficiency, tool uti-
lity, and exploration diversity, enabling adaptive and cost-
aware reasoning across tree-structured trajectories.

Algorithm 1 ToT-RPO Rollout

Require: An array of video-query pairs $Q = \{(q_1, \mathcal{F}_{1,0}), (q_2, \mathcal{F}_{2,0}), \dots, (q_n, \mathcal{F}_{n,0})\}$
Ensure: Rollout responses T of tree for all $(q, v) \in Q$.

- 1: $P \leftarrow Q$ ▷ Init root node of tree
- 2: **while** $P \neq \emptyset$ **do**
- 3: $S \leftarrow \text{INFERENCE}(P)$ ▷ Reason with action
- 4: $p^{\text{last}} \leftarrow P$
- 5: $P \leftarrow \emptyset$ ▷ Clean up the node queue
- 6: **for** s_k **in** S **do** ▷ Recongnize node states
- 7: **if** $a_k \in \mathcal{A}_{\text{terminal}}$ **or** $a_k \in \mathcal{A}_{\text{answer}}$ **then**
- 8: $T \leftarrow T \oplus \{p_k^{\text{last}} \oplus s_k\}$ ▷ Build tree
- 9: **else**
- 10: **for** $a_{k,l}$ **in** s_k **do** ▷ Expand child node
- 11: $s'_{k,l} \leftarrow \text{EXECUTE}(a_{k,l}, s_k)$ ▷ Execute
- 12: $P \leftarrow P \cup \{s'_{k,l}\}$
- 13: **end for**
- 14: **end if**
- 15: **end for**
- 16: **end while**
- 17: **return** T ▷ Return the final responses of tree

Tree-structured Rollout Process. In TreeReasoner, each
reasoning episode can be viewed as a stochastic tree-
search process, where the policy π_θ governs both reason-
ing and tool-usage actions across the tree nodes. Given a
video–query pair (v, q) , the model generates a reasoning
tree $\mathcal{T} = \{\tau_1, \tau_2, \dots, \tau_{|\mathcal{T}|}\}$ consisting of $|\mathcal{T}|$ branch tra-
jectories $\tau_i = \{(s_t, a_t)\}_{t=1}^{T_i}$, each representing a distinct
hypothesis exploration path ending with either an answer or
early termination signal. Unlike linear chain rollouts, *ToT-
RPO* treats these trajectories collectively during optimiza-
tion, leveraging their structural interdependence. To op-
erationalize this process, the reasoning tree is constructed
through breadth-first expansion: starting at the root node
 N_0 , a queue P is maintained to manage nodes at each
depth level, and for every node in P , the model proac-
tively generates reasoning analysis along with an associ-
ated action. These actions encompass tool-calling actions,

Table 1. **Performance of our TreeReasoner on video understanding benchmarks.** We compare our **TreeReasoner** against baseline models (IO), chain-based **Tool-call** method, and other SOTA models on 6 datasets. The 🔥 denotes models trained with both SFT and RL.

Model	Method	VideoMME	LongVideoBench	EgoSchema	MLVU	VSIBench	TOMATO	Average
OpenAI GPT-4o [15]	IO	71.9	66.7	72.2	64.6	34.0	37.7	57.9
Gemini2.5-Flash-Lite [6]	IO	65.0	60.2	–	69.3	27.0	–	–
Qwen2.5-VL-7B [1]	IO	65.1	56.0	65.0	70.2	34.2	27.6	53.0
Qwen2.5-VL-72B [1]	IO	73.3	60.7	76.2	74.6	37.2	31.2	58.9
Llama3.2-11B-Inst[25]	IO	46.0	45.5	54.3	44.4	20.6	21.5	38.7
Gemma3-12B-IT[34]	IO	58.2	51.5	56.9	52.3	32.4	28.6	46.7
MiMo-VL 7B [33]	IO	70.3	60.4	67.4	71.3	40.8	36.1	57.7
	Tool-call	72.3	63.1	69.2	74.1	45.1	38.4	60.4
	Tool-call 🔥	75.1	65.2	71.0	78.2	47.8	40.3	62.9
	TreeReasoner	73.2	63.3	69.6	77.1	46.7	38.5	61.4
	TreeReasoner 🔥	75.9	67.8	73.7	81.8	49.6	41.5	65.0
Qwen3-VL 8B [56]	IO	71.4	58.4	73.2	78.1	59.4	34.9	62.5
	Tool-call	73.3	61.5	75.6	80.8	61.8	37.3	65.1
	Tool-call 🔥	76.2	64.8	77.8	84.3	63.1	40.5	67.8
	TreeReasoner	73.7	62.2	76.4	81.5	62.6	37.8	65.7
	TreeReasoner 🔥	77.2	65.9	80.0	85.5	64.6	41.5	69.1
Qwen2.5-VL 32B [1]	IO	70.4	57.8	69.2	71.3	37.1	29.8	55.9
	Tool-call	72.3	59.1	71.8	75.4	43.7	31.5	59.0
	Tool-call 🔥	74.8	62.1	74.6	79.8	46.1	35.2	62.1
	TreeReasoner	72.5	60.2	73.0	76.4	44.7	33.2	60.0
	TreeReasoner 🔥	76.4	62.9	77.6	80.4	47.6	37.5	63.7

341 answer-generation actions, and early-stopping actions (see
 342 Sec. 3.2.1). When a tool-calling action is correctly pro-
 343 duced, the corresponding tool is executed and its feedback
 344 is appended to the next-layer nodes, which are then added
 345 to P ; otherwise, that branch halts. Each node can expand
 346 into at most W child nodes, and the overall depth of the
 347 tree is capped at D , ensuring a controlled yet expressive ex-
 348 ploratory reasoning structure.

349 **Hierarchical Reward Design.** During *ToT-RPO* training,
 350 the rollout process naturally produces a structured reason-
 351 ing tree rather than a single linear trajectory. This struc-
 352 tural property prevents direct application of conventional
 353 trajectory-level policy gradient updates. To address this,
 354 we employ a hierarchical reward formulation that evaluates
 355 both local branch trajectories and the global reasoning tree,
 356 enabling stable and efficient optimization.

357 At the branch-trajectory level, we define a reward func-
 358 tion that jointly captures prediction correctness and search
 359 efficiency:

$$360 \mathcal{R}_{\text{branch}}(\tau_i) = \begin{cases} \frac{1}{|\tau_i|}, & \text{if } \hat{y}_i = y^*, \\ 0, & \text{if the branch terminates,} \\ -|\tau_i|, & \text{if } \hat{y}_i \neq y^*, \end{cases} \quad (6)$$

361 where $|\tau_i|$ denotes the number of nodes in the branch tra-
 362 jectory and \hat{y}_i is the predicted answer. This formulation
 363 encourages the discovery of correct answers with minimal

reasoning depth while heavily penalizing long yet incorrect
 exploration paths.

Beyond branch evaluation, we introduce a tree-level re-
 ward to assess the overall quality of the reasoning tree \mathcal{T}_j .
 This reward integrates three complementary aspects: (1) *di-*
versity, encouraging exploration of distinct reasoning hy-
 potheses through average branch dissimilarity; (2) *perfor-*
mance, measured by the success rate of branch trajectories
 arriving at correct answers; and (3) *efficiency*, reflected by
 the depth of the reasoning tree. Together, these components
 ensure that the model learns to balance explorative breadth
 and computational economy.

To stabilize optimization, both branch-trajectory rewards
 and tree-level rewards are standardized across the batch.
 The combined advantage $E(\tau_{i,j})$ is defined as

$$E(\tau_{i,j}) = \hat{\mathcal{R}}_{\text{branch}}(\tau_{i,j}) + k_j \cdot \hat{\mathcal{R}}_{\text{tree}}(\mathcal{T}_j), \quad (7)$$

where

$$\begin{cases} \hat{\mathcal{R}}_{\text{branch}}(\tau_{i,j}) = \frac{\mathcal{R}_{\text{branch}}(\tau_{i,j}) - \mu(\mathcal{R}_{\text{branch}}(\cdot, j))}{\sigma(\mathcal{R}_{\text{branch}}(\cdot, j))}, \\ \hat{\mathcal{R}}_{\text{tree}}(\mathcal{T}_j) = \frac{\mathcal{R}_{\text{tree}}(\mathcal{T}_j) - \mu(\mathcal{R}_{\text{tree}})}{\sigma(\mathcal{R}_{\text{tree}})}. \end{cases} \quad (8)$$

Here, $\mu(\cdot)$ and $\sigma(\cdot)$ denote the empirical mean and standard
 deviation, respectively, computed over the collection of val-
 ues indicated within the parentheses.

The optimization objective follows the standard GRPO framework:

$$\mathcal{L}_{\text{ToT-RPO}} = -\mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{k=1}^K A_{\text{group}}(\tau) \log \pi_{\theta}(a_k | s_k) - \lambda_{\text{KL}} D_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}) \right]. \quad (9)$$

The group-normalized advantage is computed as

$$A_{\text{group}}(\tau) = E(\tau) \cdot \text{clip} \left(\frac{\pi_{\theta}(\tau)}{\pi_{\text{old}}(\tau)}, 1 - \delta, 1 + \delta \right) \quad (10)$$

where the clipping operation prevents excessively large updates and contributes to training robustness.

4. Experiments

4.1. Settings

Models and baselines. We select three representative families of SOTA open-source VLMs, including Qwen2.5-VL [1], Qwen3-VL [56], and MiMo-VL [33] to evaluate our TreeReasoner framework. For each selected model family, we first assess its baseline performance, then enhance its reasoning capabilities through both standard chain-based tool-calling and our proposed TreeReasoner paradigm, with and without SFT and end-to-end RL. Additionally, we include several strong closed-source models and other competitive open-source models in our comparison directly to comprehensively demonstrate the effectiveness and robustness of our approach.

Benchmarks. We evaluate on several challenging video understanding and reasoning benchmarks, including one general video understanding task (VideoMME [11]), three long-video understanding tasks (MLVU [68], LongVideoBench [50], EgoSchema [23]) and two complex video reasoning tasks (TOMATO [29], VSIBench [58]).

4.2. Performances of TreeReasoner

Table 1 shows the main performance of different baselines and our proposed TreeReasoner methods. Compared with the IO and naive Tool-Call baselines, TreeReasoner delivers consistent and substantial performance improvements across all six datasets. This demonstrates that TreeReasoner’s parallel hypothesis exploration and early evidence verification provide immediate benefits over linear Tool-call or single-pass inference. When equipped with SFT and RL training, the TreeReasoner framework yields additional performance gains and enables multiple backbone models to reach state-of-the-art results across several benchmarks. Case study for interpretability of TreeReasoner can be found in the appendix.

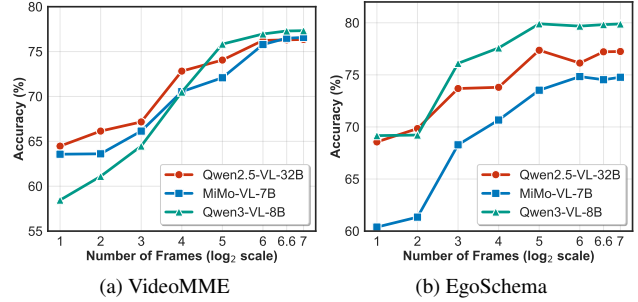


Figure 5. **Ablation of Frame Numbers.** (a) Accuracy with various Frame Numbers on VideoMME. (b) Accuracy with various Frame Numbers on EgoSchema.

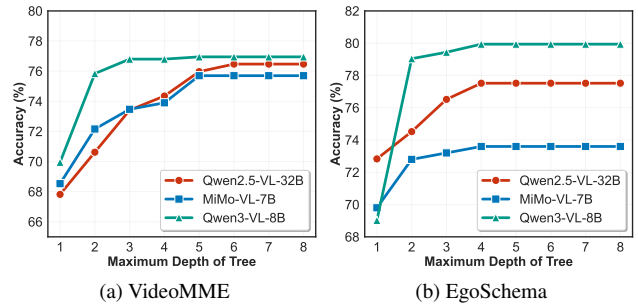


Figure 6. **Ablation of Max Tree Depth of TreeReasoner.** (a) Accuracy with various tree depth budgets on VideoMME. (b) Accuracy with various tree depth budgets on EgoSchema.

4.3. Analysis

Ablation on training stages. Both SFT and RL are compatible with the proposed TreeReasoner framework. To elucidate the impact of different training stages on final performance, we conducted an ablation study comparing various training configurations of TreeReasoner. As shown in Table 2, we can find that incorporating the RL stage consistently improves performance across various base models and benchmarks. These results indicate that TreeReasoner’s RL process effectively refines temporal reasoning and cross-frame consistency, providing stable performance boosts regardless of model scale.

Ablation on limit of exploring depth and frame numbers for each node. Tree depth and the number of frames per node are two critical hyperparameters in our TreeReasoner reasoning framework, playing a pivotal role in balancing inference accuracy and computational efficiency. We conduct comprehensive ablation studies on these key parameters to identify the optimal configuration and to demonstrate their adaptability and robustness across diverse tasks. As shown in Fig. 5, increasing the number of sampled frames per node improves performance until saturation. Similarly, Fig. 6 shows that increasing the maximum tree depth yields consistent gains before convergence. These results indi-

Table 2. **Ablation Study on Training Stages with Different Base Models.** Accuracy(%) comparison across six video benchmarks for comparison with different base models at both SFT stage and the combined SFT with RL (SFT+RL) stage.

Base Model	Stage	VideoMME	LongVideoBench	EgoSchema	MLVU	VSIBench	TOMATO	Average
MiMo-VL 7B	SFT	73.9	65.6	71.6	79.7	47.6	39.6	63.0
	SFT+RL	75.9	67.8	73.7	81.8	49.6	41.5	65.0
Qwen3-VL 8B	SFT	75.2	63.9	78.2	83.4	63.2	39.4	67.3
	SFT+RL	77.2	65.9	80.0	85.5	64.6	41.5	69.1
Qwen2.5-VL 32B	SFT	74.4	61.0	75.4	78.4	45.7	35.6	61.8
	SFT+RL	76.4	62.9	77.6	80.4	47.6	37.5	63.7

Table 3. **Hallucination.** Accuracy(%) on Halluciner benchmark.

Model	Method	Basic	Hallucinated	Overall
MiMo-VL 7B	IO	82.8	66.8	57.6
	Tool-call	84.0	67.6	58.6
	TreeReasoner	84.3	69.7	60.7
	Tool-call 🔥	87.7	72.6	65.2
	TreeReasoner 🔥	89.9	75.9	69.1
Qwen3-VL 8B	IO	79.4	78.2	61.8
	Tool-call	81.1	81.2	64.6
	TreeReasoner	83.5	83.0	66.5
	Tool-call 🔥	84.5	85.1	69.8
	TreeReasoner 🔥	86.0	87.1	72.9
Qwen2.5-VL 32B	IO	75.2	80.4	62.0
	Tool-call	76.2	81.8	62.8
	TreeReasoner	77.9	82.3	64.6
	Tool-call 🔥	79.4	87.0	67.5
	TreeReasoner 🔥	82.6	88.7	71.8

450 cate that enriching frame-level visual diversity enables more
 451 comprehensive evidence aggregation and strengthens the
 452 model’s structured visual reasoning process, while redun-
 453 dant frames offer limited additional benefit.

454 **Transferability on Hallucination.** A key challenge in
 455 the reasoning process of MLLM lies in their propensity
 456 to hallucinate—particularly when reasoning chains become
 457 overly long or complex. As attention distributions become
 458 sparse, the model may drift away from visual grounding,
 459 resulting in unfaithful or fabricated content. To assess the
 460 robustness and transferability of our reasoning framework
 461 under such conditions, we evaluate models on the Hal-
 462 luciner benchmark [45]. Each sample in Halluciner consists
 463 of a Basic question (testing factual understanding) and a
 464 Hallucinated counterpart (containing unverifiable or false
 465 premises). A point is awarded only when the model cor-
 466 rectly answers both questions in a pair, ensuring that robust-
 467 ness against hallucination is jointly assessed with factual
 468 consistency. Experimental results show that our proposed
 469 TreeReasoner framework achieves substantially lower hal-
 470 lucination rates compared to standard IO and naive Tool-
 471 call, demonstrating its stronger transferability and robust-
 472 ness in mitigating hallucination. As shown in Table 3,
 473 TreeReasoner consistently mitigates hallucination across
 474 various backbone models and configurations.

Table 4. **Efficiency.** Efficiency analysis on LongVideoBench.

Method	Output/Input tokens	Accuracy
IO (best of 6)	9.1k / 21.1k	63.5
Tool call (best of 6) 🔥	13.2k / 49k	67.4
TreeReasoner 🔥	13.6k / 36.2k	67.8

475 **Efficiency.** To evaluate the computational efficiency of
 476 our approach, we conduct a comprehensive comparison be-
 477 tween TreeReasoner and naive chain-based Tool Call on the
 478 LongVideoBench under the Pass@K setting, where K de-
 479 notes the ratio of token consumption relative to the naive
 480 Tool Call baseline. As shown in Table 4, TreeReasoner
 481 achieves superior performance while maintaining compara-
 482 ble token consumption to naive Tool Call Pass@K, demon-
 483 strating its efficiency advantage. This improvement stems
 484 from TreeReasoner’s ability to effectively explore multi-
 485 ple “local-global” video understanding paths in parallel,
 486 thereby circumventing the accumulated errors and global
 487 misjudgments that often arise from naive Tool Call’s se-
 488 quential, single-path exploration strategy.

489 5. Conclusion

490 In this work, we present TreeReasoner, a novel framework
 491 that recasts long-video understanding as an active search
 492 problem over temporal hypotheses. By maintaining mul-
 493 tiple parallel reasoning trajectories in a tree structure and
 494 strategically invoking frame-level tools—temporal zoom-
 495 ing, jumping, and sliding—our approach efficiently discov-
 496 ers minimal evidence chains without exhaustive frame pro-
 497 cessing. Trained end-to-end via *ToT-RPO*, TreeReasoner
 498 achieves state-of-the-art performance across six challeng-
 499 ing benchmarks while requiring substantially fewer frames
 500 than existing methods. Using MiMo-VL 7B as example,
 501 our experiments demonstrate consistent improvements of
 502 7.3% over corresponding standard IO version and 2.1%
 503 over naive chain-based tool-calling approaches with train-
 504 ing, validating that tree-structured exploration with tool
 505 augmentation provides a principled and efficient solution
 506 for long-video understanding. The interpretable emergent
 507 behaviors—including temporal bracketing and adaptive res-
 508 olution adjustment—further confirm that our framework
 509 learns sophisticated temporal reasoning strategies that mir-
 510 ror human video comprehension processes.

511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567

References

- [1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhao-hai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. 1, 6, 7
- [2] Boyu Chen, Zhengrong Yue, Siran Chen, Zikang Wang, Yang Liu, Peng Li, and Yali Wang. Lvagent: Long video understanding by multi-round dynamical collaboration of mllm agents. *arXiv preprint arXiv:2503.10200*, 2025. 1, 2
- [3] Xinlong Chen, Yuanxing Zhang, Yushuo Guan, Weihong Lin, Zekun Wang, Bohan Zeng, Yang Shi, Sihan Yang, Qiang Liu, Pengfei Wan, Liang Wang, and Tieniu Tan. Vidbridge-r1: Bridging qa and captioning for rl-based video understanding models with intermediate proxy tasks, 2025. 2
- [4] Yukang Chen, Wei Huang, Baifeng Shi, Qinghao Hu, Hanrong Ye, Ligeng Zhu, Zhijian Liu, Pavlo Molchanov, Jan Kautz, Xiaojuan Qi, Sifei Liu, Hongxu Yin, Yao Lu, and Song Han. Scaling rl to long videos. 2025. 14
- [5] Yukang Chen, Wei Huang, Baifeng Shi, Qinghao Hu, Hanrong Ye, Ligeng Zhu, Zhijian Liu, Pavlo Molchanov, Jan Kautz, Xiaojuan Qi, et al. Scaling rl to long videos. *arXiv preprint arXiv:2507.07966*, 2025. 2
- [6] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 6
- [7] Jisheng Dang, Jingze Wu, Teng Wang, Xuanhui Lin, Nannan Zhu, Hongbo Chen, Wei-Shi Zheng, Meng Wang, and Tat-Seng Chua. Reinforcing video reasoning with focused thinking. *arXiv preprint arXiv:2505.24718*, 2025. 2
- [8] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shiron Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiusi Du, Ruiqi Ge, Ruison Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan,
- S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wan-jia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. 2, 5
- [9] Jiajun Fei, Dian Li, Zhidong Deng, Zekun Wang, Gang Liu, and Hui Wang. Video-ccam: Enhancing video-language understanding with causal cross-attention masks for short and long videos. *arXiv preprint arXiv:2408.14023*, 2024. 1, 2
- [10] Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Junfei Wu, Xiaoying Zhang, Benyou Wang, and Xiangyu Yue. Video-r1: Reinforcing video reasoning in mllms. *arXiv preprint arXiv:2503.21776*, 2025. 2
- [11] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui Lin, Sirui Zhao, Ke Li, Tong Xu, Xiawu Zheng, Enhong Chen, Caifeng Shan, Ran He, and Xing Sun. Video-mme: The first-ever comprehensive evaluation benchmark of multimodal llms in video analysis, 2025. 7
- [12] Hongcheng Gao, Jiashu Qu, Jingyi Tang, Baolong Bi, Yue Liu, Hongyu Chen, Li Liang, Li Su, and Qingming Huang. Exploring hallucination of large multimodal models in video understanding: Benchmark, analysis and mitigation, 2025. 2
- [13] Google. Gemini 2.5 pro: A multimodal reasoning model for video, audio, and text. <https://deepmind.google/models/gemini/pro/>, 2025. 1, 2, 14
- [14] Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang, Jiawei Wang, Jingji Chen, Jingjia Huang, Kang Lei, Liping Yuan, Lishu Luo, Pengfei Liu, Qinghao Ye, Rui Qian, Shen Yan, Shixiong Zhao, Shuai Peng, Shuangye Li, Sihang Yuan, Sijin Wu, Tianheng Cheng, Weiwei Liu, Wenqian Wang, Xi-anhan Zeng, Xiao Liu, Xiaobo Qin, Xiaohan Ding, Xiaojun Xiao, Xiaoying Zhang, Xuanwei Zhang, Xuehan Xiong, Yanghua Peng, Yangrui Chen, Yanwei Li, Yanxu Hu, Yi Lin, Yiyuan Hu, Yiyuan Zhang, Youbin Wu, Yu Li, Yudong Liu, Yue Ling, Yujia Qin, Zhanbo Wang, Zhiwu He, Aoxue Zhang,

568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625

- 626 Bairen Yi, Bencheng Liao, Can Huang, Can Zhang, Chaorui
627 Deng, Chaoyi Deng, Cheng Lin, Cheng Yuan, Chenggang
628 Li, Chenhui Gou, Chenwei Lou, Chengzhi Wei, Chundian
629 Liu, Chunyuan Li, Deyao Zhu, Donghong Zhong, Feng Li,
630 Feng Zhang, Gang Wu, Guodong Li, Guohong Xiao, Haibin
631 Lin, Haihua Yang, Haoming Wang, Heng Ji, Hongxiang
632 Hao, Hui Shen, Huixia Li, Jiahao Li, Jialong Wu, Jianhua
633 Zhu, Jianpeng Jiao, Jiashi Feng, Jiase Chen, Jianhui Duan,
634 Jihao Liu, Jin Zeng, Jingqun Tang, Jingyu Sun, Joya Chen,
635 Jun Long, Junda Feng, Junfeng Zhan, Junjie Fang, Junting
636 Lu, Kai Hua, Kai Liu, Kai Shen, Kaiyuan Zhang, Ke Shen,
637 Ke Wang, Keyu Pan, Kun Zhang, Kunchang Li, Lanxin Li,
638 Lei Li, Lei Shi, Li Han, Liang Xiang, Liangqiang Chen, Lin
639 Chen, Lin Li, Lin Yan, Liying Chi, Longxiang Liu, Mengfei
640 Du, Mingxuan Wang, Ningxin Pan, Peibin Chen, Pengfei
641 Chen, Pengfei Wu, Qingqing Yuan, Qingyao Shuai, Qiuyan
642 Tao, Renjie Zheng, Renrui Zhang, Ru Zhang, Rui Wang, Rui
643 Yang, Rui Zhao, Shaoqiang Xu, Shihao Liang, Shipeng Yan,
644 Shu Zhong, Shuaishuai Cao, Shuangzhi Wu, Shufan Liu,
645 Shuhan Chang, Songhua Cai, Tenglong Ao, Tianhao Yang,
646 Tingting Zhang, Wanjun Zhong, Wei Jia, Wei Weng, Weihao
647 Yu, Wenhao Huang, Wenjia Zhu, Wenli Yang, Wenzhi Wang,
648 Xiang Long, XiangRui Yin, Xiao Li, Xiaolei Zhu, Xiaoy-
649 ing Jia, Xijin Zhang, Xin Liu, Xincheng Zhang, Xinyu Yang,
650 Xiongcai Luo, Xiuli Chen, Xuantong Zhong, Xuefeng Xiao,
651 Xujing Li, Yan Wu, Yawei Wen, Yifan Du, Yihao Zhang,
652 Yining Ye, Yonghui Wu, Yu Liu, Yu Yue, Yufeng Zhou,
653 Yufeng Yuan, Yuhang Xu, Yuhong Yang, Yun Zhang, Yun-
654 hao Fang, Yuntao Li, Yurui Ren, Yuwen Xiong, Zehua Hong,
655 Zehua Wang, Zewei Sun, Zeyu Wang, Zhao Cai, Zhaoyue
656 Zha, Zhecheng An, Zhehui Zhao, Zhengzhuo Xu, Zhipeng
657 Chen, Zhiyong Wu, Zhuofan Zheng, Zihao Wang, Zilong
658 Huang, Ziyu Zhu, and Zuquan Song. Seed1.5-v1 technical
659 report, 2025. 1
- [15] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perel-
660 man, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda,
661 Alan Hayes, Alec Radford, et al. Gpt-4o system card.
662 *arXiv preprint arXiv:2410.21276*, 2024. 6
- [16] Chaoya Jiang, Yongrui Heng, Wei Ye, Han Yang, Haiyang
663 Xu, Ming Yan, Ji Zhang, Fei Huang, and Shikun Zhang.
664 Vlm-r³: Region recognition, reasoning, and refinement for
665 enhanced multimodal chain-of-thought, 2025. 2
- [17] Xin Lai, Junyi Li, Wei Li, Tao Liu, Tianjian Li, and Heng-
666 shuang Zhao. Mini-o3: Scaling up reasoning patterns and
667 interaction turns for visual search, 2025. 2
- [18] Bangzheng Li, Ximeng Sun, Jiang Liu, Ze Wang, Jialian Wu,
668 Xiaodong Yu, Hao Chen, Emad Barsoum, Muhao Chen, and
669 Zicheng Liu. Latent visual reasoning, 2025. 2
- [19] Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhang
670 Zhu, Haihan Huang, Jianfei Gao, Kunchang Li, Yinan He,
671 Chenting Wang, et al. Videochat-flash: Hierarchical com-
672 pression for long-context video modeling. *arXiv preprint*
673 *arXiv:2501.00574*, 2024. 1, 2
- [20] Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu
674 Zeng, Yinan He, Yali Wang, Yu Qiao, Yi Wang, and
675 Limin Wang. Videochat-r1: Enhancing spatio-temporal
676 perception via reinforcement fine-tuning. *arXiv preprint*
677 *arXiv:2504.06958*, 2025. 2
- [21] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An
678 image is worth 2 tokens in large language models. In *ECCV*,
679 pages 323–340. Springer, 2024. 1, 2
- [22] Ziyu Liu, Yuhang Zang, Yushan Zou, Zijian Liang, Xiaoyi
680 Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi
681 Wang. Visual agentic reinforcement fine-tuning, 2025. 2
- [23] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra
682 Malik. Egoschema: A diagnostic benchmark for very long-
683 form video language understanding. *Advances in Neural In-*
684 *formation Processing Systems*, 36:46212–46244, 2023. 7
- [24] Damiano Marsili, Rohun Agrawal, Yisong Yue, and Geor-
685 gia Gkioxari. Visual agentic ai for spatial reasoning with a
686 dynamic api, 2025. 2
- [25] Meta AI. The llama 3.2 herd: Multilingual & multimodal
687 llms for edge and vision. [https://arxiv.org/abs/](https://arxiv.org/abs/2407.21783)
688 [2407.21783](https://arxiv.org/abs/2407.21783), 2024. 6
- [26] OpenAI. Openai o3 and o4-mini. [https://openai.](https://openai.com/index/introducing-o3-and-o4-mini/)
689 [com/index/introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/),
690 2025. 1
- [27] Yingzhe Peng, Gongrui Zhang, Miaosen Zhang, Zhiyuan
691 You, Jie Liu, Qipeng Zhu, Kai Yang, Xingzhong Xu, Xin
692 Geng, and Xu Yang. Lmm-r1: Empowering 3b llms with
693 strong reasoning abilities through two-stage rule-based rl,
694 2025. 2
- [28] Viorica Pătrăucean, Lucas Smaira, Ankush Gupta, Adrià Re-
695 casens Contintente, Larisa Markeeva, Dylan Banarse, Skanda
696 Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang,
697 Carl Doersch, Tatiana Matejovicova, Yury Sulsky, Antoine
698 Miech, Alex Frechette, Hanna Klimczak, Raphael Koster,
699 Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osin-
700 dero, Dima Damen, Andrew Zisserman, and João Carreira.
701 Perception test: A diagnostic benchmark for multimodal
702 video models. In *Advances in Neural Information Process-*
703 *ing Systems*, 2023. 14
- [29] Ziyao Shangquan, Chuhan Li, Yuxuan Ding, Yanan Zheng,
704 Yilun Zhao, Tesca Fitzgerald, and Arman Cohan. Tomato:
705 Assessing visual temporal reasoning capabilities in multi-
706 modal foundation models, 2025. 7
- [30] Wei Shen, Jiangbo Pei, Yi Peng, Xuchen Song, Yang Liu,
707 Jian Peng, Haofeng Sun, Yunzhuo Hao, Peiyu Wang, Jian-
708 hao Zhang, and Yahui Zhou. Skywork-r1v3 technical report,
709 2025. 1
- [31] Yan Shu, Peitian Zhang, Zheng Liu, Minghao Qin, Junjie
710 Zhou, Tiejun Huang, and Bo Zhao. Video-xl: Extra-long
711 vision language model for hour-scale video understanding.
712 *arXiv preprint arXiv:2409.14485*, 2024. 1, 2
- [32] Huajie Tan, Yuheng Ji, Xiaoshuai Hao, Xiansheng Chen,
713 Pengwei Wang, Zhongyuan Wang, and Shanghang Zhang.
714 Reason-rft: Reinforcement fine-tuning for visual reasoning
715 of vision language models, 2025. 2
- [33] Core Team, Zihao Yue, Zhenru Lin, Yifan Song, Weikun
716 Wang, Shuhuai Ren, Shuhao Gu, Shicheng Li, Peidian
717 Li, Liang Zhao, Lei Li, Kainan Bao, Hao Tian, Hailin
718 Zhang, Gang Wang, Dawei Zhu, Cici, Chenhong He,
719 Bowen Ye, Bowen Shen, Zihan Zhang, Zihan Jiang, Zhix-
720 ian Zheng, Zhichao Song, Zhenbo Luo, Yue Yu, Yudong
721 Wang, Yuanyuan Tian, Yu Tu, Yihan Yan, Yi Huang, Xu
722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740

- 741 Wang, Xinzhe Xu, Xingchen Song, Xing Zhang, Xing Yong,
742 Xin Zhang, Xiangwei Deng, Wenyu Yang, Wenhan Ma, Wei-
743 wei Lv, Weiji Zhuang, Wei Liu, Sirui Deng, Shuo Liu, Shi-
744 mao Chen, Shihua Yu, Shaohui Liu, Shande Wang, Rui Ma,
745 Qiantong Wang, Peng Wang, Nuo Chen, Menghang Zhu,
746 Kangyang Zhou, Kang Zhou, Kai Fang, Jun Shi, Jinhao
747 Dong, Jiebao Xiao, Jiaming Xu, Huaqiu Liu, Hongshen Xu,
748 Heng Qu, Haochen Zhao, Hanglong Lv, Guoan Wang, Duo
749 Zhang, Dong Zhang, Di Zhang, Chong Ma, Chang Liu, Can
750 Cai, and Bingquan Xia. MIMO-v1 technical report, 2025. 6,
751 7
- [34] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya
752 Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Ta-
753 tiana Matejovicova, Alexandre Ramé, Morgane Rivière,
754 et al. Gemma 3 technical report. *arXiv preprint*
755 *arXiv:2503.19786*, 2025. 6
756
- [35] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao
757 Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun
758 Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding,
759 Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du,
760 Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao,
761 Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu,
762 Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xi-
763 aoru Hao, Tianhong He, Weiran He, Wenyang He, Chao
764 Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi
765 Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin,
766 Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang
767 Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei
768 Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin,
769 Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu,
770 Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei
771 Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu,
772 Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling
773 Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei,
774 Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin,
775 Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan
776 Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung,
777 Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu
778 Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiax-
779 ing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang,
780 Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi
781 Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu
782 Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu,
783 Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing
784 Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu,
785 Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi
786 Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang,
787 Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye,
788 Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hong-
789 bang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang,
790 Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang,
791 Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang,
792 Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao,
793 Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou,
794 Zaida Zhou, Zhen Zhu, Weiyou Zhuang, and Xinxing Zu.
795 Kimi k2: Open agentic intelligence, 2025. 2
796
- [36] Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen
797 Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chen-
798 zhuang Du, Chu Wei, Congcong Wang, Dehao Zhang,
799 Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang
800 Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao
801 Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Hao-
802 tian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin
803 Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng,
804 Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Jun-
805 jie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu,
806 Mengfan Dong, Mengnan Dong, Nuo Xu, Pengyu Cheng,
807 Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu,
808 Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao
809 Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe
810 Wu, Xinhao Li, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y.
811 Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen,
812 Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen,
813 Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu,
814 Yuhao Dong, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan,
815 Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin
816 Yang, Zhiqi Huang, Zihao Huang, Zijia Zhao, Ziwei Chen,
817 and Zongyu Lin. Kimi-v1 technical report, 2025. 1
818
- [37] V Team, Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo
819 Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi,
820 Junhui Ji, Lihang Pan, Shuaiqi Duan, Weihang Wang, Yan
821 Wang, Yean Cheng, Zehai He, Zhe Su, Zhen Yang, Ziyang
822 Pan, Aohan Zeng, Baoxu Wang, Bin Chen, Boyan Shi,
823 Changyu Pang, Chenhui Zhang, Da Yin, Fan Yang, Guoqing
824 Chen, Jiazheng Xu, Jiale Zhu, Jiali Chen, Jing Chen, Jin-
825 hao Chen, Jinghao Lin, Jinjiang Wang, Junjie Chen, Leqi
826 Lei, Letian Gong, Leyi Pan, Mingdao Liu, Mingde Xu,
827 Mingzhi Zhang, Qinkai Zheng, Sheng Yang, Shi Zhong,
828 Shiyu Huang, Shuyuan Zhao, Siyan Xue, Shangqin Tu,
829 Shengbiao Meng, Tianshu Zhang, Tianwei Luo, Tianxiang
830 Hao, Tianyu Tong, Wenkai Li, Wei Jia, Xiao Liu, Xiaoh-
831 an Zhang, Xin Lyu, Xinyue Fan, Xuancheng Huang, Yan-
832 ling Wang, Yadong Xue, Yanfeng Wang, Yanzi Wang, Yi-
833 fan An, Yifan Du, Yiming Shi, Yiheng Huang, Yilin Niu,
834 Yuan Wang, Yuanchang Yue, Yuchen Li, Yutao Zhang, Yut-
835 ing Wang, Yu Wang, Yuxuan Zhang, Zhao Xue, Zhenyu Hou,
836 Zhengxiao Du, Zihan Wang, Peng Zhang, Debing Liu, Bin
837 Xu, Juanzi Li, Minlie Huang, Yuxiao Dong, and Jie Tang.
838 glm-4.5v and glm-4.1v-thinking: Towards versatile multi-
839 modal reasoning with scalable reinforcement learning, 2025.
840 1
841
- [38] Haozhe Wang, Alex Su, Weiming Ren, Fangzhen Lin, and
842 Wenhua Chen. Pixel reasoner: Incentivizing pixel-space rea-
843 soning with curiosity-driven reinforcement learning, 2025. 2
844
- [39] Kangrui Wang, Pingyue Zhang, Zihan Wang, Yaning Gao,
845 Linjie Li, Qineng Wang, Hanyang Chen, Chi Wan, Yiping
846 Lu, Zhengyuan Yang, Lijuan Wang, Ranjay Krishna, Jiajun
847 Wu, Li Fei-Fei, Yejin Choi, and Manling Li. Vagen: Re-
848 inforcing world model reasoning for multi-turn vlm agents,
849 2025. 2
850
- [40] Qi Wang, Yanrui Yu, Ye Yuan, Rui Mao, and Tianfei Zhou.
851 Videorf: Incentivizing video reasoning capability in mlms
852 via reinforced fine-tuning. *arXiv preprint arXiv:2505.12434*,
853 2025. 2
854
- [41] Weiyun Wang, Zhangwei Gao, Lianjie Chen, Zhe Chen, Jin-
855 guo Zhu, Xiangyu Zhao, Yangzhou Liu, Yue Cao, Shenglong
856

- 857 Ye, Xizhou Zhu, Lewei Lu, Haodong Duan, Yu Qiao, Jifeng
858 Dai, and Wenhai Wang. Visualprm: An effective process
859 reward model for multimodal reasoning, 2025. 2
- 860 [42] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long
861 Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Sheng-
862 long Ye, Jie Shao, et al. Internvl3.5: Advancing open-source
863 multimodal models in versatility, reasoning, and efficiency.
864 *arXiv preprint arXiv:2508.18265*, 2025. 1
- 865 [43] Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-
866 Levy. Videoagent: Long-form video understanding with
867 large language model as agent. In *European Conference on*
868 *Computer Vision*, pages 58–76. Springer, 2024. 1, 2
- 869 [44] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yanan He,
870 Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong
871 Shi, et al. Internvideo2: Scaling foundation models for mul-
872 timodal video understanding. In *European Conference on*
873 *Computer Vision*, pages 396–416. Springer, 2024. 1
- 874 [45] Yuxuan Wang, Yueqian Wang, Dongyan Zhao, Cihang Xie,
875 and Zilong Zheng. Videohalluc: Evaluating intrinsic
876 and extrinsic hallucinations in large video-language models,
877 2024. 8
- 878 [46] Yi Wang, Xinhao Li, Ziang Yan, Yanan He, Jiashuo Yu, Xi-
879 angyu Zeng, Chenting Wang, Changlian Ma, Haiyan Huang,
880 Jianfei Gao, et al. Internvideo2. 5: Empowering video
881 mllms with long and rich context modeling. *arXiv preprint*
882 *arXiv:2501.12386*, 2025. 1, 2
- 883 [47] Zikang Wang, Boyu Chen, Zhengrong Yue, Yi Wang, Yu
884 Qiao, Limin Wang, and Yali Wang. Videochat-a1: Thinking
885 with long videos by chain-of-shot reasoning. *arXiv preprint*
886 *arXiv:2506.06097*, 2025. 1, 2
- 887 [48] Zhenhailong Wang, Xuehang Guo, Sofia Stoica, Haiyang
888 Xu, Hongru Wang, Hyeonjeong Ha, Xiusi Chen, Yangyi
889 Chen, Ming Yan, Fei Huang, and Heng Ji. Perception-aware
890 policy optimization for multimodal reasoning, 2025. 2
- 891 [49] Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong
892 Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal.
893 Videotree: Adaptive tree-based video representation for llm
894 reasoning on long videos. In *Proceedings of the Computer*
895 *Vision and Pattern Recognition Conference*, pages 3272–
896 3283, 2025. 1, 2
- 897 [50] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li.
898 Longvideobench: A benchmark for long-context interleaved
899 video-language understanding, 2024. 7
- 900 [51] Jiaer Xia, Yuhang Zang, Peng Gao, Sharon Li, and Kaiyang
901 Zhou. Visionary-r1: Mitigating shortcuts in visual reasoning
902 with reinforcement learning, 2025. 2
- 903 [52] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua.
904 Next-qa: Next phase of question-answering to explaining
905 temporal actions. In *Proceedings of the IEEE/CVF Confer-*
906 *ence on Computer Vision and Pattern Recognition (CVPR)*,
907 pages 9777–9786, 2021. 14
- 908 [53] Guowei Xu, Peng Jin, Ziang Wu, Hao Li, Yibing Song,
909 Lichao Sun, and Li Yuan. Llava-cot: Let vision language
910 models reason step-by-step, 2025. 2
- 911 [54] Yi Xu, Chengzu Li, Han Zhou, Xingchen Wan, Caiqi Zhang,
912 Anna Korhonen, and Ivan Vulić. Visual planning: Let’s think
913 only with images, 2025. 2
- [55] Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu,
Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian Tang, Shang
Yang, Zhijian Liu, et al. Longvila: Scaling long-context
visual language models for long videos. *arXiv preprint*
arXiv:2408.10188, 2024. 1, 2
- [56] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen
Huang, Chenxu Lv, Chuji Zheng, Dayiheng Liu, Fan Zhou,
Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jia-
long Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai
Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei
Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin
Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tian-
hao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu
Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su,
Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun
Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
Qiu. Qwen3 technical report, 2025. 2, 6, 7
- [57] Biao Yang, Bin Wen, Boyang Ding, Changyi Liu, Chen-
glong Chu, Chengru Song, Chongling Rao, Chuan Yi, Da
Li, Dunju Zang, Fan Yang, Guorui Zhou, Guowang Zhang,
Han Shen, Hao Peng, Haojie Ding, Hao Wang, Haonan Fan,
Hengrui Ju, Jiaming Huang, Jiangxia Cao, Jiankang Chen,
Jingyun Hua, Kaibing Chen, Kaiyu Jiang, Kaiyu Tang, Kun
Gai, Muhao Wei, Qiang Wang, Ruitao Wang, Sen Na, Sheng-
nan Zhang, Siyang Mao, Sui Huang, Tianke Zhang, Tingting
Gao, Wei Chen, Wei Yuan, Xiangyu Wu, Xiao Hu, Xingyu
Lu, Yi-Fan Zhang, Yiping Yang, Yulong Chen, Zeyi Lu,
Zhenhua Wu, Zhixin Ling, Zhuoran Yang, Ziming Li, Di
Xu, Haixuan Gao, Hang Li, Jing Wang, Lejian Ren, Qigen
Hu, Qianqian Wang, Shiyao Wang, Xinchun Luo, Yan Li,
Yuhang Hu, and Zixing Zhang. Kwai keye-v1 1.5 technical
report, 2025. 1, 2
- [58] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han,
Li Fei-Fei, and Saining Xie. Thinking in space: How mul-
timodal large language models see, remember, and recall
spaces, 2025. 7
- [59] Zeyuan Yang, Xueyang Yu, Delin Chen, Maohao Shen, and
Chuang Gan. Machine mental imagery: Empower multi-
modal reasoning with latent visual tokens, 2025. 2
- [60] Linli Yao, Haoning Wu, Kun Ouyang, Yuanxing Zhang,
Caiming Xiong, Bei Chen, Xu Sun, and Junnan Li. Gen-
erative frame sampler for long video understanding. *arXiv*
preprint arXiv:2503.09146, 2025. 1, 2
- [61] Jinhui Ye, Zihan Wang, Haosen Sun, Keshigeyan Chan-
drasegaran, Zane Durante, Cristobal Eyzaguirre, Yonatan
Bisk, Juan Carlos Niebles, Ehsan Adeli, Li Fei-Fei, et al. Re-
thinking temporal search for long-form video understanding.
In *Proceedings of the Computer Vision and Pattern Recogni-*
tion Conference, pages 8579–8591, 2025. 1, 2
- [62] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Ji-
ajun Wu, Antonio Torralba, and Joshua B. Tenenbaum.
Clevrer: Collision events for video representation and rea-
soning, 2020. 14
- [63] En Yu, Kangheng Lin, Liang Zhao, Jisheng Yin, Yana Wei,
Yuang Peng, Haoran Wei, Jianjian Sun, Chunrui Han, Zheng

- 971 Ge, Xiangyu Zhang, Daxin Jiang, Jingyu Wang, and Wen-
972 bing Tao. Perception-rl: Pioneering perception policy with
973 reinforcement learning, 2025. [2](#)
- 974 [64] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng,
975 Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan,
976 Chunyuan Li, and Ziwei Liu. Long context transfer from
977 language to vision. *CoRR*, abs/2406.16852, 2024. [1](#), [2](#)
- 978 [65] Xintong Zhang, Zhi Gao, Bofei Zhang, Pengxiang Li, Xi-
979 aowen Zhang, Yang Liu, Tao Yuan, Yuwei Wu, Yunde Jia,
980 Song-Chun Zhu, and Qing Li. Chain-of-focus: Adaptive
981 visual search and zooming for multimodal reasoning via rl,
982 2025. [2](#)
- 983 [66] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Zi-
984 wei Liu, and Chunyuan Li. Video instruction tuning with
985 synthetic data, 2024. [14](#)
- 986 [67] Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao,
987 Guohai Xu, Le Yang, Chao Shen, and Xing Yu. Deep-
988 eyes: Incentivizing "thinking with images" via reinforce-
989 ment learning, 2025. [2](#)
- 990 [68] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang
991 Liang, Shitao Xiao, Minghao Qin, Xi Yang, Yongping
992 Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu:
993 Benchmarking multi-task long video understanding, 2025. [7](#)
- 994 [69] Orr Zohar, Xiaohan Wang, Yonatan Bitton, Idan Szpektor,
995 and Serena Yeung-levy. Video-star: Self-training enables
996 video instruction tuning with any supervision. In *arXiv*
997 *preprint arXiv:2407.06189*, 2024. [14](#)

TreeReasoner: Reinforcing Tool-Augmented Tree-of-Videos Reasoning

Supplementary Material

998 A. Case Study

999 In the introduction, we posited that our agentic frame-
1000 work enhances reasoning by actively navigating the video
1001 timeline—a departure from passive frame processing—and
1002 highlighted the interpretability of its search strategies. In
1003 this section, we provide qualitative examples (Fig. 7 and
1004 8) to substantiate these claims, specifically demonstrating
1005 the model’s search behaviors, such as temporal bracketing,
1006 causal chain verification, and adaptive temporal resolution
1007 adjustment.

1008 **Temporal bracketing.** As shown in Fig. 7a, consider
1009 a query requiring multi-step reasoning, such as “Identify
1010 the color of the pants worn by the young woman talk-
1011 ing to an elderly man on the subway.” Initially, TreeRea-
1012 soner employs a coarse-grained search with a large tem-
1013 poral stride to rapidly scan the entire timeline, identifying
1014 broad temporal regions where key objects (e.g., subway
1015 scenes, elderly man) appear. Upon detecting these rele-
1016 vant contexts, the model automatically switches to a fine-
1017 grained temporal mode, progressively narrowing down the
1018 time window by examining adjacent segments. This al-
1019 lows it to precisely bracket the specific temporal span where
1020 the target event—“a young woman conversing with an el-
1021 derly man”—occurs. Within this bracketed time window,
1022 the model then analyzes the fine-grained visual details to
1023 extract the answer. By “bracketing” the event temporally
1024 in this coarse-to-fine manner, the model avoids processing
1025 the redundant background footage, directly contributing to
1026 the efficiency gains mentioned in our experiments. Another
1027 similar case can be found in Fig. 8a.

1028 **Causal chain verification.** For complex queries such as
1029 “Find the color of the pants worn by the young girl who
1030 talks to the old man on the subway”, the model exhibits a
1031 hierarchical reasoning strategy. TreeReasoner first searches
1032 for video segments containing the objects mentioned in the
1033 query, including segments with subway scenes, old men,
1034 and young girls. After identifying segments with the speci-
1035 fied objects, TreeReasoner further examines the consistency
1036 between video content and the query-specified events, ulti-
1037 mately pinpointing the segment where “the young girl talks
1038 to the old man” occurs. It then comprehends the object in-
1039 formation in this segment and provides the answer. This
1040 search strategy enables TreeReasoner to retrieve query-
1041 relevant video segments at a relatively fine-grained level,
1042 forming a complete chain of visual evidence and delivering
1043 accurate responses grounded in fine-grained visual details.
1044 This behavior demonstrates that TreeReasoner validates re-

1045 lationships by progressively refining its search through mul-
1046 tiple verification steps, ensuring the answer is grounded in
1047 visual evidence rather than hallucination. Another similar
1048 case can be found in Fig. 8b.

1049 **Adaptive temporal resolution adjustment.** In the final
1050 reasoning tree, TreeReasoner exhibits variable temporal res-
1051 olutions across different video segments through adaptive
1052 tool invocation. As shown in Fig. 7c, for static, repeti-
1053 tive, or query-irrelevant portions of the video, TreeReasoner
1054 dynamically increases its temporal stride, skipping large
1055 chunks of uninformative frames. Conversely, for segments
1056 closely related to the query—such as those containing the
1057 target individuals and their interaction—the model immedi-
1058 ately invokes tools to extract short video clips and applies
1059 much smaller step sizes to capture high-frequency details
1060 of the girl’s appearance and clothing. This adaptive mech-
1061 anism explains how our method achieves superior accuracy
1062 while processing substantially fewer frames than fixed-rate
1063 baselines. Another similar case can be found in Fig. 8c.

B. Dataset Details 1064

1065 To endow TreeReasoner with robust tool-use capabilities
1066 and ensure generalization across diverse visual domains, we
1067 constructed a unified training corpus by aggregating queries
1068 from multiple open-source datasets, including CLEVRER
1069 [62], LLaVA-Video-178K [66], NEX-T-QA [52], Percep-
1070 tionTest [28], Video-STaR [69], and LongVideo-Reason
1071 [4].

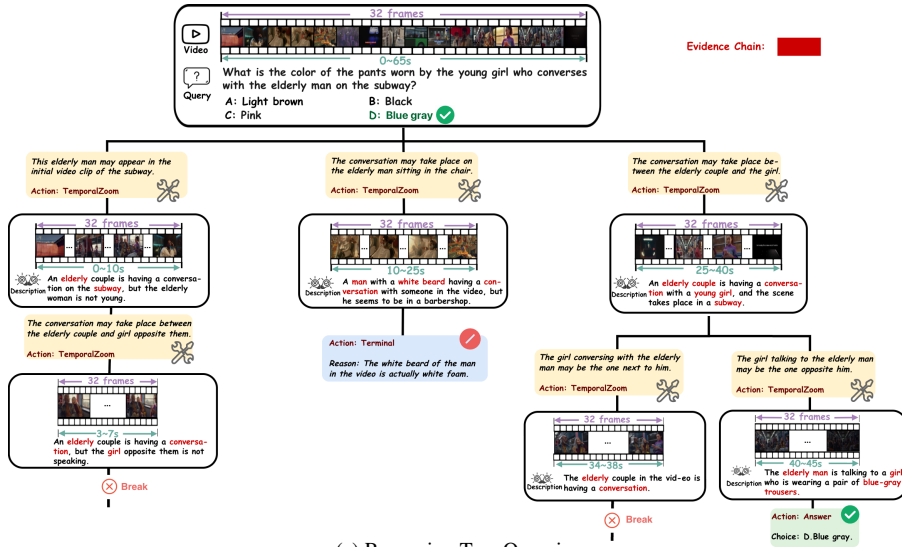
1072 **SFT Data Construction.** From this broad collection, we
1073 first sampled a balanced subset of 20,000 video-question
1074 pairs to serve as the source for Supervised Fine-Tuning
1075 (SFT). Since these datasets lack explicit annotations for
1076 tool-augmented reasoning chains, we employed a knowl-
1077 edge distillation strategy utilizing Gemini-2.5-Pro [13] as
1078 the teacher model. Specifically, we prompted the teacher to
1079 generate multi-branch Tree-of-Tool reasoning paths based
1080 on our defined toolset (zoom, jump, slide). To align with
1081 our framework’s objective of identifying a *minimal yet suf-*
1082 *ficient* chain of evidence, we applied an *efficiency-oriented*
1083 *filtering strategy*. For each query, we ranked the success-
1084 ful reasoning paths by length and retained up to the top- k
1085 ($k = 2$) shortest trajectories. Additionally, we explicitly
1086 preserved a subset of trajectories where the model actively
1087 triggered termination to ensure robust exploration capabili-
1088 ties. This process resulted in a final corpus of 33,724 high-
1089 quality trajectories for behavioral cloning.

1090 **RL Data Construction.** Distinct from the SFT dataset, we

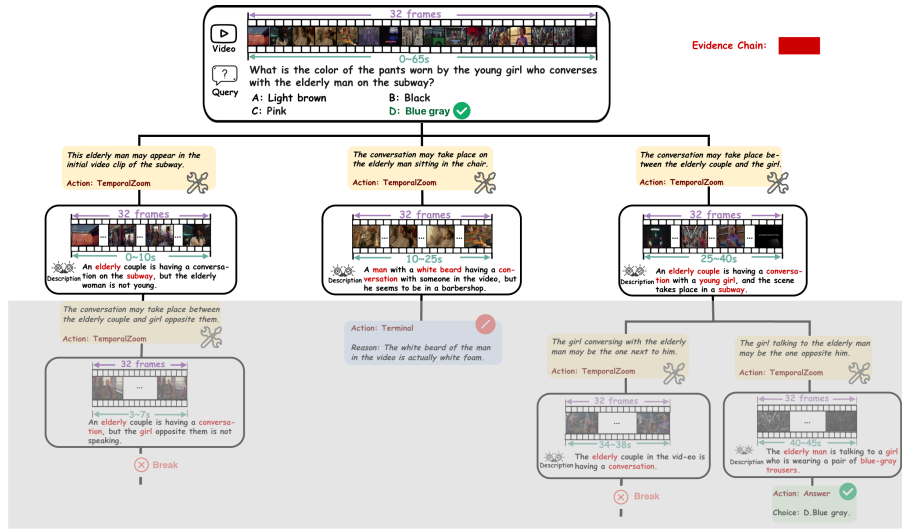
1091 sampled an *additional, non-overlapping subset* of 40,000
1092 video-question pairs from the aggregated corpus to drive the
1093 RL stage. By using raw queries instead of teacher trajectories,
1094 the model generates its own reasoning paths optimized
1095 via ToT-RPO, balancing accuracy with efficiency and en-
1096 hancing generalization.

1097 **C. Limitations**

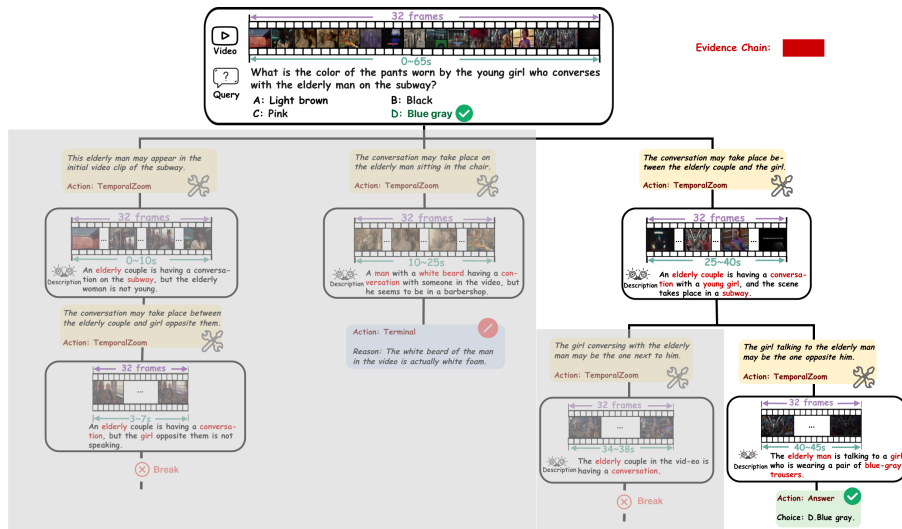
1098 Despite the demonstrated effectiveness of TreeReasoner in
1099 handling complex video reasoning tasks, we acknowledge
1100 several limitations. First, the reliance on a tree search mech-
1101 anism introduces a computational overhead compared to
1102 standard end-to-end IO of MLLMs. While our Reinforce-
1103 ment Learning strategy explicitly optimizes for shorter rea-
1104 soning paths, the inference process still requires multiple
1105 forward passes to explore and prune branches, which may
1106 hinder deployment in strictly real-time scenarios. Second,
1107 our current toolset is limited to visual navigation primitives
1108 (zoom, jump, slide). While these are sufficient for spa-
1109 tiotemporal grounding, the model lacks tools for processing
1110 other modalities (e.g., audio analysis). Finally, our training
1111 paradigm depends on knowledge distillation from a teacher
1112 model (Gemini-2.5-Pro). Although we apply rigorous fil-
1113 tering to the distilled data, the student model’s upper bound
1114 is inevitably influenced by the teacher’s capabilities.



(a) Reasoning Tree Overview

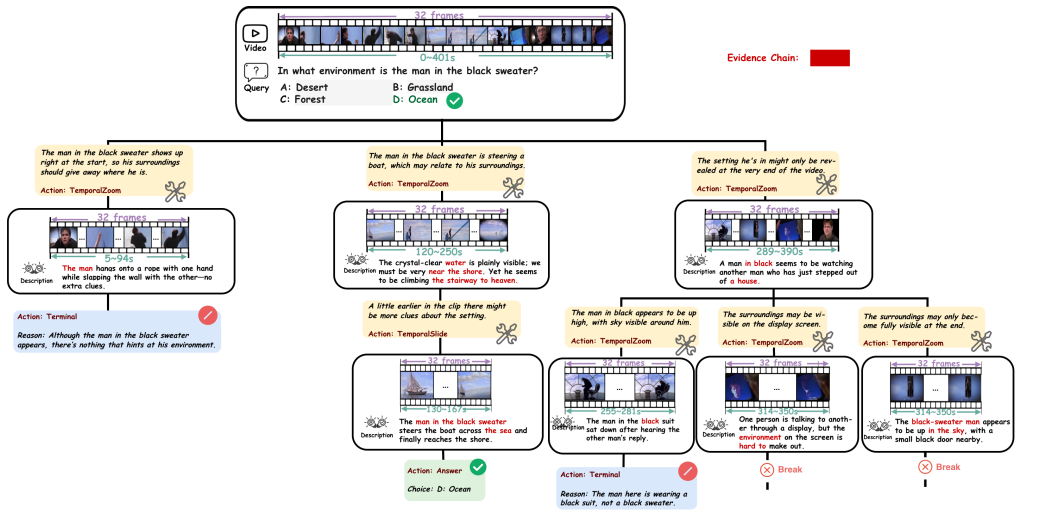


(b) Temporal Bracketing

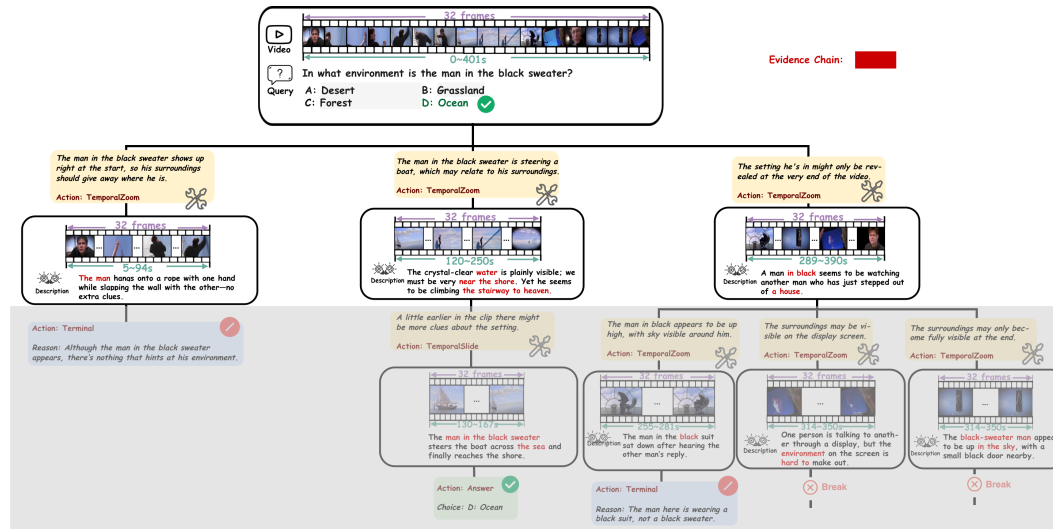


(c) Causal Chain Verification

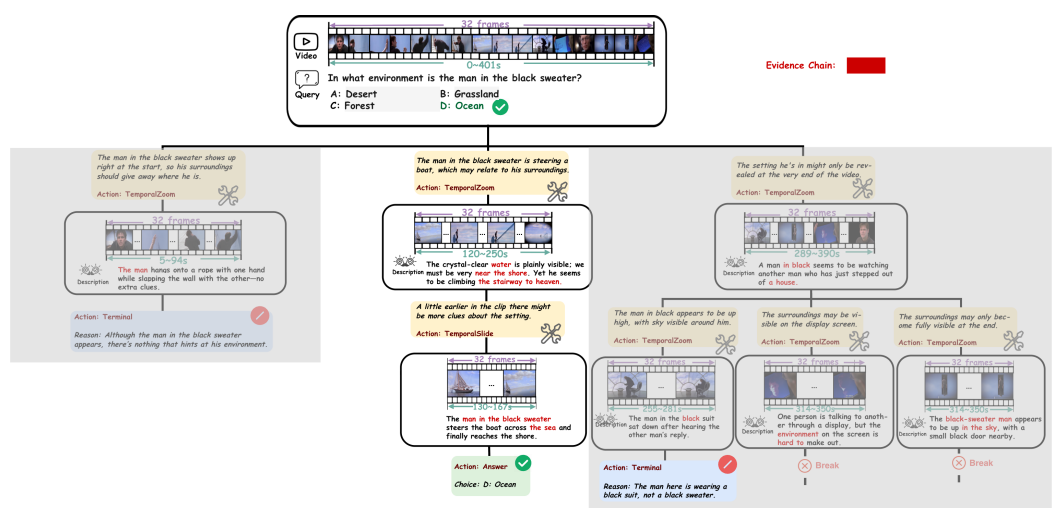
Figure 7. A case study on "identifying a girl's pants color query in subway scene": (a) shows an overview of the complete reasoning tree; (b) illustrates the temporal bracketing mechanism; (c) presents a branch trajectory along with its corresponding visual evidence chain. Irrelevant regions are masked in (b) and (c) to highlight key content.



(a) Reasoning Tree Overview



(b) Temporal Bracketing



(c) Causal Chain Verification

Figure 8. A case study on “the environment query of a man in black sweater”: (a) shows an overview of the complete reasoning tree; (b) illustrates the temporal bracketing mechanism; (c) presents a branch trajectory along with its corresponding visual evidence chain. Irrelevant regions are masked in (b) and (c) to highlight key content.