

UNLOCKING VIDEO-LLM VIA AGENT-OF-THOUGHTS DISTILLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper tackles the problem of video question answering (VideoQA), a task that often requires multi-step reasoning and a profound understanding of spatial-temporal dynamics. While large generative video-language models perform well on benchmarks, they often lack explainability and spatial-temporal grounding. In this paper, we propose **Agent-of-Thoughts Distillation (AoTD)**, a method that enhances generative models by incorporating automatically generated Chain-of-Thoughts (CoTs) into the instruction-tuning process. Specifically, we leverage an agent-based system to decompose complex questions into sub-tasks, and address them with specialized vision models, the intermediate results are then treated as reasoning chains. We also introduce a verification mechanism using a large language model (LLM) to ensure the reliability of generated CoTs. Extensive experiments demonstrate that AoTD improves the performance on multiple-choice and open-ended benchmarks.

1 INTRODUCTION

Video Question Answering (VideoQA) is a critical task in the computer vision community, offering a natural interface for human-machine interaction through language (Yu et al., 2019; Wu et al., 2021; Xiao et al., 2021; Pătrăucean et al., 2023). This synergy of visual content and language enhances the accessibility of AI systems for the general public, allowing users to query complex visual content with everyday language. By encompassing tasks such as action recognition, object detection, and scene understanding, VideoQA serves as a comprehensive benchmark for evaluating AI’s ability to interpret videos, addressing the fundamental questions of ‘who,’ ‘what,’ ‘when,’ and ‘where’ that are crucial to understand daily life activities, pushing the boundaries of what AI systems can interpret from dynamic visual content.

Recent literature in VideoQA has highlighted two key directions. The first focuses on training large generative video-language models (Video-LLMs) through direct instruction-tuning, where videos are only paired with questions and answers (Alayrac et al., 2022; Lin et al., 2024; Maaz et al., 2024; Cheng et al., 2023). While these models have shown exceptional performance on public benchmarks, they often lack explainability and struggle with spatio-temporal grounding. This limitation hinders their ability to provide clear reasoning, which is essential for real-world applications where transparency and interpretability are critical (Mitra et al., 2023).

In contrast, an emerging approach focuses on agent-based systems (Surís et al., 2023; Gupta & Kembhavi, 2023; Hu et al., 2024b), which break down complex questions into simpler sub-tasks. Each sub-task is handled by specialized tools, and the results are aggregated to generate a final answer. This approach theoretically offers greater interpretability, as the reasoning process is divided into explainable steps that can be independently assessed. However, our experiments indicate that current video understanding tools are not strong enough for building reliable agent-based systems. Additionally, the high memory demands and time-consuming nature of these systems present significant challenges for their practical use.

In this paper, we propose enhancing the capabilities of large generative video-language models by incorporating automatically generated Chain-of-Thoughts (CoTs) into the instruction-tuning process. Our approach draws inspiration from agent-based systems, which break down complex questions into a sequence of sub-tasks, each handled by specialized models (Fan et al., 2024; Mahmood et al., 2024; Min et al., 2024). We use the outputs from these specialized models to construct CoTs that

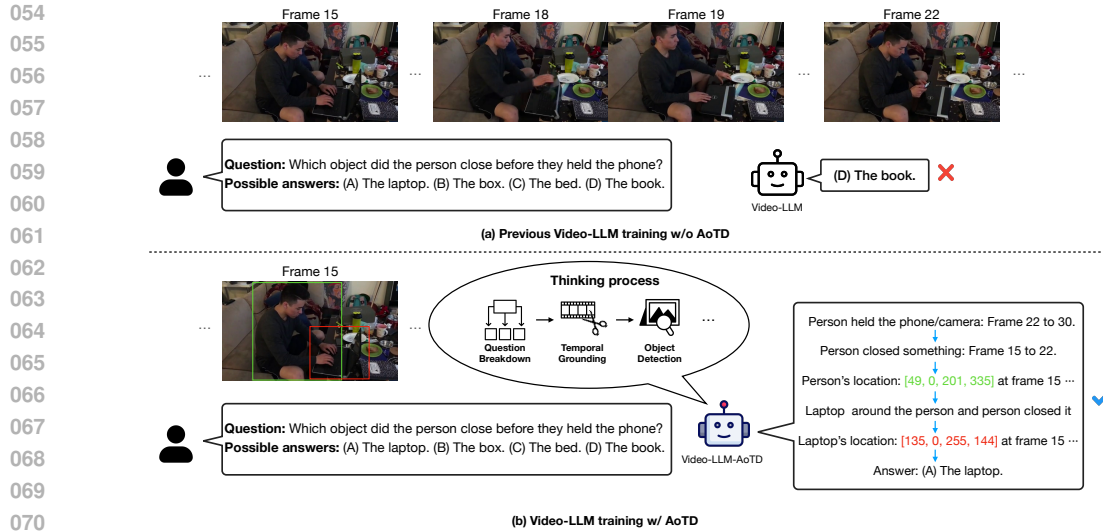


Figure 1: Our method, **AoTD**, distills multi-step reasoning and spatio-temporal understanding into a single generative video-language model. When addressing complex VideoQA tasks, the model trained with AoTD (as shown in (b)) enables to generate a step-by-step reasoning to get the correct answer. In contrast, previous models trained solely on question-answer pairs (as in (a)) generate only a final answer, often without intermediate reasoning, which can lead to incorrect conclusions.

explicitly represent step-by-step reasoning paths, capturing the reasoning processes that generative models typically struggle to model independently.

To ensure the reliability of the constructed CoTs, we systematically evaluate existing models and tools for atomic video understanding tasks, such as action recognition (Weng et al., 2023; Wang et al., 2024) and language grounding (Lin et al., 2023), using a well-annotated dataset. This allows us to identify the best-performing tools for each sub-task, preparing for effective CoTs distillation. This process also serves as an evaluation of the broader capabilities of visual models in more general and complex scenes, offering guidance for future exploration in the computer vision community. Additionally, we introduce a verification mechanism with a large language model (LLM), to assess whether the generated CoTs follow a clear, step-by-step reasoning process and contain useful information for answering the question. This filters out low-quality or logically inconsistent reasoning paths. The verified, high-quality CoTs are then distilled into large generative video-language models, enhancing both performance and the interpretability of their outputs. By combining the strengths of both approaches, our method balances performance with transparency, leading to the development of more robust, accurate, and interpretable VideoQA systems.

In summary, our contributions are three-fold: *First*, we propose a novel approach for enhancing large generative video-language models (Video-LLMs) by distilling high-quality Chain-of-Thoughts (CoTs) into their instruction tuning. These CoTs capture step-by-step reasoning paths, improving both the model’s performance and its interpretability; *Second*, to automatically construct the CoTs for any datasets, we employ an agent-based system to decompose complex VideoQA questions into simpler sub-tasks, leveraging off-the-shelf vision models to handle each sub-task. The intermediate outputs from these models can therefore be collected as CoTs for addressing the corresponding visual question; *Third*, we demonstrate through extensive experiments that our distilled model outperforms existing methods across both multiple-choice and open-ended VideoQA benchmarks, enabling to deliver not only accurate answers but also clear and comprehensive reasoning explanations.

2 AGENT-OF-THOUGHTS DISTILLATION

In this paper, we propose a novel approach, termed Agent-of-Thought Distillation (AoTD), to enhance the Video-LLMs by training them with multi-step chain-of-thoughts (CoTs). Specifically, we

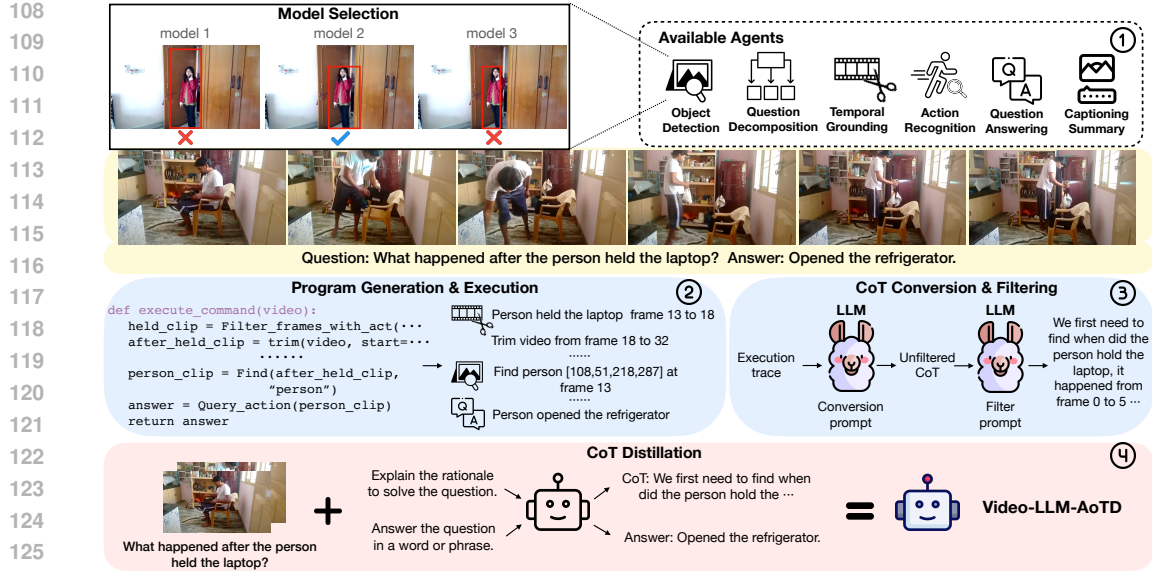


Figure 2: Overview on Agent-of-Thoughts Distillation (AoTD). **Step 1:** Selecting best-performing agents for each sub-task to construct an agent-based system. **Step 2:** Decomposing question into executable program and leveraging chosen models to solve it sequentially to generate execution trace. **Step 3:** The execution trace is converted and filtered by LLM to produce high quality natural language CoTs. **Step 4:** Distilling CoTs into Video-LLM with two forms of prompt, allowing it achieve a balance between concise answers and comprehensive rationales. The final model is Video-LLM-AoTD.

begin by developing an agent-based video understanding system to generate multi-step reasoning chains that address complex video questions. These reasoning chains are then distilled into one Video-LLM through instruction tuning. By combining the strengths of agent-based systems and large generative models, our proposed AoTD enables to build more reliable and interpretable Video Question Answering systems.

2.1 PROBLEM FORMULATION

Given a video clip with t frames, $\mathcal{V} = \{x_1, \dots, x_t\}$, and a set of n questions $\mathcal{Q} = \{q_1, q_2, \dots, q_n\}$, our goal is to train a Video-LLM capable of producing both concise answers and comprehensive rationales. Depending on the suffix prompt p_s , the model can generate different types of outputs. The process can be formulated as:

$$\{a_i, \mathcal{S}_i\} = \Phi(\mathcal{V}, q_i, p_s), \quad \text{where } \mathcal{S}_i = \{\emptyset\} \text{ or } \{s_{i,1}, s_{i,2}, \dots, s_{i,k}\}$$

where q_i denotes the i -th question, a_i is the answer in free-form text, and \mathcal{S}_i represents the rationale, consisting of the multi-step reasoning process. If the prompt specifies to only generate the answer, $\mathcal{S}_i = \{\emptyset\}$. Otherwise, if the prompt requires the generation of rationales, $\mathcal{S}_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,k}\}$, where each $s_{i,j}$ corresponds to a reasoning step.

Discussion. Unlike existing models that are typically instruction-tuned on VideoQA datasets using simple question-answer pairs, which bypass the intermediate thought process, our approach emphasizes the importance of training with rationales, or chain-of-thoughts (CoTs). In the following section, we outline the process for generating high-quality CoTs from existing VideoQA datasets.

2.2 CoTs CONSTRUCTION WITH AGENT-BASED SYSTEM

Recent work, such as STAR (Wu et al., 2021), has introduced executable symbolic programs that can directly decompose questions into sub-tasks. When combined with scene graphs that contain comprehensive video information from key frames—such as object locations, interactions, and actions—these programs facilitate the generation of concise Chain-of-Thoughts (CoTs) through the

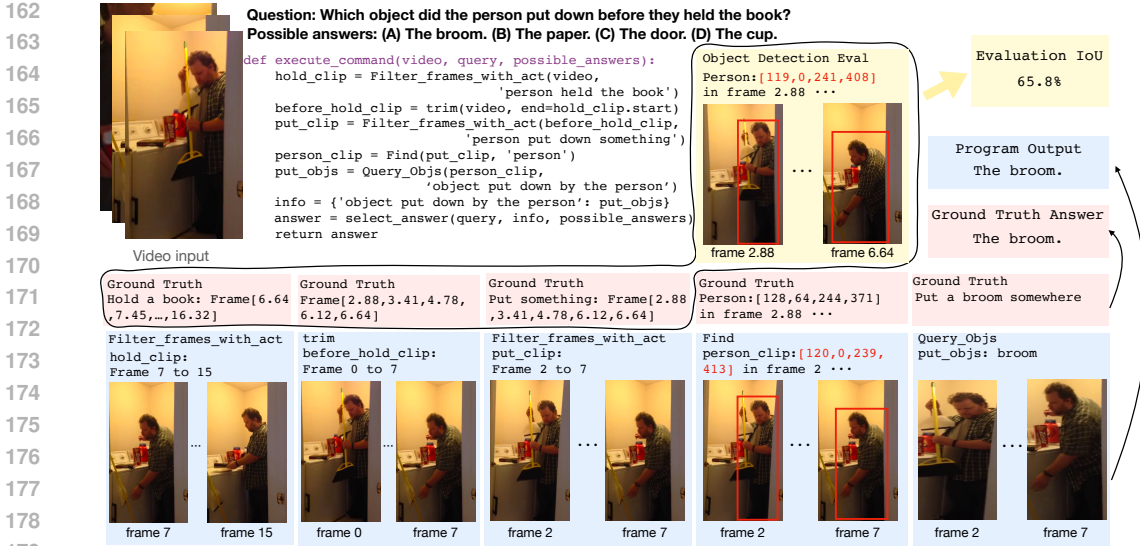


Figure 3: Program execution process in an agent-based system. We uniformly sample 32 frames from the video, and to ensure scale consistency, the frame ids of key frames are normalized into these 32 frames. The blue boxes represent the program execution steps, the red boxes denote the ground truth for each step.

direct execution of symbolic operations. However, datasets of this nature are limited in scale. In response to this limitation, we propose an agent-based system capable of breaking down complex questions into simpler sub-tasks, utilizing off-the-shelf vision models. The intermediate outputs from this system can then be employed to construct CoTs for any existing VideoQA dataset.

Agent-based VideoQA. Given a video input (\mathcal{V}), questions (\mathcal{Q}), and a set of visual models ($\mathcal{M} = \{\phi_{act}, \phi_{det}, \dots, \phi_{qa}\}$), an LLM-based agent core ($\pi(\cdot)$) processes the question along with the documentation of the visual models (\mathcal{T}), which includes variables and functionalities. The agent then decomposes the question into sub-tasks and addresses them by invoking the corresponding visual models. It is important to note that the visual models can be arranged in various orders depending on the specific question, ensuring flexibility in problem-solving.

Specifically, in the example illustrated in Fig. 3, the question is first decomposed into a series of sub-tasks, including temporal grounding, object detection, and question answering. The corresponding specialized models are then executed sequentially to address these sub-tasks, ultimately yielding the final answer:

$$\{\phi_{ground}, \phi_{det}, \phi_{qa}\} := \pi(q_i, \mathcal{T}), \quad y_i = \phi_{ground}(\mathcal{V}) \rightarrow \phi_{det}(\mathcal{V}) \rightarrow \phi_{qa}(\mathcal{V})$$

CoTs Construction. To ensure the correctness of outputs at intermediate steps, we leverage the training set from STAR for hyperparameter tuning, enabling us to identify the most effective model for each sub-task within the agent-based system. By following the provided programs, we evaluate the performance of the corresponding vision models on tasks such as object detection and action recognition. Given the availability of complete reasoning chains, we independently assess each sub-task using ground truth data for all preceding steps.

Table 1 presents the evaluation results for the various sub-tasks. For **question decomposition**, we compared several code LLMs, with DeepSeek-Coder-Instruct achieving the highest performance, outperforming even GPT-3.5-Turbo. In **object detection**, OWL-ViT v2 recorded the highest Intersection over Union (IoU) score, showcasing its superior open-vocabulary detection capability. The results for **temporal grounding** indicate that while UniVTG leads in performance, there remains a need for further advancements in this area. In **action recognition**, our evaluations showed that generative models outperformed discriminative models, likely due to the fine-grained action list provided by the STAR dataset. However, the performance of both model types reveals significant opportunities for improvement. Finally, in the **one-hop question answering** sub-task, all models

Table 1: Sub-tasks definition and evaluation results. We choose 3 model candidates for each sub-task and evaluate them in STAR training set with the corresponding metrics. Models with best performance are placed at the bottom of each column.

Sub-task name	Model name	Metric	Number (%)
Question decomposition	CodeQwen1.5-Chat (7B) (Bai et al., 2023)		52.7
	GPT-3.5-Turbo (OpenAI, 2023a)	Acc	73.1
	DeepSeek-Coder-Instruct (6.7B) (Daya et al., 2024)		85.7
Object detection	OWL-ViT v1 (Matthias et al., 2022)		47.3
	GLIP (Li* et al., 2022)	IoU	58.9
	OWL-ViT v2 (Minderer et al., 2024)		63.0
Temporal grounding	LITA (13B) (Huang et al., 2024)		11.7 / 20.2
	TimeChat (7B) (Ren et al., 2024)	IoU / Recall	13.9 / 23.1
	UniVTG (Lin et al., 2023)		24.7 / 35.3
Action recognition	InternVideo2 (1B) (Wang et al., 2024)		7.6
	Open-VCLIP (Weng et al., 2023)	Top1-Acc	8.9
	LLaVA-NeXT-Video-DPO (7B) (Zhang et al., 2024)		18.2
Question answering	LLaMA-VID (7B) (Li et al., 2024c)		43.5
	SeViLA (Yu et al., 2023a)	Acc	46.5
	LLaVA-NeXT-Video-DPO (7B) (Zhang et al., 2024)		53.4

performed admirably, with LLaVA-NeXT-Video-DPO standing out as a top performer, consistent with its strong results on other benchmarks.

With these high-performing models, we implement the agent-based approach on VideoQA datasets that consist solely of question-answer pairs. During the execution of the programs, we record all intermediate outputs to construct the CoTs. Since the outputs from these vision models vary in format—such as bounding boxes and free-form text—we employ another LLM to translate the execution trace into natural language, facilitating its use in the distillation process. Detailed examples are provided in Appendix B.

2.3 CoTs VERIFICATION

To further refine the quality of reasoning chains for VideoQA samples, we implement a two-step verification process: (i) we filter execution traces to retain only those where the program can reach correct output. For multiple-choice datasets, the output must match the correct answer exactly, while for open-ended datasets, we prompt the LLM to verify correctness, accounting for format differences; (ii) we prompt the LLM to evaluate the logical coherence and usefulness of the reasoning chains in solving the problem. The model assesses whether the CoTs follow a clear, step-by-step reasoning process and provides a binary evaluation (‘Yes’ or ‘No’) to indicate their quality (detailed prompts can be found in Appendix C). This two-step approach ensures that only accurate and high-quality CoTs are utilized for further distillation into the model.

After filtering, we provide the statistics for the generated CoTs on different datasets in Table 2. We primarily select compositional QA datasets, as these require the model to process spatial-temporal information from different events comprehensively.

2.4 DISTILL STEP BY STEP

In this section, we describe the process of distilling the generated CoTs into a Video-LLM. This distillation enhances the model’s ability for spatial-temporal video understanding and multi-step reasoning, thereby improving its performance on complex Video Question Answering (VideoQA) tasks.

Specifically, using the generated CoTs, we can build the dataset $D = \{(\mathcal{V}_j, q_j, \hat{y}_j, c_j, p_s)\}_{j=1}^N$, where N is the total number of samples in the distilling dataset, \mathcal{V}_j is the video input, q_j is the

Dataset	Description	# Labels	# CoTs
AGQA	Compositional	25.0K	5.4K
ANetQA	Compositional	25.0K	3.6K
STAR	Compositional	45.7K	11.2K
NExT-QA	Temporal & Causal	34.1K	12.1K
CLEVRER	Spatial & Temporal	21.0K	-
EgoQA	Ego-centric	7.8K	-
Total		158.6K	32.3K

Table 2: Dataset statistics. The column “# Labels” indicates the number of VideoQA pairs, which include the video, query, possible answers (multiple-choice), and the correct answer. “# CoTs” refers to the number of CoTs generated using our agent-based system for each dataset.

question, \hat{y}_j is the ground-truth answer, c_j is the generated CoT, p_s is the task-specific suffix prompt, to distinguish different tasks, for example, for multiple-choice VQA, the prompt is “Answer with the option’s letter from the given choices directly and only give the best option”, and for open-ended VQA, the prompt is “Answer in one word or phrase”. Detailed prompts are provided in Appendix C.

At distillation stage, we minimize the cross-entropy loss of predicting both the answer and the chain-of-thoughts, we replace the suffix prompt p_s with “Explain the rationale to answer the question” to control whether we want a question answer or a CoT to explain the thinking steps. Thus, our optimization objective is:

$$\mathcal{L} = \mathcal{L}_{\text{label}} + \lambda \mathcal{L}_{\text{rationale}} = \sum_{j=1}^N \ell(\Phi(\mathcal{V}_j, q_j, p_s), \hat{y}_j) + \lambda \ell(\Phi(\mathcal{V}_j, q_j, p_s), c_j)$$

Here we set λ to 1 to ensure the importance of answer and rationale are equally considered, which can not only keep the capacity to predict the short question answer but also expand the ability to generate the rationale to solve the question. Notice that not all the QA pairs can generate qualified CoT. In that case, the $\mathcal{L}_{\text{rationale}}$ will be set to 0.

Table 3: Training and evaluation datasets statics.

Dataset	Avg		Size		Type	Train	Eval
	Duration (s)	train	train	eval			
MC-VQA							
STAR (Wu et al., 2021)	11.6	45.7K	7.1K	Compositional	✓	✓	
NExT-QA (Xiao et al., 2021)	44	34.1K	5.0K	Temporal & Causal	✓	✓	
CLEVRER (Yi et al., 2020)	5	21.0K	-	Spatial-temporal	✓	✗	
Perception-Test (Pătrăucean et al., 2023)	23	-	11.5K	General	✗	✓	
MVBench (Li et al., 2024b)	5-35	-	2.0K	General	✗	✓	
VideoMME (Fu et al., 2024)	1010	-	2.7K	General	✗	✓	
OE-VQA							
AGQA (Grunde-McLaughlin et al., 2021)	30	25.0K	2.0K	Compositional	✓	✓	
ANetQA (Yu et al., 2023b)	180	25.0K	2.0K	Compositional	✓	✓	
EgoQA (Grauman et al., 2022)	6.4	7.8K	-	Ego-centric	✓	✗	
Activitynet-QA (Yu et al., 2019)	112	-	8.0K	General	✗	✓	
Video-ChatGPT (Maaz et al., 2024)	108	-	3.0K	General	✗	✓	

3 EXPERIMENTS

In this section, we present the experimental setup (Sec. 3.1) and comparison results on various VideoQA benchmarks (Sec. 3.2). Extensive ablation studies are also undertaken to further examine the contributions of our approach in Sec. 3.3, and an evaluation on the quality of rationales generated by the distilled model is made in Sec. 3.4.

Table 4: Comparison with Video-LLMs on MC-VQA benchmarks. LLaVA-NeXT-Video-AoTD outperforms all other baselines the and the version without CoT distillation.

Model	MVBench (Acc.)	VideoMME (Acc.)	STAR (Acc.)	NEXT-QA (Acc.)	Perception-Test (Acc.)
Proprietary Models					
Gemini 1.0 Pro (Google, 2023)	-	-	-	-	51.1
Gemini 1.0 Ultra (Google, 2023)	-	-	-	-	54.7
Gemini 1.5 Pro (Google, 2024)	-	75.7	-	-	-
GPT4-V (OpenAI, 2023b)	43.7	60.7	-	-	-
GPT4-O (OpenAI, 2024)	-	66.2	-	-	-
Open-source Models					
LLaMA-VID (7B) (Li et al., 2024c)	41.9	25.9	-	-	44.6
Video-LLaVA (7B) (Lin et al., 2024)	41.0	40.4	-	-	44.3
VideoChat2 (7B) (Li et al., 2024b)	51.1	33.7	59.0	68.6	47.3
VideoLLaMA2 (7B) (Cheng et al., 2024)	53.4	44.0	58.5	62.3	49.6
LLaVA-NeXT-Video (7B) (Zhang et al., 2024)	46.5	41.0	52.4	61.6	47.5
LLaVA-NeXT-Video-Instruct (7B)	53.4	43.2	72.2	77.1	50.3
LLaVA-NeXT-Video-AoTD (7B)	55.6	45.0	74.3	77.6	50.6

3.1 EXPERIMENTAL SETUP

Base model. We use LLaVA-NeXT-Video (7B) (Zhang et al., 2024) (LNV for short) as base Video-LLM, which has shown remarkable performance on image-centric tasks, for example image QA (Yue et al., 2024). We present comparison on naive instruction tuning with video question answering dataset or with additional CoTs distillation. For CoT conversion and verification, we prompt LLaMA-3.1-8B with the manually-designed instruction and some in-context examples. Detailed prompts are provided in Appendix C.

Instruction tuning. As shown in Table 2, we utilize both multiple-choice and open-ended QA data, along with the generated CoTs, to fine-tune the base video question answering model. The resulting distilled model is named **LLaVA-NeXT-Video-AoTD** (LNV-AoTD for short). Additionally, as baseline, we also train another version of the model using only the basic QA data, which we refer to as **LLaVA-NeXT-Video-Instruct** (LNV-Instruct for short).

Evaluation benchmarks. We conduct extensive evaluations on Multiple-Choice Video QA (MC-VQA) and Open-Ended Video QA (OE-VQA). We report the top-1 accuracy for all MC benchmarks, which means the proportion of the output equal to the answer. For the evaluation on AGQA and ANetQA, we sample subsets from them, due to the large volume of test set. We report a GPT-assessed accuracy and score with the help of GPT-3.5-turbo-0613 for all OE benchmarks. The accuracy is a binary right or wrong choice and the score means similarity of output to the answer. We evenly select the benchmark in-domain and out-of-domain for testing to ensure a comprehensive and reasonable evaluation of the model capability. Detailed statistics for evaluation benchmarks are shown in Table 3.

3.2 QUANTITATIVE RESULTS

We divide the comparison into two parts: the first focuses on comparing the distilled model with other baselines, while the second examines the difference between the instruction-tuned model and the AoTD version. Note that, as the base model continues improving with more data and compute, we expect our proposed idea can be used to enhance the performance of any model.

MC-VQA performance. As shown in Table 4, our LLaVA-NeXT-Video-AoTD achieves superior performance across all benchmarks. Several key observations can be made: (i) Compared to the base model, even a simple instruction-tuning on certain VideoQA datasets significantly enhances the model’s question-answering performance. This improvement is notable since the base model

Table 5: Comparison with Video-LLMs on OE-VQA benchmarks. LLaVA-NeXT-Video-AoTD improves performance in all open-ended benchmarks compared with the Instruct version.

Model	ANetQA (Acc./Score)		AGQA (Acc./Score)		Video-ChatGPT (Score) Corr. Deta. Cont. Temp. Cons.				ActivityNet (Acc./Score)
Proprietary Models									
Gemini 1.0 Pro (Google, 2023)	-	-	-	-	-	-	-	-	49.8/-
Gemini 1.0 Ultra (Google, 2023)	-	-	-	-	-	-	-	-	52.2/-
Gemini 1.5 Pro (Google, 2024)	-	-	-	-	-	-	-	-	56.7/-
GPT4-V (OpenAI, 2023b)	-	-	4.09	3.88	4.37	3.94	4.02	-	59.5/-
GPT4-O (OpenAI, 2024)	-	-	-	-	-	-	-	-	61.9/-
Open-Source Models									
VideoLLaMA (7B) (Cheng et al., 2023)	-	-	1.96	2.18	2.16	1.82	1.79	-	12.4/1.1
Video-ChatGPT (7B) (Maaz et al., 2024)	-	-	2.50	2.57	2.69	2.16	2.20	-	35.2/2.7
LLaMA-VID (7B) (Li et al., 2024c)	-	-	2.96	3.00	3.53	2.46	2.51	-	47.4/3.3
Video-LLaVA (7B) (Lin et al., 2024)	-	-	2.87	2.94	3.44	2.45	2.49	-	45.3/3.3
VideoChat2 (7B) (Li et al., 2024b)	-	-	3.02	2.88	3.51	2.66	2.81	-	49.1/3.3
VideoLLaMA2 (7B) (Cheng et al., 2024)	-	-	3.09	3.09	3.68	2.63	3.25	-	49.9/3.3
LLaVA-NeXT-Video (7B) (Zhang et al., 2024)	46.4/3.3	27.4/2.2	3.26	3.22	3.77	2.47	2.99	-	54.3/3.2
LLaVA-NeXT-Video-Instruct (7B)	47.1/3.1	59.3/3.4	2.96	2.81	3.35	2.42	2.82	-	50.0/3.3
LLaVA-NeXT-Video-AoTD (7B)	53.9/3.4	60.9/3.6	3.11	3.00	3.60	2.41	2.91	-	53.2/3.4

was primarily trained on static images and struggled with video understanding. (ii) Our model, trained with CoTs distillation, demonstrates further performance enhancements across all benchmarks, particularly on the compositional VideoQA benchmark (STAR) and comprehensive benchmarks (VideoMME, MVBench). This suggests that our AoTD method effectively improves the model’s ability to address complex problems and interpret spatial-temporal scenes. (iii) The distilled model consistently outperforms all other baselines across all benchmarks, even when compared to more powerful models. This finding illustrates that our method effectively bridges performance gaps created by varying model components.

OE-VQA performance. As shown in Table 5, LLaVA-NeXT-Video-AoTD outperforms the Instruct variant across all open-ended VideoQA benchmarks. Notably, it achieves a greater percentage increase compared to the Multiple-Choice (MC-VQA) benchmarks, suggesting that CoTs distillation may be more effective for open-ended generation than for multiple-choice selection. While the distilled model scores higher than most models listed in the table, it does not surpass LLaVA-NeXT-Video on certain benchmarks. We conjecture this is due to the model’s extensive training on images, that can also benefit the question answering without requiring complex reasonings, as also suggested by the findings in VideoLLaMA2 (Cheng et al., 2024). Additionally, the inherent challenges of evaluating open-ended VQA may influence the results. Assessments conducted by GPT can be biased or inaccurate, and the metrics we employ primarily indicate general trends rather than providing absolute accuracy.

3.3 ABLATION STUDY

Analysis on CoT filtering. To demonstrate the effectiveness of our filtering mechanism, we trained an alternative model without CoTs filtering while maintaining all other settings. The amount of CoTs distillation data increased to 36.3K. As shown in Table 6, the model’s performance declines significantly on both the Multiple-Choice (MC-VQA) and Open-Ended VQA (OE-VQA) benchmarks when the CoT filtering mechanism is not utilized. This confirms that employing large language models (LLMs) to filter CoTs is an crucial for enhancing data quality.

Analysis on model transferability. As AoTD is a distillation method that leverages Chain-of-Thoughts (CoTs), it can theoretically be applied to any Video-LLMs. To assess the transferability of our method, we conduct experiments on another very recent model, LLaVA-OneVision(7B) (Li et al., 2024a). As shown in Table 6, our method still achieves significant improvements on the benchmarks, demonstrating both the transferability and robustness of the approach. Due to the rapid

advancements in the computer vision field, evaluating all models and benchmarks is prohibitively infeasible. Thus, we focus on assessing a single model against selected benchmarks to provide a representative evaluation.

3.4 EVALUATION ON RATIONALES

To verify whether the model has effectively learned multi-step reasoning through CoTs distillation, we analyze the rationales generated by the model. Specifically, we extract and evaluate the temporal and spatial information embedded within these rationales. This approach extends beyond merely assessing the correctness of the final answer, which could be influenced by biases or other external factors. By examining the reasoning process in detail, we gain a more accurate understanding of the model’s ability to perceive and reason about spatial and temporal relationships.

Evaluation protocols. We randomly select 200 samples from the STAR validation set and perform inference on this subset using the suffix prompt, recording the generated rationales. From these rationales, we extract the predicted temporal windows and bounding boxes, comparing them to the ground truth. For the spatial evaluation, we calculate the IoU between the predicted and ground truth bounding boxes. For the temporal evaluation, we compute both IoU and Recall, leveraging the frame-level scene graph annotations provided in the STAR dataset.

Evaluation results. Table 7 presents the evaluation results. For comparison, we also test UniVTG for temporal reasoning and OWL-ViT v2 for spatial reasoning. The results show that LLaVA-NeXT-Video-Instruct struggles to generate valid rationales, even when using the suffix prompt. In contrast, LLaVA-NeXT-Video-AoTD demonstrates comparable performance to specialized models in both spatial and temporal reasoning, indicating that the model successfully acquired these abilities through the distillation process.

Table 6: Ablation results of CoT filtering and model transferability.

Model	Filtering	MVBench (Acc.)	STAR (Acc.)	AGQA (Acc. / Score)
LNV-AoTD	✗	53.7	73.3	59.5/3.5
LNV-AoTD	✓	55.6	74.3	60.9/3.6
Onevision	-	58.0	65.9	39.0/3.0
Onevision-Instruct	-	59.2	75.8	65.6/3.7
Onevision-AoTD	✓	60.5	76.6	65.7/3.7

Table 7: Temporal and spatial abilities evaluation result.

Model	Temporal IoU (%)	Grounding Recall (%)	Spatial Grounding IoU (%)
UniVTG	22.8	31.0	-
OWL-ViT v2	-	-	64.7
LNV-Instruct	✗	✗	✗
LNV-AoTD	21.7	34.0	45.2

4 RELATED WORK

Video-language models (Video-LLMs). Most existing Video-LLMs are composed of a pre-trained visual encoder (like CLIP (Radford et al., 2021) or SigLIP (Zhai et al., 2023)) to encode video frames into a sequence of visual features, an adapter to transfer the visual features to tokens which can be understood by the language model, and a pre-trained LLM to output the final response. These models achieve strong ability for general vision-language tasks like Video question-answering (think the task as auto-regressive generation with question as prompt prefix). More recent works such as VideoLLaMA2 (Cheng et al., 2024), LLaVA-NeXT-Video (Zhang et al., 2024) and Videochat2 (Li et al., 2024b), with their excellent architecture design and reasonable instruction tuning data collection, have achieved a new level of zero-shot results in Video QA task. However, current end-to-end models still lack of interpretability for questions, as well as the ability to think and visually process complex problems in multiple steps, leads to their weakness in real complex scenarios, which is an important part for embodied learning and autonomous driving.

Visual Programing and Agents. With the progress of LLMs, some recent works (Gupta & Kembhavi, 2023; Surís et al., 2023) begin to try to use LLM as planner to solve the complex reasoning task in real scenarios. They attempt to decompose the question into some easier sub-questions, and use different specialist models as agents to solve these sub-questions, and finally gather them to get the answer of the raw question. MoReVQA (Min et al., 2024) proposes a multi-stage system, consisting

486 of an event parser, a grounding module, and a reasoning module with an external memory, getting a
487 strong zero-shot Video QA ability while is able to create interpretable intermediate outputs. VURF
488 (Mahmood et al., 2024) proposes a self-refinement method to resolve the LLM hallucinations to get
489 a more concise program based on the context cues. These models demonstrate a strong ability to
490 obtain trustworthy answers based on the intermediate evidence they get, but they lag far behind the
491 end to end model in terms of inference speed, and often require some in-context examples to assist
492 them in solving problems, which undoubtedly brings a lot of trouble to the use of these agent-based
493 models.

494 **Chain-of-Thought (CoT).** Recent advancements in Chain-of-Thought (Wei et al., 2022; Yao et al.,
495 2024) have made significant improvements in boosting the capabilities of LLMs. Though there have
496 been several works enhancing the power of LLMs through distilling the CoT into the model, we still
497 note a lack of research focused on applying CoT to video scenarios, as videos often have complex
498 spatio-temporal relationships, and multi-step thinking is needed to solve the problems happened in
499 these scenes. MotionEpic (Fei et al., 2024) develops a Video-of-Thought reasoning framework by
500 integrating video spatial-temporal scene graph. But it requires explicit training on the graph encoder,
501 which needs additional graph data, and cannot be directly migrated to other common Video-LLMs.
502 Thus, we construct natural language CoTs which are involved with spatial-temporal information to
503 adapt to any different models.

504 **Visual CoT.** The potential of Chain-of-Thought (CoT) reasoning extends beyond NLP to the visual
505 domain. Several studies (Zhang et al., 2023; Mitra et al., 2024; Shao et al., 2024; Gao et al., 2024b)
506 have applied CoT to visual understanding tasks, using powerful MLLMs for CoT generation or tool-
507 based architectures for step-by-step problem solving. However, these methods face limitations, such
508 as errors in CoT generation by MLLMs or high time and memory costs for tool-based approaches.
509 Recent works like Visual Program Distillation (VPD) (Hu et al., 2024a) and Fact (Gao et al., 2024a)
510 aim to maintain CoT accuracy and diversity while leveraging MLLMs to directly generate CoTs.
511 These methods decompose complex tasks through code programs, invoking expert models to ad-
512 dress sub-tasks, and use the generated CoTs as training data for fine-tuning visual-language models,
513 thereby improving the model’s ability to generate rationales directly. While all these methods focus
514 on image-based domains, they overlook the video domain, where CoT is especially suitable due
515 to the complex spatio-temporal nature of video understanding tasks. To bridge this gap, we pro-
516 pose AoTD, a method inspired by VPD and Fact, tailored to the video domain. Video-STaR (Zohar
517 et al., 2024) also constructs CoTs using videos and existing labels for instruction tuning, without
518 developing an agent-based system.

519 5 CONCLUSION & LIMITATION

520
521 In this work, we present Agent-of-Thought Distillation (AoTD), a novel approach aimed at distill-
522 ing multi-step reasoning and spatial-temporal understanding into a large generative video-language
523 model (Video-LLM). Our method introduces an agent-based system that automates the generation of
524 Chain-of-Thoughts (CoTs) from various Video Question Answering (VideoQA) datasets by break-
525 ing down complex questions into manageable sub-tasks that can be addressed by specialized vision
526 models. Extensive experiments validate that the distilled model significantly enhances performance
527 on both Multiple-Choice (MC-VQA) and Open-Ended VQA (OE-VQA) benchmarks, underscoring
528 the effectiveness of our approach.

529 Despite these advancements, several limitations remain and we leave them as future work: (i) Simi-
530 lar to prior approaches, the effectiveness of our agent-based system is contingent upon the progress
531 of the underlying visual model components. Enhancing its ability to generalize across diverse
532 datasets is essential for broader applicability. (ii) While our primary focus has been on compositional
533 VideoQA tasks, and we have demonstrated improvements across a series of benchmarks, achieving
534 holistic enhancements will require further exploration into creating a more balanced distribution of
535 training data. (iii) Furthermore, our agent-based framework has the potential to address additional
536 video-related tasks, such as video captioning and referring segmentation. We aim to expand our
537 methodology to these domains, which could yield even more robust and versatile applications in the
538 future. Overall, we believe AoTD represents a promising future direction for advancing multimodal
539 reasoning abilities in Video-LLMs.

REFERENCES

- 540
541
542 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
543 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
544 model for few-shot learning. In *Advances in Neural Information Processing Systems*, 2022.
- 545 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
546 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,
547 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi
548 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng
549 Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi
550 Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang
551 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint*
552 *arXiv:2309.16609*, 2023.
- 553 Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi
554 Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and
555 audio understanding in video-llms. In *Proceedings of the Conference on Empirical Methods in*
556 *Natural Language Processing*, 2023.
- 557 Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi
558 Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing spatial-temporal
559 modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024.
- 560 Guo Daya, Zhu Qihao, Yang Dejian, Dong Zhenda Xie, Kai, Zhang Wentao, Chen Guanting,
561 Bi Xiao, Y. Wu, Y.K. Li, Luo Fuli, and Liang Yingfei, Xiongand Wenfeng. Deepseek-coder:
562 When the large language model meets programming – the rise of code intelligence. *arXiv preprint*
563 *arXiv:2401.14196*, 2024.
- 564 Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. Videoagent: A
565 memory-augmented multimodal agent for video understanding. *arXiv preprint arXiv:2403.11481*,
566 2024.
- 567 Hao Fei, Shengqiong Wu, Wei Ji, Hanwang Zhang, Meishan Zhang, Mong-Li Lee, and Wynne Hsu.
568 Video-of-thought: Step-by-step video reasoning from perception to cognition. In *Proceedings of*
569 *the International Conference on Machine Learning*, 2024.
- 570 Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu
571 Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evalua-
572 tion benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
- 573 Minghe Gao, Shuang Chen, Liang Pang, Yuan Yao, Jisheng Dang, Wenqiao Zhang, Juncheng Li,
574 Siliang Tang, Yueting Zhuang, and Tat-Seng Chua. Fact: Teaching mllms with faithful, concise
575 and transferable rationales. In *ACM Multimedia*, 2024a.
- 576 Timin Gao, Peixian Chen, Mengdan Zhang, Chaoyou Fu, Yunhang Shen, Yan Zhang, Shengchuan
577 Zhang, Xiawu Zheng, Xing Sun, Liujuan Cao, et al. Cantor: Inspiring multimodal chain-of-
578 thought of mllm. In *ACM Multimedia*, 2024b.
- 579 Gemini Team Google. Gemini: a family of highly capable multimodal models. *arXiv preprint*
580 *arXiv:2312.11805*, 2023.
- 581 Gemini Team Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens
582 of context. *arXiv preprint arXiv:2403.05530*, 2024.
- 583 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-
584 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in
585 3,000 hours of egocentric video. In *Proceedings of the IEEE Conference on Computer Vision and*
586 *Pattern Recognition*, 2022.
- 587 Madeleine Grunde-McLaughlin, Ranjay Krishna, and Maneesh Agrawala. Agqa: A benchmark for
588 compositional spatio-temporal reasoning. In *Proceedings of the IEEE Conference on Computer*
589 *Vision and Pattern Recognition*, 2021.
- 590
591
592
593

- 594 Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning
595 without training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023.
596
597
- 598 Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ran-
599 jay Krishna, and Ariel Fuxman. Visual program distillation: Distilling tools and programmatic
600 reasoning into vision-language models. In *Proceedings of the IEEE Conference on Computer
601 Vision and Pattern Recognition*, 2024a.
- 602 Ziniu Hu, Ahmet Iscen, Chen Sun, Kai-Wei Chang, Yizhou Sun, David Ross, Cordelia Schmid, and
603 Alireza Fathi. Avis: Autonomous visual information seeking with large language model agent. In
604 *Advances in Neural Information Processing Systems*, 2024b.
- 605 De-An Huang, Shijia Liao, Subhashree Radhakrishnan, Hongxu Yin, Pavlo Molchanov, Zhiding Yu,
606 and Jan Kautz. Lita: Language instructed temporal-localization assistant. In *Proceedings of the
607 European Conference on Computer Vision*, 2024.
608
- 609 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li,
610 Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024a.
- 611 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,
612 Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In
613 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2024b.
614
- 615 Liunian Harold Li*, Pengchuan Zhang*, Haotian Zhang*, Jianwei Yang, Chunyuan Li, Yiwu Zhong,
616 Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao.
617 Grounded language-image pre-training. In *Proceedings of the IEEE Conference on Computer
618 Vision and Pattern Recognition*, 2022.
- 619 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language
620 models. In *Proceedings of the European Conference on Computer Vision*, 2024c.
621
- 622 Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual
623 representation by alignment before projection. In *Proceedings of the Conference on Empirical
624 Methods in Natural Language Processing*, 2024.
- 625 Kevin Qinghong Lin, Pengchuan Zhang, Joya Chen, Shraman Pramanick, Difei Gao, Alex Jin-
626 peng Wang, Rui Yan, and Mike Zheng Shou. Univtg: Towards unified video-language temporal
627 grounding. In *Proceedings of the International Conference on Computer Vision*, 2023.
628
- 629 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:
630 Towards detailed video understanding via large vision and language models. In *Association for
631 Computational Linguistics*, 2024.
- 632 Ahmad Mahmood, Ashmal Vayani, Muzammal Naseer, Salman Khan, and Fahad Khan. Vurf: A
633 general-purpose reasoning and self-refinement framework for video understanding. *arXiv preprint
634 arXiv:2403.14743*, 2024.
- 635 Minderer Matthias, Gritsenko Alexey, Stone Austin, Neumann Maxim, Weissenborn Dirk, Doso-
636 vitskiy Alexey, Mahendran Aravindh, Arnab Anurag, Dehghani Mostafa, Shen Zhuoran, Wang
637 Xiao, Zhai Xiaohua, Kipf Thomas, and Houlby Neil. Simple open-vocabulary object detection
638 with vision transformers. In *Proceedings of the European Conference on Computer Vision*, 2022.
639
- 640 Juhong Min, Shyamal Buch, Arsha Nagrani, Minsu Cho, and Cordelia Schmid. Morevqa: Exploring
641 modular reasoning models for video question answering. In *Proceedings of the IEEE Conference
642 on Computer Vision and Pattern Recognition*, 2024.
- 643 Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection.
644 In *Advances in Neural Information Processing Systems*, 2024.
645
- 646 Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Cudas, Clarisse Simoes, Sahaj Agar-
647 wal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, et al. Orca 2: Teaching
small language models how to reason. *arXiv preprint arXiv:2311.11045*, 2023.

- 648 Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-
649 thought prompting for large multimodal models. In *Proceedings of the IEEE Conference on*
650 *Computer Vision and Pattern Recognition*, 2024.
- 651 OpenAI. Gpt-3.5-turbo system card, 2023a. URL [https://platform.openai.com/docs/
652 models/gpt-3-5-turbo](https://platform.openai.com/docs/models/gpt-3-5-turbo).
- 653 OpenAI. Gpt-4v(ision) system card, 2023b. URL [https://openai.com/research/
654 gpt-4v-system-card](https://openai.com/research/gpt-4v-system-card).
- 655 OpenAI. Gpt-4o system card, 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- 656 Viorica Pătrăucean, Lucas Smaira, Ankush Gupta, Adrià Recasens Contente, Larisa Markeeva,
657 Dylan Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, Carl Doersch,
658 Tatiana Matejovicova, Yury Sulsky, Antoine Miech, Alex Frechette, Hanna Klimczak, Raphael
659 Koster, Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osindero, Dima Damen, Andrew
660 Zisserman, and João Carreira. Perception test: A diagnostic benchmark for multimodal video
661 models. In *Advances in Neural Information Processing Systems*, 2023.
- 662 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
663 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
664 models from natural language supervision. In *Proceedings of the International Conference on*
665 *Machine Learning*, 2021.
- 666 Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multimodal
667 large language model for long video understanding. In *Proceedings of the IEEE Conference on*
668 *Computer Vision and Pattern Recognition*, 2024.
- 669 Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hong-
670 sheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset and
671 benchmark for chain-of-thought reasoning. In *Advances in Neural Information Processing Sys-*
672 *tems*, 2024.
- 673 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for
674 reasoning. In *Proceedings of the International Conference on Computer Vision*, 2023.
- 675 Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Chenting Wang, Guo Chen, Baoqi Pei,
676 Rongkun Zheng, Jilan Xu, Zun Wang, et al. Internvideo2: Scaling video foundation models
677 for multimodal video understanding. In *Proceedings of the European Conference on Computer*
678 *Vision*, 2024.
- 679 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
680 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances*
681 *in Neural Information Processing Systems*, 2022.
- 682 Zejia Weng, Xitong Yang, Ang Li, Zuxuan Wu, and Yu-Gang Jiang. Open-vclip: Transforming clip
683 to an open-vocabulary video model via interpolated weight optimization. In *Proceedings of the*
684 *International Conference on Machine Learning*, 2023.
- 685 Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. Star: A benchmark for
686 situated reasoning in real-world videos. In *Advances in Neural Information Processing Systems*,
687 2021.
- 688 Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-
689 answering to explaining temporal actions. In *Proceedings of the IEEE Conference on Computer*
690 *Vision and Pattern Recognition*, 2021.
- 691 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
692 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Ad-*
693 *vances in Neural Information Processing Systems*, 2024.
- 694 Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B.
695 Tenenbaum. CLEVRER: collision events for video representation and reasoning. In *Proceedings*
696 *of the International Conference on Learning Representations*, 2020.

702 Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model
703 for video localization and question answering. In *Advances in Neural Information Processing*
704 *Systems*, 2023a.

705
706 Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-
707 qa: A dataset for understanding complex web videos via question answering. In *Proceedings of*
708 *the AAAI Conference on Artificial Intelligence*, 2019.

709
710 Zhou Yu, Lixiang Zheng, Zhou Zhao, Fei Wu, Jianping Fan, Kui Ren, and Jun Yu. Anetqa: A large-
711 scale benchmark for fine-grained compositional reasoning over untrimmed videos. In *Proceedings*
712 *of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023b.

713
714 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
715 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal
716 understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE Conference*
717 *on Computer Vision and Pattern Recognition*, 2024.

718
719 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
720 image pre-training. In *Proceedings of the International Conference on Computer Vision*, 2023.

721
722 Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and
723 Chunyuan Li. Llava-next: A strong zero-shot video understanding model, 2024. URL <https://llava-vl.github.io/blog/2024-04-30-llava-next-video/>.

724
725 Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal
726 chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.

727
728 Orr Zohar, Xiaohan Wang, Yonatan Bitton, Idan Szepktor, and Serena Yeung-Levy. Video-star: Self-
729 training enables video instruction tuning with any supervision. *arXiv preprint arXiv:2407.06189*,
730 2024.

731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A EXPERIMENTAL DETAILS

A.1 TRAINING DETAILS

For all models, their projection layers and language model are finetuned and visual encoder is frozen. We use a cosine learning rate schedule, with warm up ratio 0.03 and learning rate $4e-5$. For both Instruct and AoTD setting, we finetune the model with batch size 48 and totally 1 epoch. We believe that longer training will get a better performance on in-domain benchmarks but maybe a destroy on out-of-domain benchmarks.

A.2 SPECIALIZED MODELS EVALUATION DETAILS

In this section we will show the details about each sub-task’s evaluation from data preparation to evaluation metric.

Question Decomposition. Since there may be multiple valid ways to decompose the same problem, we evaluate only the accuracy of the final output in this sub-task. Specifically, the model takes the query and instruction as input and generates an executable program. We replace all intermediate outputs within the program and focus on whether the final output matches the correct answer. If the decomposition is correct, the final output must align with the answer. Any programs that cannot be executed or that lead to an incorrect answer are considered failures.

Object Detection. To evaluate the performance of detection models, we sample frames with scene graph annotations from the input video clip and provide them, along with the text query, as input to the model. The model then outputs a series of bounding boxes that exceed a confidence threshold. We select the bounding box with the highest confidence as the final output and calculate the IoU to assess accuracy.

Temporal Grounding. Since scene graphs provide both the start and end frame IDs, as well as key frame IDs for each event, we use IoU and Recall as metrics to capture different aspects of model performance. The model takes the video clip and text query as input and outputs the predicted start and end frame IDs. We calculate IoU based on the alignment between the predicted and annotated start and end frame IDs, and we compute Recall using the key frame ID annotations to evaluate how well the model captures important frames.

Action Recognition. For discriminative models, we provide the video clip and a list of action labels as input to complete a classification task. For generative models, we provide the video clip along with an instruction prompt, asking the model to generate five actions most relevant to the video, ranked by likelihood. We then use the top-ranked output from each model to calculate the Top-1 accuracy for both approaches.

Question Answering. The evaluation of question answering follows a similar approach to previous methods. The model takes the video clip and question as input and returns an answer, from which we directly calculate the accuracy. The key difference between this sub-task and a standard QA task is that the answers are based on a series of information collected by preceding agents, allowing for a more accurate assessment of the model’s pure question-answering ability.

B MORE RESULTS

Here we introduce some examples to show the process from query to Chain-of-Thought using our agent-based pipeline. We can find that our pipeline is able to decompose complex questions into easier sub-tasks and the final CoT retains step-by-step problem-solving ideas and spatial-temporal information representing video understanding ability.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863



Figure 4: Example form NExT-QA (Xiao et al., 2021)



Figure 5: Example form ANetQA (Yu et al., 2023b)

C PROMPTS

In this section we present the prompts used in our agent-based pipeline for generating program, converting execution trace and filtering rationales.

C.1 PROMPT FOR PROGRAM GENERATION

For each video and query, we call a LLM to decompose the query to a Python program under the guidance of the prompt below. We modify the ViperGPT (Surís et al., 2023) prompt to adapt to the visual agents we use.

```
1 def Query_Objs(clip, query):
2     """
3     Query the objects that appear in video clip and match the query descriptions.
4     Parameters
5     -----
6     clip:
7         a list of video frames.
8     query:
9         Description of the target object.
10    Returns
11    -----
12    a list of bounding boxes of the objects that match the query.
13    Examples
14    -----
15    #return white_objs
16    def execute_command(video_clip):
17        white_objs = Query_Objs(video_clip, "white object")
18        return white_objs
```

```

918 19 """
919 20
920 21 def Query_Actions(clip, obj=None):
921 22 """
922 23 Find the actions happened in the video clip, if obj is not None, query the actions related
923 24 to it.
924 25 Parameters
925 26 -----
926 27 clip:
927 28     a list of the video frames.
928 29 obj:
929 30     object class which is used to query the actions related to it.
930 31 Returns
931 32 -----
932 33 a list of actions classes happened in the video clip.
933 34 Examples
934 35 -----
935 36 #return actions
936 37 def execute_command(video_clip, query, possible_answers):
937 38     actions = Query_Actions(video_clip)
938 39     return actions
939 40 """
940 41 def Filter_frames_with_act(clip, action):
941 42 """
942 43 filter a new video clip containing the time period in which the target action occurred
943 44 Parameters
944 45 -----
945 46 clip:
946 47     a list of video frames.
947 48 action:
948 49     the target action which is used to filter frames.
949 50 Returns
950 51 -----
951 52 a new video clip containing the time period in which the target action occurred.
952 53 Examples
953 54 -----
954 55 #return jump_clip
955 56 def execute_command(video_clip, query, possible_answers):
956 57     jump_clip = Filter_frames_with_act(video_clip, "person is jumping")
957 58     return jump_clip
958 59 """
959 60
960 61 def Filter_frames_with_obj(clip, obj):
961 62 """
962 63 filter a new video clip that the target object occurred.
963 64 Parameters
964 65 -----
965 66 clip:
966 67     a list of video frames.
967 68 obj:
968 69     class or description about the target object.
969 70 Returns
970 71 -----
971 72 a new video clip that the target object occurred in it.
972 73 Examples
973 74 -----
974 75 #return shoe_clip
975 76 def execute_command(video_clip, query, possible_answers):
976 77     shoe_clip = Filter_frames_with_obj(video_clip, "shoe")
977 78     return shoe_clip
978 79 """
979 80
980 81 def trim(clip, start=None, end=None):
981 82 """
982 83 Returns a new video clip containing a trimmed version of the original video at the [start,
983 84 end] clip.
984 85 Parameters
985 86 -----
986 87 clip:
987 88     a list of video frames.
988 89 start : Union[int, None]
989 90     An int describing the starting frame in this video clip with respect to the original
990 91 video.
991 92 end : Union[int, None]
992 93     An int describing the ending frame in this video clip with respect to the original
993 94 video.
994 95 Returns
995 96 -----
996 97 a new video clip with start and end.

```

```

972 96 """
973 97 def Find(clip, obj):
974 98 """
975 99 find all bounding boxes around a certain object in the video clip,
100 100 and collates them into a collection of frames.
976 101 Parameters
102 102 -----
977 103 clip:
978 104     a list of video frames.
105 105 obj:
979 106     the object to look for.
107 107 Returns
108 108 -----
980 109 a new video clip composed of crops of the object.
981 110 Examples
111 111 -----
982 112 # Return the shoe_clip
983 113 def execute_command(video_clip, query, possible_answers):
984 114     shoe_clip = Find(video_clip, "shoe")
985 115     return shoe_clip
986 116 """
117 117
987 118 def select_answer(query, info, possible_answers):
988 119 """
989 120 Uses a language model to choose the option that best answers the question given the input
121 121 information.
990 122 Parameters
123 123 -----
991 124 query:
992 125     the input question.
126 126 info:
993 127     Any useful information to answer the question.
128 128 possible_answers:
994 129     a list of possible answers to the question.
130 130 Returns
995 131 -----
996 132 one answer chosen from the possible answers.
997 133 Examples
998 134 -----
999 135 # Return the answer
1000 136 def execute_command(video_clip, query, possible_answers):
1001 137     clip_summary = Video_summary(video_clip)
1002 138     info = {
1003 139         "summary of the target video": clip_summary
1004 140     }
1005 141     answer = select_answer(query, info, possible_answers)
1006 142     return answer
1007 143 """
1008 144 def exist(clip, query):
1009 145 """
1010 146 judge whether a object exists in the video.
1011 147 Parameters
1012 148 -----
1013 149 clip:
1014 150     a list of video frames.
1015 151 query:
1016 152     query to the object class.
1017 153 Returns
1018 154 -----
1019 155 Return True if the object specified by query is found in the video, and False otherwise.
1020 156 Examples
1021 157 -----
1022 158 # Return the flag
1023 159 def execute_command(video_clip, query, possible_answers):
1024 160     flag = exist(video_clip, "shoe")
1025 161     return flag
1026 162 """
1027 163 def Video_summary(clip, query):
1028 164 """
1029 165 give a brief summary of the video clip related to the query.
1030 166 Parameters
1031 167 -----
1032 168 clip:
1033 169     a list of video frames.
1034 170 query:
1035 171     a question about the video.
1036 172 Returns
1037 173 -----
1038 174 return a brief summary of the video clip.
1039 175 Examples
1040 176 -----

```

```

1026
1027 176 # Return the clip_summary
1028 177 def execute_command(video_clip, query, possible_answers):
1029 178     clip_summary = Video_summary(video_clip, query)
1030 179     return clip_summary
1031 180 """
1032 181 Write a function using Python and the functions (above) that could be executed to provide an
1033 182 answer to the query.
1034 183 Consider the following guidelines:
1035 184 - Use base Python (comparison, sorting) for basic logical operations, start/end, math, etc.
1036 185 - Objects with mutiple names like "phone/camera", "cup/glass/bottle" with slash, input them as
1037 186 a whole object name.
1038 187 - Just use the class and function appear above except for some base python operations.
1039 188 - Only answer with a function starting def execute_command, do not answer any extra words and
1040 189 symbols before and after the function.
1041 190 - No text that is not related to function can appear.
1042 191 - the answer only begins with "def execute_command" and ends with "return answer".
1043 192 Here are some examples of the function you should write:
1044 193 -----
1045 194 question: What else is the person able to do with the door?
1046 195 possible answers: ["Hold the door.", "Put down the door.", "Close the door.", "Open the door."
1047 196 ]
1048 197 def execute_command(video_clip, query, possible_answers):
1049 198     door_clip = Filter_frames_with_obj(video_clip, "door")
1050 199     person_clip = Find(door_clip, "person")
1051 200     clip_summary = Video_summary(person_clip, query)
1052 201     door_actions = Query_Actions(person_clip, "door", possible_answers=possible_answers)
1053 202     door_actions =
1054 203     info = {
1055 204         "actions the person able to do with the door else": door_actions,
1056 205         "summary of the target video": clip_summary
1057 206     }
1058 207     answer = select_answer(query, info, possible_answers)
1059 208     return answer
1060 209 -----
1061 210 Query: INSERT_QUERY_HERE
1062 211 possible answers: INSERT_POSSIBLE_ANSWERS_HERE

```

1052

1053

1054 C.2 PROMPT FOR EXECUTION TRACE CONVERSION

1055

1056 After getting the execution trace by running the program step by step, we use a LLM to convert
 1057 the trace into a natural language CoT. The LLM takes query, execution trace, possible answers (in
 1058 MC-VQA) and execution trace as input. The instruction prompt is as follow:

```

1059 1 Given a video and a question, I wrote the function execute_command using Python, and the other
1060 2 functions above that could be executed to provide an answer to the query.
1061 3 As shown in the code, the code will print execution traces.
1062 4 I need you to rewrite the execution trace into a natural language rationale that leads to the
1063 5 answer.
1064 6 Consider the following guidelines:
1065 7 - Use all the bounding box information in the rationale, do not use words like "so on" to omit
1066 8 the bounding box, just write all of them into the rationale.
1067 9 - Referencing the execution trace, write a reasoning chain that leads to the most common human
1068 10 answer. Notice that the output should be the same as the human answer, not necessarily the
1069 11 program output.
1070 12 - If some part of the rationale lacks logic, add reasonable content to make it logical.
1071 13 Here are some examples of the rationale you should write:
1072 14 -----
1073 15 Question: What did the person do with the table?
1074 16 def execute_command(video_clip, query, possible_answers, time_wait_between_lines, syntax):
1075 17     table_clip = Filter_frames_with_act(video_clip, 'person interacting with table')
1076 18     person_clip = Find(table_clip, 'person')
1077 19     table_bboxes = Find(table_clip, 'table')
1078 20     clip_summary = Video_summary(person_clip)
1079 21     person_action = Query_Actions(person_clip, 'table', possible_answers=possible_answers)
1080 22     info = {'actions the person do with the table': person_action, 'summary of the target
1081 23 video': clip_summary}
1082 24     answer = select_answer(query, info, possible_answers)
1083 25     return answer
1084 26 Execution trace:
1085 27 call Filter_frames_with_act
1086 28 filter action person interacting with table
1087 29 find action from frame 2 to frame 11
1088 30 call function Find

```

```

1080 28 finding person
1081 29 find person at [139, 141, 229, 342] in frame 2
1082 30 find person at [151, 123, 242, 349] in frame 3
1083 31 find person at [153, 121, 242, 274] in frame 4
1084 32 find person at [158, 123, 255, 261] in frame 5
1085 33 find person at [163, 124, 270, 262] in frame 6
1086 34 find person at [153, 121, 242, 351] in frame 7
1087 35 find person at [95, 113, 196, 316] in frame 8
1088 36 find person at [83, 113, 196, 285] in frame 9
1089 37 find person at [112, 116, 201, 332] in frame 10
1090 38 call function Find
1091 39 finding table
1092 40 find table at [183, 140, 269, 257] in frame 2
1093 41 find table at [194, 131, 269, 255] in frame 3
1094 42 find table at [227, 129, 269, 252] in frame 4
1095 43 find table at [226, 165, 269, 258] in frame 5
1096 44 find table at [233, 170, 270, 259] in frame 6
1097 45 find table at [217, 129, 269, 256] in frame 7
1098 46 find table at [217, 122, 270, 254] in frame 8
1099 47 find table at [221, 123, 269, 256] in frame 9
1100 48 find table at [225, 125, 270, 263] in frame 10
1101 49 call function Video_summary
1102 50 summary result: The video shows a man in a kitchen, bending over and holding an orange object,
1103 51 surrounded by various kitchen items and furniture, with a focus on his actions and the
1104 52 domestic setting.
1105 53 call function Query_Actions
1106 54 Query table
1107 55 Answer: tidied up.
1108 56 call function select_answer
1109 57 the information used: - actions the person do with the table: tidied up.
1110 58 - summary of the target video: The video shows a man in a kitchen, bending over and holding an
1111 59 orange object, surrounded by various kitchen items and furniture, with a focus on his actions
1112 60 and the domestic setting.
1113 61 program output: Tidied up.
1114 62 Rationale:
1115 63 To solve this question, we first have to find when did the person interact with the table.
1116 64 From the video, we can see that the person is interacting with the table from frame 2 to frame
1117 65 11.
1118 66 In this time period, we can find person at [139, 141, 229, 342] in frame 2, [151, 123, 242,
1119 67 349] in frame 3, [153, 121, 242, 274] in frame 4 and so on.
1120 68 Table can also be found at [183, 140, 269, 257] in frame 2, [194, 131, 269, 255] in frame 3,
1121 69 [227, 129, 269, 252] in frame 4 and so on.
1122 70 By analyzing the person and table bounding box region, we can see that the person is holding
1123 71 an orange object to clean the table in the kirchen environment.
1124 72 So the answer should be tidied up.
1125 73 -----
1126 74 Now, look the question, program and execution trace, please transfer these information to a
1127 75 rationale.
1128 76 Question: INSERT_QUESTION_HERE
1129 77 INSERT_PROGRAM_HERE
1130 78 Execution trace:
1131 79 INSERT_EXECUTION_TRACE_HERE
1132 80 Rationale:

```

1116

1117

1118

1119

1120

1121 C.3 PROMPT FOR COT FILTERING

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

```

1 I will give you a question and a rationale to solve the question, you need to judge whether
2 the rationale is thinking step by step and helpful to solve the question.
3 If yes, return True, If not, return False. no need to explain.
4 Here is the question and rationale:
5 Question: INSERT_QUESTION_HERE
6 Rationale: INSERT_RATIONALE_HERE

```

1134 C.4 PROMPT FOR INFERENCE
1135

1136 Question: question content
1137 Options:
1138 (A) option content
1139 (B) option content
1140 (C) option content
1141 (D) option content
1142 Answer with the option's letter from the given choices directly and only give the best option. /
1143 Explain the rationale to answer the question.
1144

1145
1146 Question: question content
1147 Answer in one word or phrase. / Explain the rationale to answer the question.
1148

1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187