

---

# Adaptive Video Understanding Agent: Enhancing efficiency with dynamic frame sampling and feedback-driven reasoning

---

Sullam Jeoung<sup>1\*</sup> Goeric Huybrechts<sup>2</sup> Bhavana Ganesh<sup>2</sup>  
Aram Galstyan<sup>2</sup> Sravan Bodapati<sup>2</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign <sup>2</sup>Amazon AGI  
sjeoung2@illinois.edu {huybrech, ganesh, argalsty, sravanb}@amazon.com

## Abstract

Understanding long-form video content presents significant challenges due to its temporal complexity and the substantial computational resources required. In this work, we propose an agent-based approach to enhance both the efficiency and effectiveness of long-form video understanding by utilizing large language models (LLMs) and their tool-harnessing ability. A key aspect of our method is query-adaptive frame sampling, which leverages the reasoning capabilities of LLMs to process only the most relevant frames in real-time, and addresses an important limitation of existing methods which typically involve sampling redundant or irrelevant frames. To enhance the reasoning abilities of our video-understanding agent, we leverage the self-reflective capabilities of LLMs to provide verbal reinforcement to the agent, which leads to improved performance while minimizing the number of frames accessed. We evaluate our method across several video understanding benchmarks and demonstrate that not only it enhances state-of-the-art performance but also improves efficiency by reducing the number of frames sampled.

## 1 Introduction

Recent advancements in video understanding have been significantly driven by end-to-end pretrained large transformer models, particularly those built upon large language models (LLMs) [16, 15], known as multimodal LLMs. Despite these advancements, comprehending long form videos remains a considerable challenge due to prohibitive computational costs and suboptimal performance [2]. Various approaches have been proposed to extend the temporal context of video transformers, including techniques such as masking, attention approximations, and parametric memory modules (e.g. [36], [24]). However, these methods often add complexity by necessitating specialized architectures and training paradigms [30].

Efficient video processing requires strategic selection of relevant frames from the total video sequence [5, 13]. Traditionally, methods in this domain mostly rely on uniform sampling [43, 30] or selective retrieval from a subset of sampled frames [3, 33]. While these techniques improve processing efficiency by reducing the number of frames, they often lack adaptability, leading to potential redundancy.

To address the above shortcomings, here we propose a novel approach that leverages LLMs as adaptive agents for video understanding tasks. Our method utilizes the advanced reasoning, planning, and tool-use capabilities of LLMs ([20, 48, 28]) to enhance sampling efficiency while maintaining robust performance in video understanding tasks. Specifically, our approach leverages a LLM-based agent that dynamically determines which frames to sample based on the specific context and query.

---

\*Work done during an internship at Amazon AGI

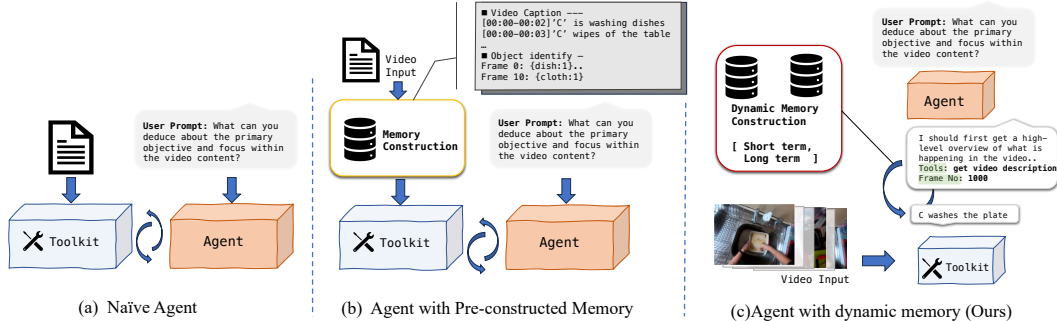


Figure 1: Comparison of methods: Our proposed method (c) is **query adaptive**, dynamically selecting frames based on query and video input to construct a responsive memory. In contrast, previous methods, including (a) Naïve agents and (b) Agents with pre-constructed memory, do not adapt to specific queries or utilize memory dynamically. We demonstrate that dynamically sampling frames have advantage over different set of benchmarks.

This method ensures that frame selection is both relevant and efficient, effectively mitigating the limitations of static sampling methods.

Our approach draws inspiration from research indicating that humans strategically allocate attention and filter out irrelevant details based on the task at hand [11, 8, 26]. For example, when asked "What is the main goal of the camera wearer in this video?" versus "What is the color of the bird that appears at the beginning?", humans deploy distinct strategies: the former may necessitate a review of the entire video to understand its context, whereas the latter would involve focusing solely on the video’s initial segment to identify the bird’s color.

Our proposed framework adaptively samples and processes video frames in response to specific queries (see Figure 1 c). While previous approaches rely on static process which is independent of the query in extracting information such as captions [3, 33], our approach attends to the given query and reasons strategically which frames to process during inference time without having to go through whole set of frames.

Our findings indicate that LLM agents, when used without guidance, exhibit suboptimal reasoning performance in terms of selecting the most informative frames. To enhance the reasoning ability of LLMs, we leverage the self-reflective capabilities of LLMs to provide insightful feedback [29, 21]. Specifically, reflective statements serve as a form of verbal reinforcement, enabling the agents to develop an updated policy that facilitates more nuanced and sophisticated reasoning. Furthermore, our framework integrates long-term memory to store and utilize past experiences. The reasoning trajectories and the refinement is stored in the memory per instance. The key rationale behind adopting the memory is that retrieving past experiences that are relevant and semantically similar to a given query can significantly enhance the reasoning behavior of the LLM.

We validate the generalizability of our framework by evaluating it across a range of benchmarks, demonstrating its effectiveness and adaptability in various video understanding tasks. The results indicate that the proposed method outperforms existing approaches, achieving higher accuracy while maintaining a lowest number of frames accessed.

## 2 Related Work

### 2.1 Long Context Multimodal Agents

Several approaches have been developed to handle multimodal inputs through agent-based reasoning [4, 39, 3, 33]. These methods leverage agents’ reasoning abilities along with their tool-calling capabilities. For instance, [39] employs Monte Carlo Tree Search for reasoning combined with tool-calling techniques, while [4] utilizes ReAct [41] for flexible video input processing.

Recent advancements have also focused on long-context videos [3, 33]. For example, [3] uses memory retrieval during inference to address specific queries, which can be effective for localizing detailed information but may become redundant depending on the query type. Similarly, [33] relies on predefined sampling methods, necessitating extensive frame access for caption generation. [31]

aims to reduce frame access by using a predefined number of frames and dynamic sampling, but primarily addresses short-form videos and straightforward question-answering scenarios.

Existing methods for addressing long-context processing using agent based approach (see Fig 1, b) involves preprocessing and extracting relevant information from frames during a pre-processing stage, with the agent retrieving memory dynamically based on the question during runtime [3, 33]. Although this approach can be effective, it is resource-intensive in terms of memory and processing time. Additionally, it operates in a static manner, irrespective of the specific question, which can be redundant.

## 2.2 Frame Sampling Methods

Several methods have been proposed to enhance the efficiency of video frame handling by selectively subsampling relevant frames based on the content of the question or text, rather than using uniform sampling [5, 13, 42, 21]. For example, [27] use CLIP model to retrieve pertinent frames through text prompts, while [7] propose a sampling technique that selects the most significant frames based on learned patterns. Although these approaches are effective, they often require pre-defined number of frames to sample or accessing to near all video frames to identify the relevant ones. These static ways of sampling frames may induce inefficiency as the video length becomes longer with exhaustive number of frames.

In contrast, our approach is inspired by human cognitive processes, which adaptively focus on information pertinent to the task at hand [11, 8, 26, 9]. We propose an agent that reasons about which frames to select based on the information from the question or previously extracted information, thereby improving the efficiency of information processing. While our method is similar to [34] in its query-adaptive nature, our method avoids the need for preprocessing (e.g., KNN clustering), thereby mitigating time-consuming operations.

Model	Long-Context	Query Adaptive Sampling	Long-term Memory	Reasoning
AssistGPT [4]	✗	✗	✗	ReAct
DoraemonGPT [39]	✗	✗	✗	MCTS
VideoAgent [3]	✓	✗	✗	ReAct
LifelongMemory [33]	✓	✗	✗	Prediction Ensemble
Ours	✓	✓	✓	Refinement + ReAct

Table 1: Comparison of existing methods. Previous approaches attempted to handle long-form video agents, however, our approach focuses on addressing long-context videos, adopting query adaptive sampling and long-term memory.

## 3 Adaptive Video Understanding Agent

We propose an AVUA: Adaptive Video Understanding Agent, which reasons which frames to process based on the observations and interactions made between the tools. Specifically, inspired by recent advancements in self-reflective ability of LLMs [10, 22, 29], we utilize the error feedback of LLMs to enhance the reasoning of the agent. We formulate the task likewise: The dataset  $\mathcal{D} = (Q, A, V)$  consists of question  $Q$ , answer  $A$ , and corresponding  $V$ . The agent  $\mathcal{L}$  is equipped with available actions  $\mathcal{A}$ . The agent  $\mathcal{L}$  has only access to the meta-data of the video  $V'$  (e.g. the total number of frames).

**Generating Policy** As illustrated in Figure 2, the initial step involves generating a policy  $\pi$  based on the question and the details of the video. This policy encompasses an analysis of the question type and a detailed question analysis, which includes a sampling strategy and identification of key elements that the agent should focus on during the reasoning process. The policy serves a dual purpose: it guides the agent in planning and reasoning, and it can be abstracted and utilized in long-term memory. The rationale behind this approach is that, while the specifics of the question may vary, the abstracted high-level question type can be retained and leveraged in a manner similar to how humans utilize their generalized experiences.

**Planning/tool invoking** At time step  $t$ , the agent  $\mathcal{L}$  selects an action  $a_t$  and action input  $x_t$  based on policy  $\pi$  in solving problem  $\mathcal{D}$ . The actions  $\mathcal{A}$  are the invocable tools, which are pre-defined and callable functions from the agent. The action input  $x_t$  is typically the frame number, indicating

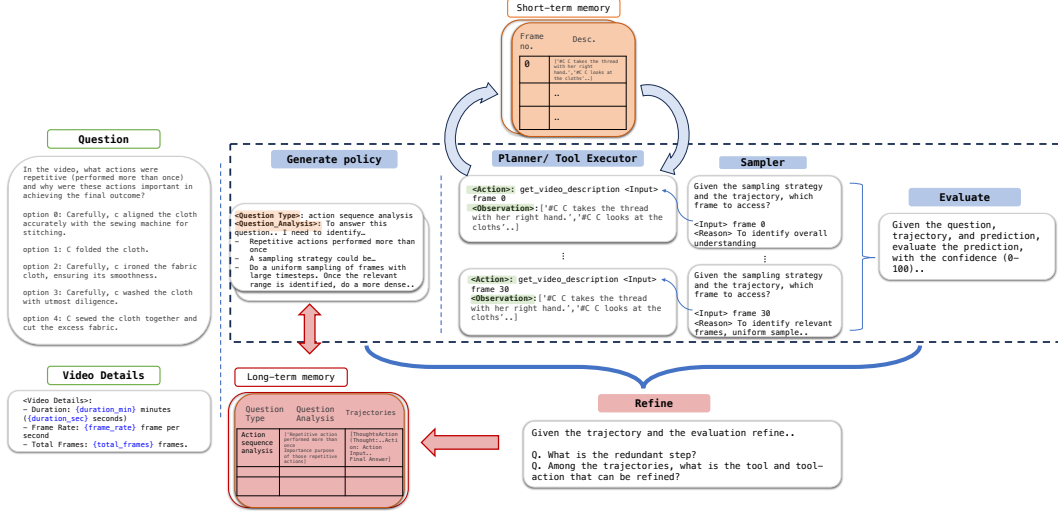


Figure 2: **Overall Framework.** The video metadata and question are provided to the agent to **generate policy**, which includes analyzing the question type and determining the task-solving strategy, including the sampling strategy. **Planner/Tool Executor**, based on the ReAct-style reasoning, generates thought processes, actions, and action inputs, and receives observations from the tools. During this stage, the sampler may suggest improved frames. After formulating the final answer, the evaluator and refiner are applied. The final result is then stored in **long-term memory**.

which frames the tools should access. The input often includes extra arguments, for example the question to query the tools (e.g. Frame index 0, what is happening in the frame?). Once the tools are invoked, it returns a observation  $\mathcal{O}$  which is the extracted information of the selected frame. The agent  $\mathcal{L}$  considers the previous observation-action trajectory  $\tau_t = [a_1, o_1, \dots, o_{t-1}]$ : in choosing which actions to call.

$$a_t = \mathcal{L}(\pi, \mathcal{D}, \tau_{t-1})$$

Specifically, the agent  $\mathcal{L}$  navigates search space,  $\mathcal{F} \times \mathcal{A}$ , where  $\mathcal{F}$  represents the set of frames within  $\mathcal{V}(|V| = |F| = n)$ . The main goal of the agent  $\mathcal{L}$  is to effectively prune the search space (i.e., minimize the number of the frames access) while ensuring performance (i.e., maximizing the reward  $r$ ). While making a decision of which action  $a_t$  to take along with the action inputs, the agent collaborates with the **Sampler**, another instantiated LLM, which is responsible for suggesting which frames to select. The sampler suggestions are based on the previous action-observation trajectory.

**Evaluator** We introduce an evaluator  $\mathcal{E}$ , which assesses the correctness of the prediction based on the question and the trajectory. It employs an error-feedback mechanism, iterating through trial-and-error to identify model errors. The evaluator  $\mathcal{E}$  receives the question  $Q_i$ , policy  $\pi_i$  and the trajectory  $\sqcup_j$  and makes a judgment whether the final answer made by the planner is valid or not. The evaluation is made in a binary style True or False with a confidence ranging from 0 to 100.

**Refiner** Once the evaluation is done, the refiner is given a question, policy, and the trajectory from the agent, and the evaluation to generate the refinement of the trajectory. Specifically, the refiner first generates diagnosis of the trajectory (e.g., if there is any redundant steps, or any actions or action input that can be refined). Then, it generates a refined plan. The refinement is generated regardless of the evaluation result. The reason behind this is that if the evaluation is correct, the refinement is stored along with the trajectory in the long-term memory to enhance the reasoning of future trials and if the evaluation if false, the refinement have direct purpose of refining the reasoning of the agent for the next trial.

**Long/Short Memory** We maintain the memory with Long-term memory  $\mathcal{M}_{\text{long}}$  to store experiences, short-term memory  $\mathcal{M}_{\text{short}}$  to store accessed frame information. This format allows us to utilize the long-term memory. When the The long-term memory  $\mathcal{M}_{\text{long}}$  is present, it is indexed by the question type based on their semantic similarity, retrieving the semantically similar experiences (question type, and the trajectories).

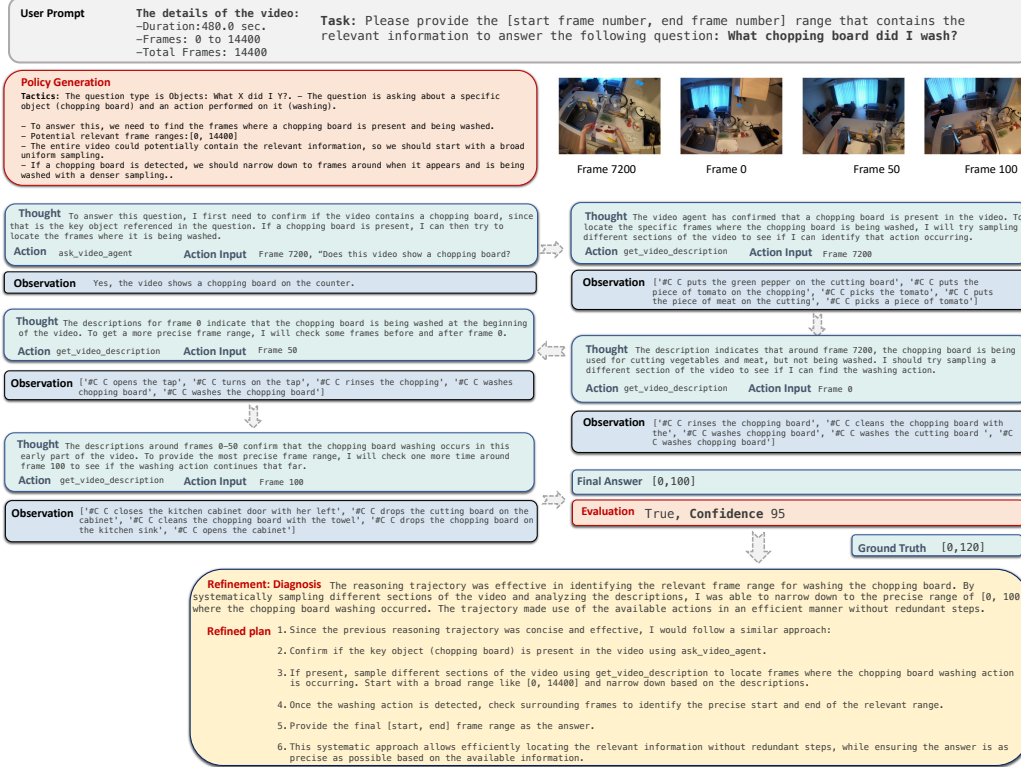


Figure 3: **Example of Ego4d NLQ Instance.** The User Prompt includes the video’s metadata and the question for the Agent to address. (1) Policy Generation: the agent generates an analysis of the question and a sampling strategy (2) Thoughts, Actions and Observation: The agent formulates a Thought based on current state, executes an Action  $\mathcal{A}$ , with Action Input, and uses tools to obtain an Observation  $\mathcal{O}$ . This process iterates until the agent comes up with the final answer. (3) Evaluation: the Final Answer is assessed. (4) Refinement: The trajectory  $\mathcal{T}$  is refined, and the results are stored in Long-term Memory  $\mathcal{M}_{\text{Long}}$ .

## 4 Experiments

### 4.1 Tools

Task	Source	Function
Video Caption Generation	LaViLa [47]	Detect actions, and objects
Video QA	Video-LLaVa [14]	Extract Information
Image QA	Claude 3 Sonnet [1]	Image description
Object Tracking	RT-DETR [46] ByteTrack [45]	Object detection
Text Caption	PaddleOCR [19]	Text caption Capture
Audio Transcription	Whisper [25]	Audio capturing

Table 2: **List of Invokable Tools.** This includes multi-modal tools, video-based tools (e.g., LaViLa, Video-LLaVa), image-based tools (e.g., Claude-3-Sonnet, PaddleOCR), and audio-based tools (e.g., Whisper)

In the experiments, the LLM used for reasoning and tool invocation is Claude-3-Sonnet [1]. The tools used in the framework are detailed in Table 2. The tools are chosen to support multi-modalities, such as video, image, or audio. The **Video Caption Generation** model, LaViLa [47], generates descriptions for selected frames. To accommodate the model’s requirement for frame sequences, we sample 3 additional frames (for a total of 4) for information extraction. Similarly, the **VideoQA** model, Video-LLaVa [14], samples 3 additional frames (totaling 4) for video frame analysis. The

**Object Tracking** model, RT-DETR [46], identifies objects with a confidence level above 0.6. The text caption tool [19] outputs text only if it is present in the frame.

## 4.2 Evaluation Datasets

Dataset	Task	Example	Avg duration	# Instances
Egoschema	Action and scene understanding, abstract reasoning	<p><b>Q: What is the overarching behavior of C and the man in the video?</b></p> <p>Option 0: C teaches the man game rules but the man seems distracted and is not paying attention</p> <p>....</p> <p>Option 5: The man shows C a new card game while C takes notes for future references</p> <p><b>A: Option 3</b></p>	3mins	0.5k
Ego4d NLQ	Temporal Localization	<p><b>Q: "What did I pick up before leaving the party?"</b></p> <p><b>A: [3410,4000]</b></p>	8.7mins	3.9k
MovieChat	Long-term video understanding	<p><b>Q: "When does the things in the video happens, ancient age, modern age or future?"</b></p> <p><b>A: "modern age"</b></p>	9.4mins	0.5k
NextQA	Causal and temporal action Reasoning	<p><b>Q: "Why was the toddler in red crying at the end of the video?"</b></p> <p><b>A: Fell backwards</b></p>	44secs	8.5k

Table 3: Overview of the evaluation Datasets. These benchmarks evaluate video understanding through a video question answering format, focusing on specific focus (denoted as Task). The average video duration varies from short form (<1min) to long form (<10min).

**EgoSchema** [18] comprises broad spectrum videos of daily human activities, three-minute egocentric video segments. Each question is associated with five possible answers, in multiple choice question answering format. To correctly answer the question, it requires long-term temporal understanding. In this paper, we use a subset of the Egoschema dataset, consisting of 500 question and answer pairs.

**Ego4D NLQ** [6] consists of egocentric videos capturing a diverse range of daily activities from individuals wearing cameras. The primary task involves temporal localizing relevant frames within these extensive video contexts (e.g. Where did I put X?). The task can be formalized, given a video  $\mathcal{V}$  and a natural language question  $Q$ , the goal is to identify a relevant frame window  $A$ , such that the answer to  $Q$  can be deduced from  $A$ . We utilize the validation set for evaluation. The average length of the video is around 8.7 minutes and the expected prediction time window is around 9.3 seconds.

**MovieChat** [30] encompasses a range of categories, including documentary and detective films. The benchmark involves questions such as identifying common objects, temporal elements (e.g., day, night), and various scenes through open-ended questions and answers. The average duration of the videos is 9.4 minutes. For our evaluation, we utilize the test set (Global mode) of this benchmark. As it involves open-ended questions, we utilized *Cluade-3.5-sonnet* as an evaluator to evaluate whether the prediction matches with the ground truth answer. To be rigorous, we made the evaluator to generate the confidence of its judgment, counting only the instances with confidence over 80 (out of 100) as correct.

**NextQA** [37] is a benchmark designed to assess various aspects of video understanding, including causal action reasoning, temporal action reasoning, and common scene comprehension. Compared to other evaluation benchmarks used in this study, NextQA focuses on relatively short video clips, with an average duration of 44 seconds. While it does not align with the long form video question-answer evaluation criteria, we include this benchmark to demonstrate the generalizability of our framework across short-form videos. Also, NextQA benchmark consists of questions with ‘textual cues’, for example, *Why was the toddler in red crying at the end of the video?*, it allows us to investigate the adaptive behavior of the agents when presented with questions with textual cues and without textual cues.

## 4.3 Baselines

We experiments with several strong baselines which are comprised of multiModal LLMs incorporating the visual components along with the textual queries as inputs. FrozenBiLM [38] learns cross modalities by training image projection layer. Similarly, InternVid [32] uses a image captioning

Model	# Frames	Accuracy
<b>MultiModalLLM</b>		
FrozenBiLM	90	26.9
InternVid	90	32.1
<b>ViT</b>		
ShortViViT	32	49.6
LongViViT	256	56.8
<b>Agent</b>		
LLoVi	180	57.6
VideoAgent	180	60.2
LifelongMemory	180	62.4
Ours	14.27	66.98
Total Avg Frames	5400 (30 fps)	

Table 4: **Egoschema** Results. The number of frames accessed and Accuracy.

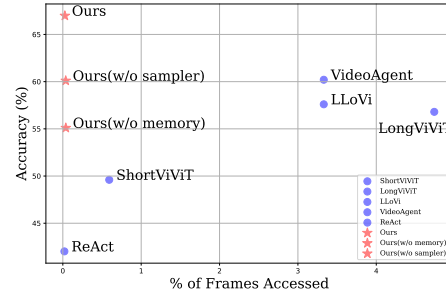


Figure 4: **Frames Accessed Ratio vs. Accuracy (%)** Our method demonstrate reduced % of frames accessed while maintaining high accuracy.

model along with transformer based text embeddings to align the image and the text. These methods work on fixed and limited number of frames. Vision transformer (ViT) based methods are based on vision transformer utilizing joint space time attention. ShortViViT and LongViViT [23] harness input masking strategy, supporting prefixed number of frames 32 frames and 256 number of frames respectively.

We also experiment with agent-based methods, which utilize language models as agents harnessing external tools to solve video question and answering task. LLoVi [43] extracts captions and LLM tackles the QA task based on the extracted captions. Analogously, LifelongMemory [33] process extracted captions and adopts voting by confidence strategy to conclude answers. VideoAgent [3] harness multiple tools to process video. These methods typically sample frames with predefined fps rate (e.g. 1fps).

## 5 Results

### 5.1 Performance Analysis

**Egoschema Results** Table 4 shows the results of evaluation on Egoschema benchmark. Our proposed method achieves accuracy of 66.98% which is more than 4% improvement over the best performing baseline method, LifelongMemory (62.4%). For the baselines, we also observe a trade-off between the number of frames accessed and accuracy. For instance, Multimodal LLMs [38] and [32] use a fixed sampling of 90 frames, but achieve relatively low accuracy (30%), whereas agent-based methods achieve significantly higher accuracy but sample twice as many frames. In contrast, our approach, which dynamically accesses relevant frames based on reasoning, reduces the number of frames accessed by approximately 93% while maintaining significant accuracy improvements. Existing methods typically use a uniform sampling strategy (1 frame per second), leading to a static number of frames. Our method avoids preprocessing all sub-sampled frames, thereby enhancing both accuracy and efficiency (Fig 4).

**Ego4d NLQ Results** We evaluate the intersection over union (IoU) at top-1 recall. (Table 5). Our method surpasses the baselines by 2% for IoU=0.3(%). Specifically, our method shows large improvement in IoU=0.5(%), which is around 10 % larger than the agent approach, and 11 % larger than the supervised approach. This may be attributed to the adaptive sampling strategy, which dynamically samples the frames, allowing both fine grained and coarse sampling. The frames are accessed on average 80% less than the agent method.

**MovieChat Results** Our method shows more than 22% increase in accuracy, while accessing only 0.1% of frames (Table 6), compared to the baseline models. This indicates that our method is more effective at processing long-form videos compared to both multimodal LLM-based (MovieChat [30]) and agent-based (VideoChatGPT, VideoLlama, VideoChat [17, 44, 12]) baselines.

		IoU=0.3(%) r@1	IoU=0.5(%) r@1	#Frames
Supervised	2D-TAN	5.04	3.12	1024
	VSLNet	5.45	6.63	1461
Agent	VideoAgent	17.38	7.47	avg 487(1fps)
	LifelongMemory	15.99	-	avg 487(1fps)
	Ours	19.5	<b>17.1</b>	avg 98 (0.002%)

Table 5: Ego4d NLQ Results.

	Accuracy	#Frames
VideoChat	57.8	32
VideoLlama	51.7	32
VideoChatGPT	47.6	100
MovieChat	62.3	2048
<b>Ours</b>	<b>84.8</b>	<b>13.59 (0.1%)</b>

Table 6: **MovieChat** Results.

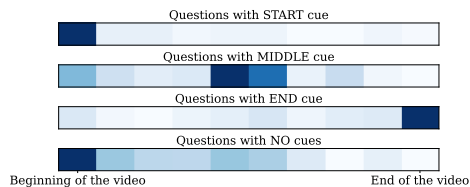


Figure 5: Frame accessed ratio based on textual cues from NextQA benchmark. Darker color corresponds to the higher ratio of access.

**NextQA Results** Our method shows a 1.4% improvement in overall average accuracy (Table 7). When analyzed by question type—temporal, causal, and descriptive—our method particularly excels in temporal reasoning tasks, providing around 6.9% absolute improvement over the next best method.

## 5.2 Ablation Analysis

**Agents Without Guidance are Suboptimal Reasoners** LLM agents using the default ReAct reasoning, without any intervention, exhibit suboptimal performance (Table 8, 9 ReAct). This approach results in both low accuracy and a reduced percentage of frames accessed. Although LLMs have the potential to examine all available frames and provide accurate answers, they often produce suboptimal results with fewer frames access. This is similar to observations where LLMs given one-shot questions demonstrate less rigorous reasoning compared to those using chain-of-thought or step-by-step interventions [35, 40]. Our framework, akin to the chain-of-thought method, enhances reasoning by incorporating internal interventions, leading to more accurate answers even if it requires accessing more frames.

**Questions including textual vs. non textual cues** Our proposed framework suggests that agents are query-adaptive, meaning they sample more efficiently when textual cues are present, as these cues guide their focus. For instance, a question like ‘Why was the toddler crying at the end of the video?’ will direct the agent to focus on the end of the video. The NextQA benchmark provides a natural testbed for evaluating whether agents leverage textual cues, as it includes both types of questions. Results indicate that the questions with textual cues result in an average of 10.56 frames accessed (.008%), compared to 12.26 (.01%) for questions without cues. Figure 8 presents a detailed analysis, showing that the ratio of frame accessed correlates with the presence of textual cues in the query. (e.g., a higher ratio of frames accessed at the beginning when ‘Start’ cues are included).

**Ablation of a component results in accuracy drop** A clear trend demonstrated across benchmarks (Table 8, 9) is that ablating any component consistently reduces accuracy. For Egoschema, the largest accuracy drop occurs when the evaluator is removed, while for Ego4D, the sampler’s removal has the greatest impact. Although accuracy trends are clear, the effect on the number of frames accessed is less consistent. For example, ablating the sampler or refiner generally increases frame access, whereas in Ego4D, it decreases. This indicates that the role of components like the sampler and refiner may vary with benchmark characteristics. Ego4D benefits from extensive frame search, while Egoschema needs a holistic video understanding. Thus, these components help balance frame access and accuracy depending on the benchmark’s requirements.

		Temporal	Causal	Descriptive	Average	# Frames (%)
Supervised	InternVid	43.4	48	65.1	49.1	19.92 (1.8%)
	SeViLA	61.3	61.5	75.6	63.6	39.85 (3.5%)
	MVU	55.4	48.1	64.1	55.2	39.85 (3.5%)
Agent	LLoVi	61	69.5	75.6	67.7	39.85 (3.5%)
	VideoAgent	64.5	72.7	81.1	71.3	39.85 (3.5%)
	<b>Ours</b>	<b>71.42</b>	69.1	77.77	<b>72.7</b>	12.37(1.1%)

Table 7: **NextQA results** The NextQA results are categorized by question types: temporal or causal reasoning and descriptive QA. Our method achieves a +1.4% higher accuracy compared to baseline methods, while accessing 2.4% fewer frames.

Model	Egoschema		Ego4d		
	# Frames(%)	Accuracy	# Frames (%)	IoU=0.3(%) r@1	IoU=0.5(%) r@1
ReAct	12.87 (.0024)	42.02	23.987(.00)	3.71	3.7
<b>Ours</b>	14.27 (.0026)	<b>66.98</b>	98 (.002)	<b>19.51</b>	<b>17.07</b>
-w/o memory	20.57 (.0038)	55.1	90.04 (.0022)	9.09	9.09
-w/o evaluator	15.69 (.003)	50.1	40.0 (.001)	5.41	4.69
-w/o sampler	19.77 (.0037)	60.1	55.67(.002)	5.01	5
-w/o refiner	20.46(0.003)	53.2	65.33(.001)	5.1	3.5

Table 8: **Ablation results on Frames Accessed and Accuracy.** The default ReAct model, with no interventions, exhibits the lowest accuracy and frame access ratio. Ablations of different components reveal varying trends in performance. The ablation results of Moviechat and NextQA can be found in Table 9

## 6 Conclusion

In this paper, we introduced a novel framework for video understanding that addresses the limitations of current methods by leveraging the daynamic reasoning capabilities of LLMs. While traditional approaches often rely on static or uniform frame sampling, which can be inefficient and redundant, our method enhances sampling efficiency by enabling the LLM based agent to adaptively select relevant frames based on specific queries. The results from extensive benchmarking validate the effectiveness and adaptability of our framework, showcasing its ability to handle diverse video understanding tasks more efficiently than traditional methods.

## 7 Limitations

While our method demonstrated effectiveness across several benchmark tasks, it is important to acknowledge its limitations. First, the performance of our approach is dependent on the capabilities of the tools it utilizes. Variations in tool performance can directly impact the overall effectiveness of the framework. Additionally, reliance on API calls introduces potential latency issues. This dependency on external APIs may affect the consistency and speed of the processing.

## 8 Broader impact

The proposed framework for video understanding presents several broader impacts with potential implications across various domains. By leveraging dynamic LLM-based agents for adaptive frame sampling, our approach offers a more efficient and effective solution to the challenges of long-form video comprehension. This advancement could significantly enhance applications in fields such as automated video content analysis, surveillance, and multimedia indexing, where processing large volumes of video data is essential.

## References

- [1] Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024.
- [2] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [3] Y. Fan, X. Ma, R. Wu, Y. Du, J. Li, Z. Gao, and Q. Li. Videoagent: A memory-augmented multimodal agent for video understanding. *arXiv preprint arXiv:2403.11481*, 2024.
- [4] D. Gao, L. Ji, L. Zhou, K. Q. Lin, J. Chen, Z. Fan, and M. Z. Shou. Assisgpt: A general multimodal assistant that can plan, execute, inspect, and learn. *arXiv preprint arXiv:2306.08640*, 2023.
- [5] D. Gao, L. Zhou, L. Ji, L. Zhu, Y. Yang, and M. Z. Shou. Mist: Multi-modal iterative spatial-temporal transformer for long-form video question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14773–14783, 2023.
- [6] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022.
- [7] W. Han, H. Chen, M.-Y. Kan, and S. Poria. Sas video-qa: Self-adaptive sampling for efficient video question-answering. *arXiv preprint arXiv:2307.04192*, 2023.
- [8] S. Heim and A. Keil. Developmental trajectories of regulating attentional selection over time. *Frontiers in Psychology*, 3:277, 2012.
- [9] S. Heim and A. Keil. Too much information, too little time: How the brain separates important from unimportant things in our fast-paced media world. *Frontiers for Young Minds*, 5(1), 2017.
- [10] E. Jang. Can llms critique and iterate on their own outputs? evjang. com. URL <https://evjang.com/2023/03/26/self-reflection.html>, 2023.
- [11] P. J. Lang, M. M. Bradley, and B. N. Cuthbert. Motivated attention: Affect, activation, and action. In *Attention and orienting*, pages 97–135. Psychology Press, 2013.
- [12] K. Li, Y. He, Y. Wang, Y. Li, W. Wang, P. Luo, Y. Wang, L. Wang, and Y. Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023.
- [13] Y. Li, X. Chen, B. Hu, and M. Zhang. Llms meet long video: Advancing long video comprehension with an interactive visual adapter in llms. *arXiv preprint arXiv:2402.13546*, 2024.
- [14] B. Lin, B. Zhu, Y. Ye, M. Ning, P. Jin, and L. Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- [15] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [16] S. Liu, H. Cheng, H. Liu, H. Zhang, F. Li, T. Ren, X. Zou, J. Yang, H. Su, J. Zhu, et al. Llava-plus: Learning to use tools for creating multimodal agents. *arXiv preprint arXiv:2311.05437*, 2023.
- [17] M. Maaz, H. Rasheed, S. Khan, and F. S. Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- [18] K. Mangalam, R. Akshulakov, and J. Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36, 2024.
- [19] PaddleOCR. Paddleocr, 2024.

- [20] V. Pallagani, B. Muppasani, K. Murugesan, F. Rossi, B. Srivastava, L. Horesh, F. Fabiano, and A. Loreggia. Understanding the capabilities of large language models for automated planning. *arXiv preprint arXiv:2305.16151*, 2023.
- [21] J. Pan, Z. Lin, Y. Ge, X. Zhu, R. Zhang, Y. Wang, Y. Qiao, and H. Li. Retrieving-to-answer: Zero-shot video question answering with frozen large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 272–283, 2023.
- [22] L. Pan, M. Saxon, W. Xu, D. Nathani, X. Wang, and W. Y. Wang. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*, 2023.
- [23] P. Papalampidi, S. Koppula, S. Pathak, J. Chiu, J. Heyward, V. Patraucean, J. Shen, A. Miech, A. Zisserman, and A. Nematzdeh. A simple recipe for contrastively pre-training video-first encoders beyond 16 frames. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14386–14397, 2024.
- [24] A. Piergiovanni, I. Noble, D. Kim, M. S. Ryoo, V. Gomes, and A. Angelova. Mirasol3b: A multimodal autoregressive model for time-aligned and contextual modalities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26804–26814, 2024.
- [25] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023.
- [26] J. E. Raymond, K. L. Shapiro, and K. M. Arnell. Temporary suppression of visual processing in an rsvp task: An attentional blink? *Journal of experimental psychology: Human perception and performance*, 18(3):849, 1992.
- [27] D. Romero and T. Solorio. Question-instructed visual descriptions for zero-shot video question answering. *arXiv preprint arXiv:2402.10698*, 2024.
- [28] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, and T. Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- [29] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [30] E. Song, W. Chai, G. Wang, Y. Zhang, H. Zhou, F. Wu, H. Chi, X. Guo, T. Ye, Y. Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024.
- [31] X. Wang, Y. Zhang, O. Zohar, and S. Yeung-Levy. Videoagent: Long-form video understanding with large language model as agent. *arXiv preprint arXiv:2403.10517*, 2024.
- [32] Y. Wang, Y. He, Y. Li, K. Li, J. Yu, X. Ma, X. Li, G. Chen, X. Chen, Y. Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2307.06942*, 2023.
- [33] Y. Wang, Y. Yang, and M. Ren. Lifelongmemory: Leveraging llms for answering queries in egocentric videos. *arXiv preprint arXiv:2312.05269*, 2023.
- [34] Z. Wang, S. Yu, E. Stengel-Eskin, J. Yoon, F. Cheng, G. Bertasius, and M. Bansal. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. *arXiv preprint arXiv:2405.19209*, 2024.
- [35] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

- [36] C.-Y. Wu, Y. Li, K. Mangalam, H. Fan, B. Xiong, J. Malik, and C. Feichtenhofer. Memvit: Memory-augmented multiscale vision transformer for efficient long-term video recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13587–13597, 2022.
- [37] J. Xiao, X. Shang, A. Yao, and T.-S. Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786, 2021.
- [38] A. Yang, A. Miech, J. Sivic, I. Laptev, and C. Schmid. Zero-shot video question answering via frozen bidirectional language models. *Advances in Neural Information Processing Systems*, 35:124–141, 2022.
- [39] Z. Yang, G. Chen, X. Li, W. Wang, and Y. Yang. Doraamongpt: Toward understanding dynamic scenes with large language models. *arXiv preprint arXiv:2401.08392*, 2024.
- [40] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [41] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [42] S. Yu, J. Cho, P. Yadav, and M. Bansal. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36, 2024.
- [43] C. Zhang, T. Lu, M. M. Islam, Z. Wang, S. Yu, M. Bansal, and G. Bertasius. A simple llm framework for long-range video question-answering. *arXiv preprint arXiv:2312.17235*, 2023.
- [44] H. Zhang, X. Li, and L. Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [45] Y. Zhang, P. Sun, Y. Jiang, D. Yu, F. Weng, Z. Yuan, P. Luo, W. Liu, and X. Wang. Bytetrack: Multi-object tracking by associating every detection box. In *European conference on computer vision*, pages 1–21. Springer, 2022.
- [46] Y. Zhao, W. Lv, S. Xu, J. Wei, G. Wang, Q. Dang, Y. Liu, and J. Chen. Detrs beat yolos on real-time object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16965–16974, 2024.
- [47] Y. Zhao, I. Misra, P. Krähenbühl, and R. Girdhar. Learning video representations from large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6586–6597, 2023.
- [48] Z. Zhao, W. S. Lee, and D. Hsu. Large language models as commonsense knowledge for large-scale task planning. *Advances in Neural Information Processing Systems*, 36, 2024.

# Appendix

## A. Case Study

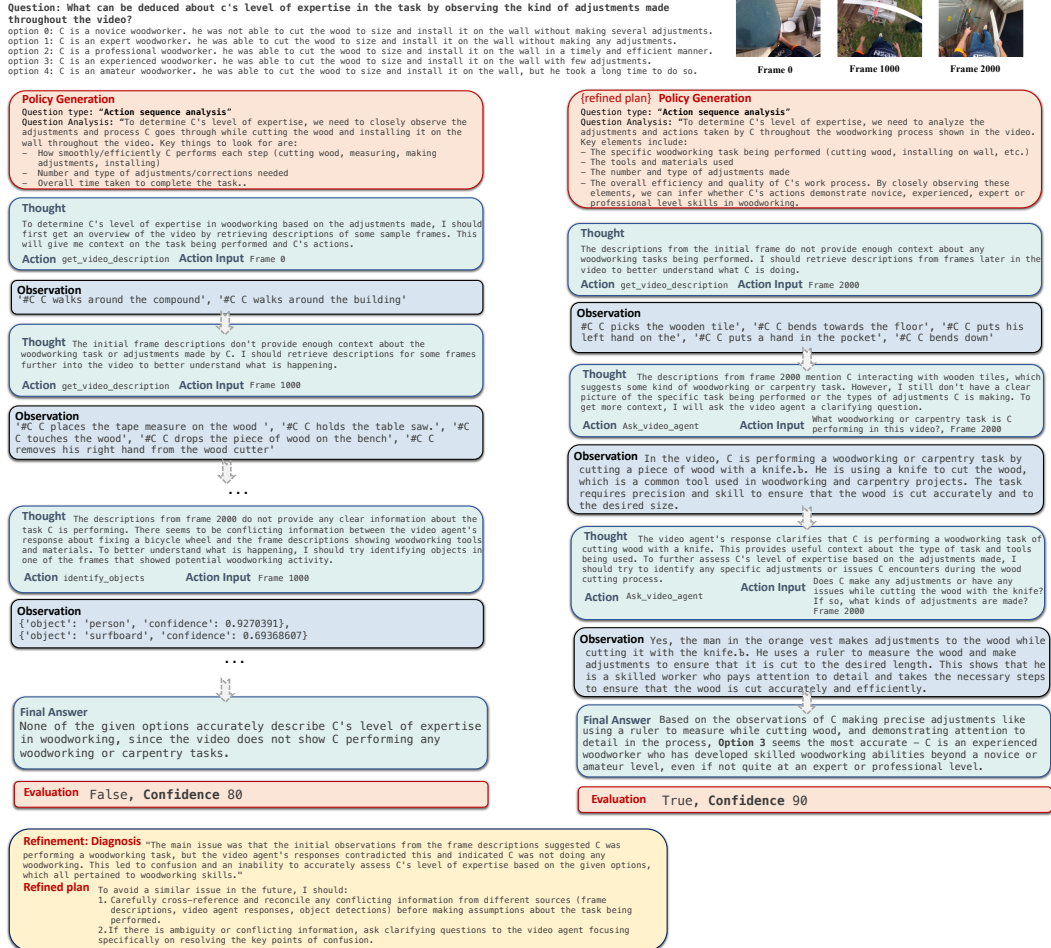


Figure 6: Example of Egoschema Refinement. Given the refinement based on the first trial (left), it attempts a second trial, with a refined policy  $\pi^*$ , which leads to the correct evaluation.

**Refined policy is more detailed and specific** What distinguishes the refined policy  $\pi^*$  from the initial generated policy? Figure 7 illustrates examples of refined policies, where the initial trial produced an incorrect prediction, while the second trial yielded a correct one. Compared to the original policy, the refined policy is notably more detailed. Specifically, it includes: 1) updates in question analysis, and 2) a more nuanced approach to sampling strategies. Although the sampling strategies in both trials were largely similar, the refined policy offers a more granular description. For instance, while the initial policy merely referenced dense and sparse sampling, the updated policy provides a comprehensive overview of when to transition between different sampling strategies.

**[Egoschema] Question:** What key transitional moments can be observed in the video where c switches from one activity to another? Explain how these transitions contribute to the overall video narrative.

**Policy Generation (1<sup>st</sup> trial)**

The question type is action sequence analysis. To answer this question, we need to identify the key transitional moments where the character switches from one activity to another in the video. We should sample frames uniformly across the video to detect these transitions. If a potential transition is detected, we can sample more densely around that region to confirm and analyze the transition in detail. The transitions and how they contribute to the overall narrative need to be explained.

**Policy Generation (2<sup>nd</sup> trial, after refinement)**

The question type is action sequence analysis. To identify the key transitional moments where the character switches activities and how they contribute to the overall narrative, I will:

1. Get a high-level overview of the main activity/process from the video agent to set the context.
2. Sample frames **densely across** the entire video duration using the `get_frame_info` tool. For each frame, analyze the specific **sub-activity happening (stirring, adding ingredients, checking progress, etc)**
3. Identify the transitional frames where the **sub-activity changes** from one to the next. These will be the "key transitional moments".
4. Describe how each identified transition contributes to the overall narrative **flow of the main cooking/food preparation process**.

---

**[Egoschema] Question:** What was the main dish prepared in the video, and which key ingredients were used?

**Policy Generation (1<sup>st</sup> trial)**

The question type is Key action/moment detection. To answer this question, we need to identify the key moments in the video where the main dish is being prepared and the ingredients being used are shown. We should sample frames densely around those key moments to observe the dish preparation process and ingredients used. For other parts of the video, a more sparse sampling may suffice.

**Policy Generation (2<sup>nd</sup> trial, after refinement)**

The question type is **Identifying main dish and key ingredients**. To answer this question, we need to:

1. Get an overview of the video content by sampling descriptions at different time points to understand if it involves cooking/preparing a dish.
2. **Analyze key frames** where active cooking/preparation seems to be happening to identify the specific dish being made and ingredients being used.
3. Keep an open mind to **different possibilities** like *burries, stir-fries, sautees* etc. based on the ingredients and cooking methods observed, instead of fixating on one dish type too early.
4. **Cross-reference the observed dish and ingredients against the multiple choice options** provided to select the best match.

The key tactics are:

- Sample video descriptions at different time points for an overview
- Analyze key cooking frames to identify dish type and ingredients
- Consider multiple dish possibilities based on evidence
- Map observations to the closest matching option without premature assumptions.

Figure 7: **Example of Refined Policy** Compared to the original policy, the refined policy is notably more detailed. Texts highlighted in yellow shows the added instruction.

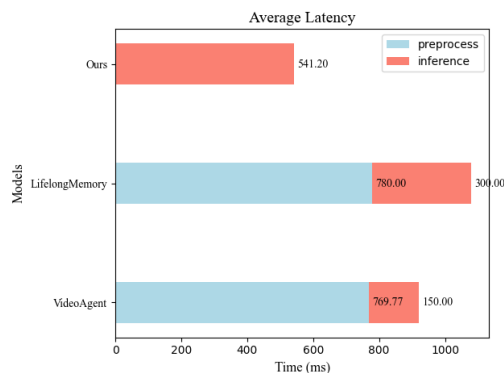


Figure 8: Latency Comparison with Other Agent Approaches: Our method reduces latency by processing videos only at runtime, compared to LifelongMemory [33] and Videoagent [3], which require preprocessing.

## B. Latency Analysis

## C. Prompts Configuration

**Policy Generation Prompt**

You are an advanced AI agent tasked with efficiently and accurately processing video question and answering tasks.

You will be given a question related to a video, and you are responsible for coming up with a set of tactics and plans based on the characteristics of each question. The questions you will encounter will vary greatly, ranging from inquires about the overall plot to specific details within the video.

To effectively handle these tasks, you must first generate a set of tactics and plans based on the characteristics of each question. You will be given a question, please analyze the question.

- Determine the type of question (e.g. purpose/goal identification, tools and materials usage, key action/moment detection, character interaction, action sequence analysis..etc)
- How should the frames be sampled to solve the question? (e.g. Uniform sampling with timestep 30. If relevant frame is detected, more uniform sample with timestep 2.)

{Question}  
{Video details}

### Agent Prompt

You are an advanced AI specialized in video question-answering tasks. Your capabilities include executing necessary tools and interpreting their outputs. Your objective is to select which frames to process and strategize which tools to deploy and use their outputs to provide accurate answers to questions related to a video.

<Video Details>:

- Duration: {duration\_min} minutes ({duration\_sec} seconds)
- Frame Rate: {frame\_rate} frame per second
- Total Frames: {total\_frames} frames.
- Frames with scene change: {scene\_list}

Among the total {total\_frames}, you will first choose sample frames to understand the context. Please use the tool 'get\_frame\_info' to get the general information of the frame. You can use the tools listed below. You can reason what's happening between frames, and what's described in the frame itself.

Use these tools to help: {tools}

Use the following format:

Thought: Consider what to do next.

Action: The action to take, using one of [{tool\_names}].

Action Input: The input for the action.

You will receive the result of the action as Observation: The result of the action. Please repeat the Thought/Action/Action Input/Observation cycle as needed.

The final answer should be provided under 'Final Answer:' You must choose one of the options among Option 0, Option 1, Option 2, Option 3, Option 4. Please start with Thought:

Begin!

### Refiner Prompt

You are an advanced reasoning agent that can improve based on self reflection. Your goal is to come up with a diagnoses and a refinement plan that is effective in making a correct prediction. You will be given a previous reasoning trial in which you were given access to execute tools to solve and an evaluation to the trial.

If the evaluation is False, you were unsuccessful in answering the question either because you guessed the wrong answer with Final Answer, or you used up your set number of reasoning steps. The optimal goal is to have concise reasoning path without having redundant actions. Even if the evaluation is True, you can improve the reasoning path by removing the redundant steps or by refining the repetitive actions. In a few sentences, Diagnose a possible reason for failure and devise a new, concise, reasoning paths that aims to mitigate the same failure. Be detailed as possible and use complete sentences.

### Evaluator Prompt

You are an advanced agent that evaluates whether the predicted answer is correct or not. You will be given a question, reasoning trajectories, and the final answer predicted by an agent. Please evaluate whether the prediction is valid or not. You can give your confidence in percentage (0-100). Remember that the reasoning and predictions are not always correct. For example, Evaluation: True, Confidence: 90

Table 9: Ablation result of Moviechat and NextQA. The results are consistent in exhibiting a drop in accuracy when ablating any component. However, the trends of the number of frames accessed are not consistent, varying across the benchmark, and the components

Model	MovieChat		NextQA	
	# Frames	Accuracy	# Frames	Accuracy
ReAct	10.62	69.4	9.87	47.27
Ours	13.59	<b>84.8</b>	12.37	<b>71.6</b>
-w/o memory	15.84	70.3	11.79	63.11
-w/o evaluator	15.31	72.43	13.59	58.97
-w/o sampler	15.51	80.2	13.87	60.2
-w/o refiner	15.51	70.2	11.8	65.42

Question: Where was the rim before I picked it?

**Sampler Suggestions**

```
> Start from the beginning of the video to identify the rim object and its initial location before being picked up <Sampler Suggestion: frame_index 0>

> Perform uniform sampling with a large timestep (7200 frames ~ 1/2 of total 14402 frames) to quickly identify if the rim object is visible around the middle portion of the video before being picked up <Sampler Suggestion: frame_index 7200>

> To identify the approximate range where the "rim" object appears before being picked up, I will sample at frame 10800 which is around 3/4 through the video. If the rim is not visible there, I can infer it likely appears in the final 1/4 of the video before being picked up.<Sampler Suggestion: frame_index 10800>

> To identify the approximate range where the rim first appears, I will sample frames closer to the start with a smaller timestep of 1800 (30 seconds at 60 fps). If the rim is visible in this frame, the relevant range likely starts around here. If not, I will continue sampling backwards from this point with an even smaller timestep. <Sampler Suggestion: frame_index 1800>

...

> To find the first appearance of the "rim" object before it is picked up, I should start watching the video sequentially from the very beginning (frame 0) since my previous sampling attempts were unable to identify when the rim first shows up. <Sampler Suggestion: frame_index 0>

> Since the rim is already visible in frame 0, I should sample a few frames ahead to see if the rim's position changes or if it gets picked up in the next few frames. Frame 30 seems like a reasonable next step to check for any movement or interaction with the rim. <Sampler Suggestion: frame_index 30>
```

Figure 9: **Sampler Example** The Sampler examples demonstrates that it is able to 1) calculate the frames in terms of sparse sampling 2) Dynamically switch sampling fps, based on previous observation 3) Densely sample relevant frames

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly states the paper's goal and contributions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: In section 7 we discuss the limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The focus of the paper is mainly on empirical experiments.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, in the section 4 and 3 section we elaborated the methods, and experimental settings.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: Yes. We use publicly accessible data. In terms of code, we plan to release the necessary code with instructions.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: Yes in the section 4 we describe the detailed information of the datasets, and the experiment settings.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.

- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: The experiments were conducted in a deterministic way (e.g. setting temperature to 0) as for reproducible purposes. The statistical significance is not applicable to our work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: In the section 4 we provide settings of the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Yes.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Yes in section 8 we discussed broader impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: There is no much risks involved in this work.

Guidelines:

- The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cited the paper or URL as long as the licenses and copyrights.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: New asset is not released in this work.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: There is no crowdsourcing involved in this work.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: There is no IRB Approval required in this work.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.