

TravelER: A Modular Multi-LMM Agent Framework for Video Question-Answering

Chuyi Shang* Amos You*
Sanjay Subramanian Trevor Darrell Roci Herzig

University of California, Berkeley

Abstract

Recently, image-based Large Multimodal Models (LMMs) have made significant progress in video question-answering (VideoQA) using a frame-wise approach by leveraging large-scale pretraining in a zero-shot manner. Nevertheless, these models need to be capable of finding relevant information, extracting it, and answering the question simultaneously. Currently, existing methods perform all of these steps in a single pass without being able to adapt if insufficient or incorrect information is collected. To overcome this, we introduce a modular multi-LMM agent framework based on several agents with different roles, instructed by a Planner agent that updates its instructions using shared feedback from the other agents. Specifically, we propose *TravelER*, a method that can create a plan to “**T**raverse” through the video, ask questions about individual frames to “**L**ocate” and store key information, and then “**E**valuate” if there is enough information to answer the question. Finally, if there is not enough information, our method is able to “**R**eplan” based on its collected knowledge. Through extensive experiments, we find that the proposed *TravelER* approach improves performance on several VideoQA benchmarks without the need to fine-tune on specific datasets. Our code is available at <https://github.com/traveler-framework/TravelER>.

1 Introduction

Over the last few years, Large Multimodal Models (LMMs) have demonstrated tremendous progress in the area of video understanding, particularly for the video question-answering (VideoQA) domain (Fu et al., 2023; Wang et al., 2024). More recently, LMMs have been able to achieve impressive results through video-based models (Lin et al., 2023a; Sun et al., 2022; Ye et al., 2023; Li et al., 2023). However, video models require a high level

of computational complexity to fine-tune, and annotations are difficult and expensive to collect. As a result, many recent approaches (Xue et al., 2023; Yu et al., 2023; Zhang et al., 2023a) operate on the frame level, leveraging large-scale image-based pretrained models in a zero-shot manner.

Moreover, these models may need to do several tasks simultaneously in a single step. In particular, they should identify the correct events in videos by understanding what information is relevant and ignoring irrelevant information. Next, they would need to extract specific and question-relevant visual details and use them to answer the question. As such, one iteration might not be enough to collect all the necessary information. For example, many current approaches use simple captioning, which is often too general to extract specific details, or they might miss important events. In these cases, current approaches cannot revisit the video to find additional information. Decomposing this process into different components allows each component to adapt to newly collected information.

To address this, we introduce a modular multi-LMM agent framework for VideoQA. A Planner agent instructs agents in charge of different tasks, such as navigating through the video, extracting visual information through a question-answering process, and reviewing this information to select an answer. Feedback from the agents is then passed on to each other and back to the Planner, who uses the feedback to update its instructions. In this way, we can revisit the video to narrow our focus or expand our search to more relevant information and extract specific details to answer the question.

Consider the example in Figure 1. Suppose we are asked why the boy turned over in the middle of the video. In the first iteration, our method uses temporal cues from the question to skip to the middle of the video and asks questions to find the relevant frames. In the next iteration, we gather more information. Asking about what the boy is doing,

*Equal contribution.

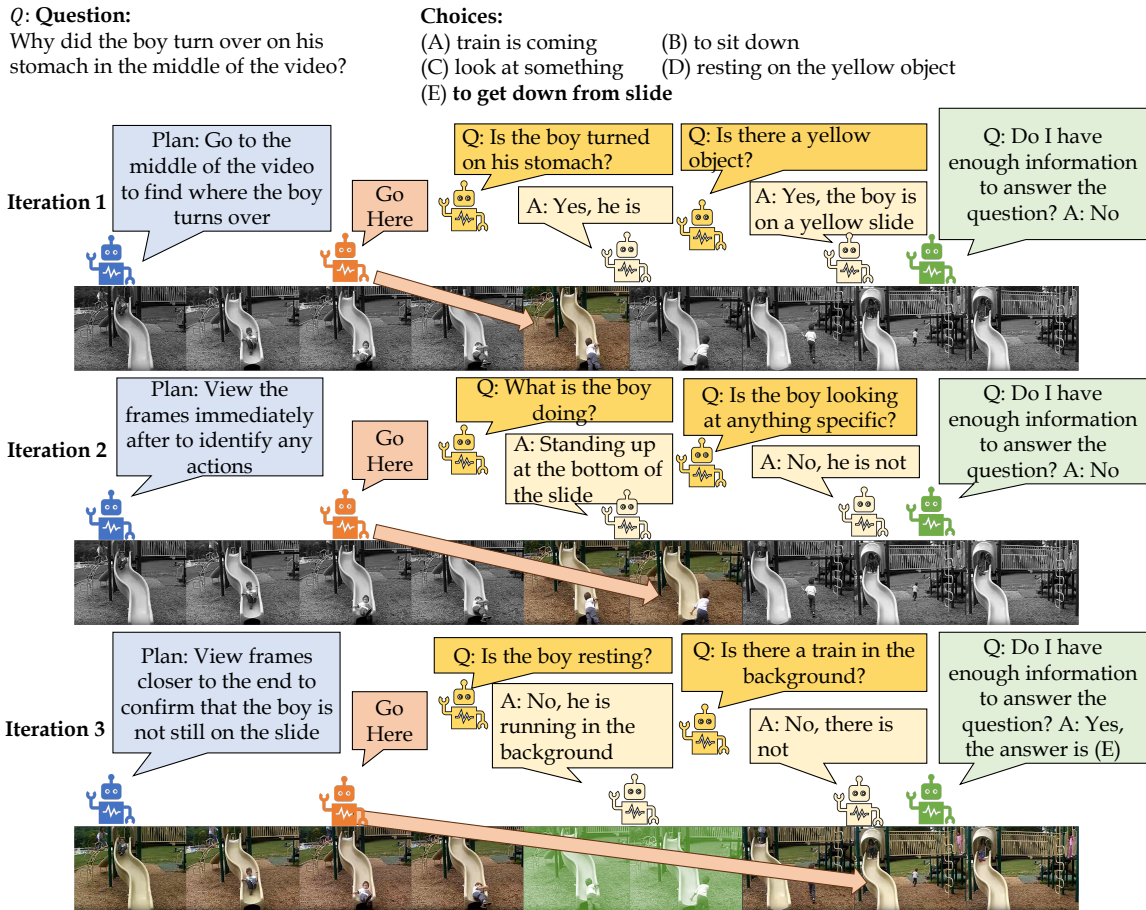


Figure 1: **A simplified overview of our TravelER framework.** Our proposed framework aims to answer the question by collecting relevant information from keyframes through interactive question-asking. To accomplish this, several agents (in colored boxes) with different roles interact (left-to-right in each row) over several iterations. TravelER creates a plan (in blue) to “traverse” (in orange) through the video, asks questions regarding individual frames (in yellow) to “locate” and store key information and, “evaluates” whether there is sufficient information to answer the question (in green), and “replans” using past collected knowledge if there is not enough information. **Click on the image to see the video.**

we learn that he is “standing up at the bottom of the slide” and is not looking at anything specific, which informs us that the boy is no longer “sitting down” (choice B) or “resting on the yellow object” (choice D). To eliminate these choices, we must confirm that the boy does not sit back down again by traveling to a timestamp near the end of the video. Finally, since we have collected enough information and followed the plans, we can select the right choice that the boy turns over to be on his stomach “to get down from slide” (Choice E).

Our proposed approach – **T**raverse, **L**ocate, **E**valuate, and **R**eplan (TravelER), has four main stages. First, in the Traversal stage, an agent creates a plan to answer the question. In the Location stage, an agent uses the plan to decide which timestamp of the video to select. The corresponding frames are then sent to another agent, which asks ques-

tions and stores the answers in a memory bank for future iterations. Finally, in the Evaluation stage, an agent reviews all collected information and decides whether to answer or create a modified plan (Replan) to start the next iteration if necessary.

To summarize, our main contributions are: (i) We introduce TravelER, a modular multi-LMM agent framework for video question-answering. (ii) Our method shows improved performance on multiple difficult video question-answering benchmarks, such as NExT-QA, EgoSchema, Perception Test, and STAR. (iii) Our method is easy to employ with different LLMs and LMMs, highlighting the effectiveness of our modular approach.

2 Related Work

Video Question-Answering. VideoQA involves answering free-form or multiple-choice questions

given an input video. Compared to image question answering, VideoQA poses unique challenges because it often requires strong temporal understanding and the ability to deal with long input sequences. Many recent works have focused on training end-to-end video-language models (Fu et al., 2023; Sun et al., 2022; Wang et al., 2024; Ye et al., 2023; Yu et al., 2022; Li et al., 2023), but doing so remains challenging due to computational constraints and difficulties in architecture scaling. As a result, many approaches adopt pretrained image models to the video domain by extracting information independently from each frame (Xue et al., 2023; Yu et al., 2023; Zhang et al., 2023a). Here, we design a framework that builds an adaptive plan to find and extract relevant information using a question-answering approach.

LMMs for Video Understanding. LMMs have been shown to be extremely useful for VideoQA. Some methods use supervised or contrastive training to perform video-LMM pretraining (Zhao et al., 2023; Yang et al., 2023; Chen et al., 2024), while others adapt existing LMMs and use instruction tuning to adapt them to the video domain (Zhang et al., 2023b; Maaz et al., 2024; Lin et al., 2023a). However, recent improvements in LMM capabilities have allowed for many strong approaches for few-shot (Alayrac et al., 2022; Wang et al., 2022) and zero-shot VideoQA (Yang et al., 2022). In particular, zero-shot methods, such as LLoVi (Zhang et al., 2023a), use pre-trained LMMs to generate captions for each frame in the video. Nevertheless, uniformly sampling frames at random may miss important visual information and focus on unimportant frames (Wu et al., 2019; Lei et al., 2021b). Recent works like SeViLA (Yu et al., 2023) addressed this problem by performing parameter-efficient finetuning using captions to identify keyframes (Lu et al., 2022; Buch et al., 2022; Qian et al., 2023), but this requires fine-tuning on specific datasets. Unlike these works, which select all keyframes in a single pass, we introduce a novel iterative and modular approach instructed by a planner.

LMM-based Agents for Videos. The strong reasoning abilities of LLMs (Brown et al., 2020; Chung et al., 2022) have made them effective in LLM-based agent approaches for videos, where an LLM performs much of the reasoning after collecting information from different modules (Chen et al., 2023; Lin et al., 2023b; Zhang et al., 2024; Zeng et al., 2023). For example, Socratic Models (Zeng et al., 2023) proposes a method to reason

about videos based on generated audio transcriptions and CLIP frame similarity scores, while other works like VideoChatCaptioner (Chen et al., 2023) proposes a way to caption videos through chat dialogues between an LLM and a LMM. Recently, there have also been works that use program generation using an LLM to answer questions (Choudhury et al., 2023; Min et al., 2024). However, these works still use a single-pass approach and provide very general captions. In contrast, our work uses an iterative question-answering process to extract specific, relevant information in the frame.

3 TravelER Framework

We begin by describing some background on the LLM and LMM architectures (Section 3.1), then introduce each component of our framework (Section 3.2) and implementation details (Section 3.3). Our method is shown in Figure 2.

3.1 Preliminaries

Large language and multimodal models. LLMs are text-conditioned generative models. Given a prompt P , they encode it into a fixed language embedding l in an embedding space $f(\cdot)$ and use this to produce text response R : $R = f(l(P))$. Similarly, Large Multimodal Models (LMMs) are adapted to jointly reason over vision and language. To map different modalities into the shared embedding space $f(\cdot)$, an image I is encoded using an encoder v , and the prompt P is encoded using a fixed language embedding l . The LMM outputs a textual response R : $R = f(v(I), l(P))$.

Video question-answering. VideoQA involves viewing a video and answering questions. The model is usually evaluated through *top-1* accuracy, which chooses the best answer out of a set of possible choices. Specifically, given a question Q , video input V consisting of a set of frames $\{I_1, \dots, I_n\}$, and a set of choices $C = \{c_1, \dots, c_n\}$, the model is asked to choose the best c_i to answer Q . Next, we introduce each component of our method.

3.2 TravelER Components

Traversal. In the Traversal stage, we create a plan for how to traverse through the video, which is a list of textual instructions that guide our approach to answering the question. To achieve this, we use the task prompt P_T , which is an instruction to create a plan for answering the question. We combine P_T with the question Q , and memory bank M , which

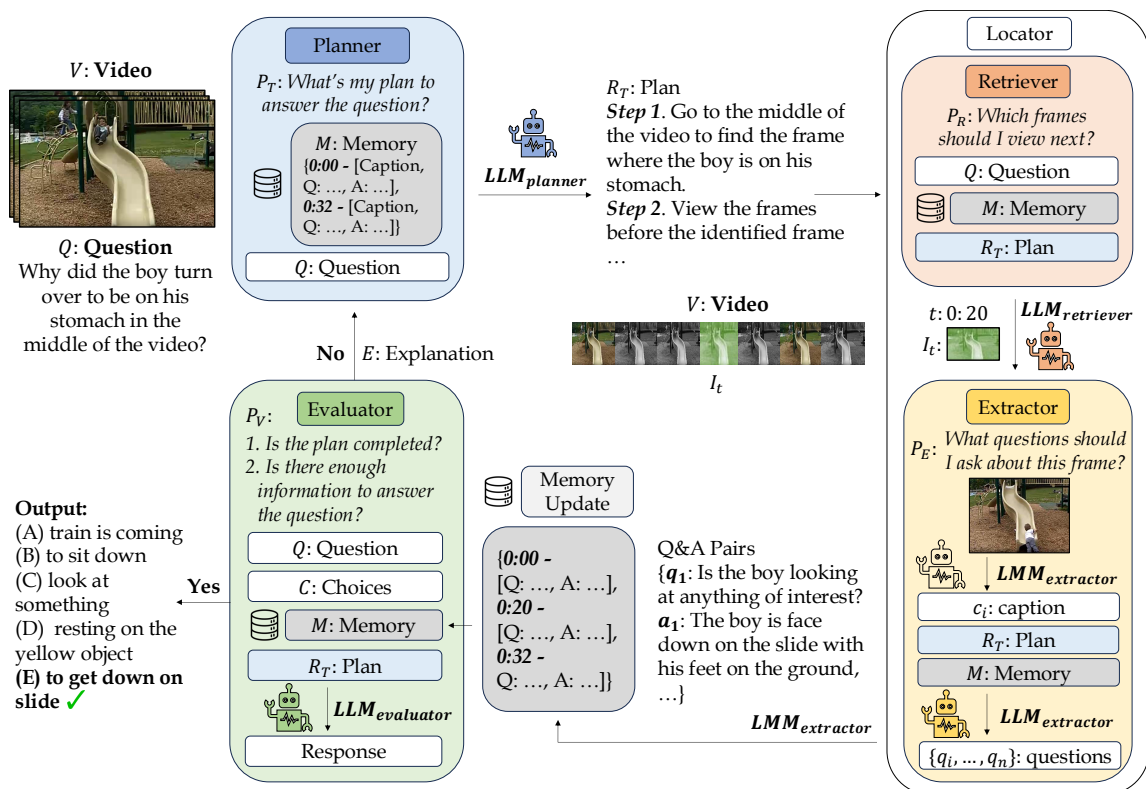


Figure 2: **TraveLER framework.** Our framework consists of four different modules, the Planner, Retriever, Extractor, and Evaluator. The Planner creates a plan and sends it to the Retriever. The Retriever uses the plan to select the next timestamp and sends this to the Extractor. The Extractor captions and generates questions about the timestamp, answers the questions, and saves the output in the memory bank. Finally, the Evaluator determines if there is enough information and if the plan has been followed. If yes, the Evaluator returns the answer, else the existing information is sent back to the Planner to begin a new iteration.

is a dictionary of collected information keyed by timestamps and containing information from the corresponding frame, to receive the final prompt $P_T^{(1)}$: $P_T^{(1)} = "[Q][M][P_T]"$.

Our method uses a memory bank M to store collected information, which allows information to persist and to be updated as we proceed through different iterations. We initialize M with captions of 5 evenly sampled frames throughout the video. We find that this memory initialization gives the model good context for the general idea of the video, and performs better than starting with an empty memory M . After the first iteration, we add information iteratively using the Extractor module, which we discuss later in this section.

Next, we input the prompt $P_T^{(1)}$ into LLM_{planner} , which returns response R_T , a step-by-step plan on how to traverse through the video and what information should be collected. This plan is revised in future iterations using collected information.

$$R_T = LLM_{\text{planner}}(l(P_T^{(1)}))$$

Our next step is to use the plan R_T in the Locator stage to locate relevant events and extract the information that we will use to answer the question.

Locator. The Locator is a component that consists of two submodules, the *Retriever* and the *Extractor*. The Retriever selects the timestamps of the next frames to view, while the Extractor extracts relevant information from these frames using a question-answering process. Next, we discuss each component in more detail.

(i) *Retriever:* The Retriever carries out the given plan R_T by selecting which frames to view. The Retriever is an LLM-based submodule that finds the next timestamp t to view, given the plan R_T , collected information M , and video metadata (frame rate, length). The task prompt P_R is an instruction that contains information about the video length and asks which timestamp to view next. Thus, we insert the question Q , the plan R_T , and the collected information M into the task prompt P_R to create the new prompt $P_R^{(1)}$: $P_R^{(1)} = "[P_R][Q][R_T][M]"$.

Given prompt $P_R^{(1)}$, the LLM in the Retriever, $\text{LLM}_{\text{retriever}}$, returns t , the next set of timestamps. Then, it retrieves frames I_t at timestamp t .

(ii) *Extractor*: The Extractor is important because it allows us to capture more relevant and question-specific details from the visual input, unlike using only captions. We pass the frames selected by the retriever I_t into the Extractor submodule, which consists of two large models: $\text{LLM}_{\text{extractor}}$, to generate context-dependent questions about the frames I_t , and a different vision-language model $\text{LMM}_{\text{extractor}}$, which extracts the desired information from the same frames. We note that we use both an LLM and LMM since the LLM is better at reasoning about the plan and the collected information, while the LMM is able to collect visual information requested by the LLM.

We first generate a general caption c_t for frame I_t using the $\text{LMM}_{\text{extractor}}$. Then, we concatenate the caption c_t , plan R_T , and memory M , and the Extractor task prompt P_E , which is an instruction that asks to use available information to create 3 questions to ask about the current frame. This results in the new prompt $P_E^{(1)}$: $P_E^{(1)} = "[P_E][c_t][R_T][M]"$.

Next, we input this new prompt $P_E^{(1)}$ into the LLM to get a set of questions $\{q_1, q_2, \dots, q_n\}$ about each frame, where n is a parameter for how many questions to ask about each frame.

$$\{q_1, q_2, \dots, q_n\} = \text{LLM}_{\text{extractor}}(l(P_E^{(1)}))$$

where l is the fixed language embedding.

In this way, the generated questions take into account both the plan R_T and information from past and future frames of the video M . We then use the frame I_t , and the corresponding questions $\{q_1, q_2, \dots, q_n\}$ as input into $\text{LMM}_{\text{extractor}}$. The $\text{LMM}_{\text{extractor}}$ then outputs a set of answers $\{a_1, a_2, \dots, a_n\}$, where each answer a_i corresponds to the question q_i .

$$\{a_1, \dots, a_n\} = \text{LMM}_{\text{extractor}}(v(I_t), l(\{q_1, \dots, q_n\}))$$

where v is the visual encoder.

Finally, to use this collected information in future iterations, we update our memory bank M . To do this, we use the timestamp t of I_t as our key and the question-answer pair list as the value, and append this to our memory M . If the memory bank dictionary M is too long, we summarize it by using the memory bank as an input to another LLM and instruct it to make the memory bank entries more

concise, while retaining the same keys and format. This output becomes our new memory bank.

Evaluator. The Evaluator decides if there is enough information and determines if the plan has been followed. We concatenate the memory information M , the plan R_T , the question Q , and the choices C with the task prompt P_V . The task prompt P_V is an instruction to evaluate if there is enough information to answer the question and if the given plan R_T has been fulfilled. Thus, we get the new prompt $P_V^{(1)}$: $P_V^{(1)} = "[P_V][Q][C][R_T][M]"$.

We use this prompt $P_V^{(1)}$ as input into the LLM in the Evaluator, $\text{LLM}_{\text{evaluator}}$, which evaluates if there is enough information to answer the question and if the plan has been completely followed. If both are true, $\text{LLM}_{\text{evaluator}}$ outputs the best choice c^* to answer the question Q . Otherwise, it provides an explanation E on why there is insufficient information and gives this explanation to the Planner to start a new iteration of the process.

Re-planning. After each iteration, if the evaluator decides that there is not enough information to answer the question Q or if the plan P has not been completed, the existing memory M will be provided to the Planner in the next iteration, in addition to an explanation E for why an answer was not chosen. The Planner then outputs a new plan, restarting the process. We also set a limit on the number of iterations a question can take to prevent infinite loops. After reaching this limit, we force the Evaluator to choose the best choice.

Summarizer. The Summarizer is an optional module used for some datasets to summarize the information. It is given the question Q , choices C , memory bank M , and task prompt P_S , which is an instruction to summarize the information in each timestamp of the memory bank M . The outputted summary is then used to replace the original memory bank M . This process helps reduce the amount of information passed to the LLM in future steps, which may struggle with very long inputs.

3.3 Implementation Details

Here, we discuss how we implement various components of our framework. More implementation details, such as prompts and dataset-specific details are in the Supplementary in Section B.

Memory bank. We represent past collected information as a Python dictionary, with the timestamp of different frames as keys and a list of extracted

Table 1: **Results on Datasets.** We show zero-shot results on different datasets. For fair comparisons, we gray out methods with fine-tuned components in their model.

(a) NExT-QA					(b) EgoSchema (Full)		(c) STAR		(d) Perception Test	
Model	Cau.	Tem.	Des.	Avg.	Model	Acc.	Model	Avg.	Model	Acc.
SeViLA	61.3	61.5	75.6	63.6	mPLUG-Owl	31.1	SeViLA	44.6	SeViLA	46.2
MC-ViT-L	-	-	-	65.0	InternVideo	32.1	Flamingo-9B	41.8	MC-ViT-B	47.0
InternVideo	43.4	48.0	65.1	49.1	LongViViT	33.3	InternVideo	41.6	MC-ViT-L	48.1
ViperGPT	-	-	-	60.0	Vamos	48.3	BLIP-2 ^{voting}	40.3	Flamingo-3B	43.6
ProViQ	-	-	-	63.8	LLoVi	50.3	BLIP-2 ^{concat}	42.2	LongViViT	45.7
LLoVi	69.5	61.0	75.6	67.7	TravelER	53.3 (+3.0)	TravelER	44.9 (+2.7)	TravelER	50.2 (+4.5)
TravelER	70.0	60.5	78.2	68.2 (+0.5)						

information from the frame as the values. This extracted information consists of a brief caption of the frame and a list of question-answer pairs. To prevent the memory bank from becoming too large, we also implement a summarizer module that instructs an LLM to summarize the memory bank and return a more concise version in the same dictionary format as before.

Agent model selection. Our modular approach has the benefit of allowing us to easily swap in different LLMs and LMMs (see Section 4.4). For our main experiments, we use LLaVA-1.6 (Liu et al., 2023) for LLM_{extractor} and GPT-3.5/GPT-4 (OpenAI et al., 2023) for LLM_{planner}, LLM_{retriever}, LLM_{extractor}, and LLM_{evaluator}.

Multi-frame selection. We also allow the Retriever to select multiple frames instead of one. This helps to capture better events that happen quickly or require more context to recognize. For example, if we want to find the action of "a woman clapping her hands", a single frame selection may cause us to incorrectly assume the woman is not clapping if we view the frame where their hands are apart. We do this by creating an optional parameter called window size. The window size refers to the number of frames the Retriever extracts each time. When the window size is non-zero, the Retriever still specifies a single timestamp to go to, but when retrieving the frame at that timestamp, we also take the number of frames specified by the window size before and after the selected frame.

4 Evaluation

We evaluated our TravelER framework on several benchmarks described in Section 4.1, and compared it to multiple baselines in Section 4.2. The results and ablations are in Section 4.3 and Sec-

tion 4.4. Additional results and ablations are in the Supplementary in Section A.

4.1 Datasets

We use the following datasets: (1) **NExT-QA (Xiao et al., 2021)** is a dataset that tests causal action reasoning and temporal understanding. Following the trend of works before us, we evaluate our method on the 5,000 questions in the NExT-QA validation set. (2) **EgoSchema (Mangalam et al., 2023)** is a challenging dataset that tests long-form video understanding. Viewers need to view 100 seconds of the video on average to answer the question correctly. (3) **STAR (Wu et al., 2021)** tests reasoning in real-world video situations. (4) **Perception Test (Patraucean et al., 2023)** is a challenging dataset that focuses on skills such as memory, abstraction, physics, and semantics and is intended to be approached in a few-shot or zero-shot manner.

4.2 Baselines

In our experiments, we compare our method to recent state-of-the-art zero-shot (ZS) methods, such as LLoVi (Zhang et al., 2023a), ProViQ (Choudhury et al., 2023), and other methods that are not necessarily ZS, such as SeViLA (Yu et al., 2023), and MC-ViT (Balazevic et al., 2024). We note that SeViLA uses fine-tuned components on QV-Highlights (Lei et al., 2021a), while MC-ViT is fine-tuned on NExT-QA for Perception Test. Additional baselines are in Supplementary in Section A.2.

4.3 Results

Our results are shown in Table 1. We use GPT-4 for NExT-QA to ensure a fair comparison with LLoVi, which is the current state-of-the-art that uses GPT-4 to uniformly caption frames across the entire video. Interestingly, our method outperforms

LLoVi despite viewing 50% fewer frames on average. Second, we also outperform SeViLA by +4.6%, although SeViLA uses a keyframe selector that is fine-tuned on a video moment retrieval and grounding task while our method is fully ZS.

For EgoSchema, we use GPT-4 with a re-trained version of LaViLA that excludes overlapping Ego4D and EgoSchema videos to prevent data leakage. We show strong performance on long-form videos, where we outperform LLoVi by +3.0%, while viewing 95% fewer clips on average.

We use GPT-3.5 for Perception Test and STAR because it is cheaper than GPT-4, but results are likely to be improved even further with GPT-4. Nevertheless, we achieve higher accuracy than LongViViT on Perception Test by +4.5% and MC-ViT by +2.1%, although it was fine-tuned on NExT-QA. We surpass both the best zero-shot approach by +2.7% and the best fine-tuned result by +0.3%.

Finally, please refer to Section A for more experimental and ablation results, and Section D for additional visualizations. For example, in Figure 8, we see that question-answering is able to extract more relevant details in comparison to simple captioning. This may explain why our method significantly outperforms the descriptive split of NExT-QA. In Figure 7, we see our method is able to reason about vague references and correctly identify relevant objects through question-answering.

4.4 Agent Ablations

We perform ablations using 1000 randomly selected questions from the NExT-QA training set (see Table 2). Unless specified, we use GPT-3.5/LLaVA-1.6 as the LLM/LMM for all agents.

Ablating the Planner. The Planner module outputs a plan, a list of instructions that guides the behavior of all other modules. We test the impact of removing it from our framework and find that it is worse. We hypothesize this is because the Planner provides many temporal cues that guide the Retriever module’s search, such as “go to the middle of the video”, and without these cues, the Retriever is not as good at selecting the next timestamp. Moreover, the Planner also helps the Evaluator better decide when to stop since in our iterative approach, the Evaluator uses the plan to determine when to stop. We also try removing re-planning by only running the Planner once in the beginning and keeping this plan fixed throughout. We find that this also reduces performance, showing the importance of adjusting plans to new information.

Table 2: Ablations on the agents in our framework.

Agent	Ablation	Avg. (%)
—	Baseline	60.4
Planner	Removal	58.1
	Fixed Plan	58.1
Retriever	Removal	56.9
Extractor	Caption Only	58.2
Evaluator	No Iteration	56.8
Summarizer	Removal	57.2

Ablating the Retriever. The Retriever module determines the next timestamps to view, which helps focus our information collection. We ablate it by uniformly sampling frames from the video at 2-second intervals, similar to other methods like LLoVi, which performs worse. We believe this is because the Retriever allows us to capture frames that might have otherwise been skipped through uniform sampling and selects fewer unimportant frames that might mislead the model.

Ablating the Extractor. Question asking is important as it allows us to capture more fine-grained and question-relevant information compared to simple caption generation, which produces a generic description. Thus, we ablate the Extractor by only allowing the LMM to caption frames. We find that this decreases performance by -2.2%, suggesting that the ability to ask specific questions about a frame is important. We notice many generated captions capture the main idea of visual information in the frame, but are lacking in fine-grained details.

Ablating the Evaluator. The Evaluator reflects on the collected information and decides if there is enough information to answer the question. To examine the impact of this reflection process, we make the Evaluator answer on the first iteration, finding it results in a -3.6% performance decrease.

Ablating the Summarizer. When collecting large amounts of information from videos, we use a Summarizer to condense the information, since long inputs can be challenging for LLMs. This has also been observed in recent work (Zhang et al., 2023a). To understand the impact of the Summarizer, we remove it. The results indicate that removing it degrades performance by -3.2%, demonstrating the advantage of more concise information.

4.5 Additional Experiments

Substituting different LLMs/LMMs. To see how the choice of the LLM and LMM affects our framework’s performance, we swap different LLMs and

Table 3: **Ablation Results.** We perform ablations on 1000 randomly selected questions from the NExT-QA training set. We report (a), (b) replacing different LLMs and LMMs, (c) selecting different numbers of frames to view in the Retriever, and (d) changing the number of questions asked in the Extractor. We use GPT-3.5/LLaVA-1.6 as the LLM/LMM, 5 frames for the Retriever, and 3 questions for the Extractor, unless otherwise specified.

(a) Replacing Diff. LLMs			(b) Replacing Diff. LMMs			(c) # of Frames		(d) # of Questions	
LLM	LMM	Accuracy	LLM	LMM	Accuracy	# Frames	Accuracy	# Questions	Accuracy
GPT-3.5		60.4	GPT-4	GPT-4V	64.7	1	59.0	0	58.1
Llama 3	LLaVA-1.6	63.9		BLIP-2	52.7	3	57.9	1	58.4
GPT-4		65.8	GPT-3.5	GPT-4V	59.5	5	60.4	3	60.4
GPT-4o		68.0		LLaVA-1.6	60.4	7	59.0	5	59.6

LMMs into our framework (see Table 3a and Table 3b). We first try different LLMs while fixing the LMM to be LLaVA-1.6. While GPT-4 performs better than GPT-3.5 by a significant margin of +5.4%, open-source model Llama 3 is very close (-1.9%), while incurring no additional cost. We also evaluate the newly released GPT-4o, which outperforms GPT-4 by +2.2% while being 61% cheaper, showing that our method can leverage better future models. Second, we use different LMMs while fixing the LLM to be GPT-3.5. We find that LLaVA-1.6 does best, GPT-4V is slightly worse (-0.9%), and BLIP-2 is significantly worse (-7.7%). Finally, we run an experiment using GPT-4V as both the LLM and LMM, and find that this does worse than GPT-4 and LLaVA-1.6 by -1.1%.

GPT-4V Baseline. To get a baseline using GPT-4V, we use a method similar to LLoVi using a subset of 500 examples. We use GPT-4V to caption the video uniformly, then ask it to answer the question given the choices, which results in a performance change of -2.0%. We note our motivation is to refrain from captioning every single frame, instead finding frames that help us answer the question.

Retriever window size. We experiment with different window sizes in Table 3c which is the number of frames the Retriever extracts centered around the selected frame. This allows us to capture better actions that occur quickly or require more context to understand. We find that choosing 5 frames yields the best results and a +1.4% increase when compared to selecting a single frame, but viewing more than 5 decreases performance. This suggests retrieving multiple frames can help the model better capture relevant information, but retrieving too many frames can lead to too much information.

Memory Initialization. In order for the Planner to create effective plans, it is beneficial to initialize the memory bank properly. Memory initialization

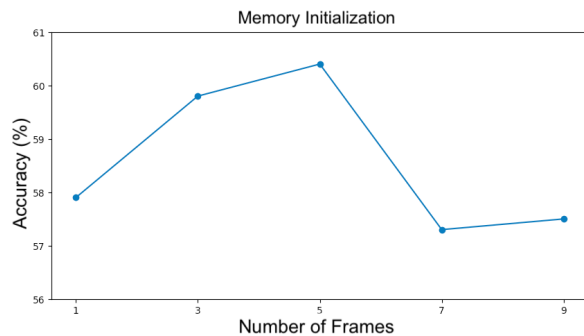


Figure 3: Comparison of different Memory Initialization (1, 3, 5 frames). 5 frames is optimal.

allows the Planner to have a high-level overview of the video, and create a corresponding plan on how to traverse the video given the initial frames. We perform three different initializations with 1, 3, 5, 7, and 9 frames and display our results in Figure 3. We observe that initializing the memory bank with 5 frames uniformly sampled from the video (0, 0.25, 0.5, 0.75, 1 for beginning, quarter, middle, three-quarter, end) yields the best result. In contrast, we notice a decrease in accuracy of -2.9% when using 9 frames, -3.1% when using 7 frames, -0.6% when initializing with 3 frames (0, 0.5, 1 for beginning, middle, end), and -2.5% when initializing with 1 frame (0.5 for middle).

Robustness to choosing incorrect frames. To see if our approach can recover from viewing incorrect frames, we choose random frames for the first 3 iterations (out of 5 total iterations) before using the Retriever. This leads to an accuracy drop of only 2.8%, showing our method is capable of recovering from viewing incorrect frames.

The number of questions. Question answering allows us to extract more specific details from our visual inputs. However, we noticed that too many questions can yield irrelevant questions and false positives. As such, we experiment with modifying

Table 4: Comparison of *TravelER* with other few-shot and zero-shot keyframe localization methods. For fair comparisons, we gray out methods with fine-tuned components in their model. The best scores are in **bold**.

Method	NExT-QA (Random Subset)			
	Temporal	Causal	Descriptive	Average
SeViLA - Localizer (Yu et al., 2023)	48.8	61.2	68.3	58.2
Moment-DETR (Lei et al., 2021a)	45.3	55.8	70.8	54.6
SigLIP (Zhai et al., 2023)	48.4	61.5	73.8	59.1
<i>TravelER</i> - Planner & Retriever (ours)	50.9	62.7	72.4	60.3

the number of questions asked for each frame by our extractor (see Table 3d). We record results for a 5-question, 3-question, and 1-question maximum. Note that 0-questions asked is equivalent to only allowing captions, which is discussed in the Extractor ablation. From our results, we find that a 3-question limit yields the best results compared to asking 1 or 5 questions (+2.0/+0.8%). This suggests that asking questions helps in extracting relevant information, but too many questions can lead to false positives or too much irrelevant information.

Memory bank initialization and formatting. Our memory bank M stores information that all modules rely on to make decisions. First, we experiment with different initializations as M must be initialized in the first iteration. We experiment with initialization of 1, 3, and 5 uniformly sampled captions. We find using 5 evenly spaced frames yields the best results, possibly because it starts the model with a general overview of the video before it starts to collect more relevant information. We also experiment with changing the format of the memory bank from JSON to markdown tables. We find this does worse by -2.9%, possibly because LLMs are better able to understand JSON formats.

Comparison with other keyframe selection methods. In Table 4, we compare our Planner and Retriever with other keyframe localization methods by replacing our Planner and Retriever with each of the other methods, and use our Extractor and Evaluator to perform the question answering. For all these methods, we use GPT-3.5 and LLaVA-1.6, and we evaluate these methods on a random subset of 1000 examples from the training set of NExT-QA. Note that other methods find keyframes in one inference iteration, whereas our inference occurs over multiple iterations. Therefore, to ensure fair comparisons, we uniformly sample 32 frames and extract out 4 keyframes in the other methods,

and we run 4 iterations of *TravelER* using the Retriever with a window size of 2 to similarly find 4 keyframes among fewer than 32 viewed frames.

We find that our Planner and Retriever surpasses other keyframe localization methods, despite considering fewer total frames (~ 25 total frames; we have 5 for memory initialization, and up to 5 frames each iteration). We would like to highlight that while our method is effective at finding keyframes, we do not need to find all keyframes to answer a question. Instead, we are often able to choose the correct answer with only a subset of the keyframes. **Cost of inference.** We test costs for GPT-3.5 and GPT-4, which have costs per inference of \$0.03 and \$0.67, respectively. We also test the open-source model Llama 3, which is free to run with only a minor performance decrease (-1.9%). Finally, we also try the new GPT-4o model, which is both better (+2.2%) and cheaper (61%) than GPT-4. We believe this trend will only continue as models become better and cheaper in the future.

5 Conclusion

We design a modular, multi-LMM agent framework for video-question answering based on several agents with different roles, instructed by a Planner agent that updates its instructions using shared feedback between the other agents. Our method creates a plan to “traverse” through the video, asking questions about individual frames to “locate” and store key information, and then “evaluate” if there is enough information to answer the question, “replanning” using new feedback if necessary. Through extensive experiments and ablations, we find that the proposed *TravelER* approach is not only easy to employ with different models but also improves performance on several video question-answering benchmarks without the need to fine-tune on specific datasets.

6 Limitations

In this work, we present a modular, zero-shot framework for video question answering (VideoQA) and demonstrate its effectiveness by improving on multiple state-of-the-art benchmarks. While TraveLER offers substantial benefits for VideoQA, it is important to recognize certain limitations that accompany our approach. Firstly, the effectiveness of our model relies heavily on the strength of the LLM and LMM. We notice that false-positives and incorrect statements from the LMM heavily impact performance. Our method’s runtime also depends on the runtime of current existing methods, and with the modularity of our method we expect this to improve with faster and better models in the future. Finally, we do not anticipate negative impacts of this work, but, as with any Machine Learning method, we recommend exercising caution.

7 Acknowledgements

We would like to thank Suzie Petryk, Chancharik Mitra, Alon Mendelson, David Chan, Assaf Arbelle, and Leonid Karlinsky for helpful feedback and discussions. This project was supported in part by DoD, including PTG and/or LwLL programs, as well as BAIR’s industrial alliance programs.

References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems*.
- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. [Neural module networks](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 39–48.
- Elad Ben Avraham, Roei Herzig, Karttikeya Mangalam, Amir Bar, Anna Rohrbach, Leonid Karlinsky, Trevor Darrell, and Amir Globerson. 2022. Bringing image scene structure to video via frame-clip consistency of object tokens. In *Thirty-Sixth Conference on Neural Information Processing Systems*.
- Ivana Balazevic, Yuge Shi, Pinelopi Papalampidi, Rahma Chaabouni, Skanda Koppula, and Olivier J Henaff. 2024. [Memory consolidation enables long-context video understanding](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 2527–2542. PMLR.
- Amir Bar, Roei Herzig, Xiaolong Wang, Anna Rohrbach, Gal Chechik, Trevor Darrell, and A. Globerson. 2021. Compositional video synthesis with action graphs. In *ICML*.
- Fabien Baradel, Natalia Neverova, Christian Wolf, Julien Mille, and Greg Mori. 2018. Object level visual reasoning in videos. In *European Conference on Computer Vision*, pages 105–121.
- Peter Battaglia, Jessica Blake Chandler Hamrick, Victor Bapst, Alvaro Sanchez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andy Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey Allen, Charles Nash, Victoria Jayne Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. 2018. [Relational inductive biases, deep learning, and graph networks](#). *arXiv*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Shyamal Buch, Cristobal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. 2022. Revisiting the “Video” in Video-Language Understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jun Chen, Deyao Zhu, Kilichbek Haydarov, Xiang Li, and Mohamed Elhoseiny. 2023. [Video chatcaptioner: Towards enriched spatiotemporal descriptions](#). Preprint, arXiv:2304.04227.
- Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. 2024. Vast: A vision-audio-subtitle-text omni-modality foundation model and dataset. *Advances in Neural Information Processing Systems*, 36.
- Rohan Choudhury, Koichiro Niinuma, Kris M. Kitani, and László A. Jeni. 2023. [Zero-shot video question answering with procedural programs](#). Preprint, arXiv:2312.00937.

- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#). *Preprint*, arXiv:2210.11416.
- Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. 2023. An empirical study of end-to-end video-language transformers with masked visual modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22898–22909.
- Jiaxin Ge, Sanjay Subramanian, Baifeng Shi, Roei Herzig, and Trevor Darrell. 2023. [Recursive visual programming](#). *Preprint*, arXiv:2312.02249.
- Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14953–14962.
- Roei Herzig, Ofir Abramovich, Elad Ben Avraham, Asaf Arbelle, Leonid Karlinsky, Ariel Shamir, Trevor Darrell, and Amir Globerson. 2024. Promptonomyvit: Multi-task prompt learning improves video transformers using synthetic scene data. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 6803–6815.
- Roei Herzig, Amir Bar, Huijuan Xu, Gal Chechik, Trevor Darrell, and Amir Globerson. 2020. Learning canonical representations for scene graph to image generation. In *European Conference on Computer Vision*.
- Roei Herzig, Elad Ben-Avraham, Karttikeya Mangalam, Amir Bar, Gal Chechik, Anna Rohrbach, Trevor Darrell, and Amir Globerson. 2022. Object-region video transformers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Roei Herzig, Elad Levi, Huijuan Xu, Hang Gao, Eli Brosh, Xiaolong Wang, Amir Globerson, and Trevor Darrell. 2019. Spatio-temporal action graph networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0.
- Roei Herzig, Alon Mendelson, Leonid Karlinsky, Asaf Arbelle, Rogerio Feris, Trevor Darrell, and Amir Globerson. 2023. Incorporating structured representations into pretrained vision & language models using scene graphs. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Roei Herzig, Moshiko Raboh, Gal Chechik, Jonathan Berant, and Amir Globerson. 2018. Mapping images to scene graphs with permutation-invariant structured prediction. In *Advances in Neural Information Processing Systems (NIPS)*.
- Achiya Jerbi, Roei Herzig, Jonathan Berant, Gal Chechik, and Amir Globerson. 2020. Learning object detection from captions via textual scene attributes. *ArXiv*, abs/2009.14558.
- Ranjay Krishna, Ines Chami, Michael S. Bernstein, and Li Fei-Fei. 2018. Referring relationships. *European Conference on Computer Vision*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Jie Lei, Tamara Lee Berg, and Mohit Bansal. 2021a. [Qvhighlights: Detecting moments and highlights in videos via natural language queries](#). In *Advances in Neural Information Processing Systems*.
- Jie Lei, Linjie Li, Luwei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. 2021b. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*.
- Jiangtong Li, Li Niu, and Liqing Zhang. 2022. From representation to reasoning: Towards both evidence and commonsense reasoning for video question-answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Linjie Li, Zhe Gan, Kevin Lin, Chung-Ching Lin, Ce Liu, Zicheng Liu, and Lijuan Wang. 2023. Lavender: Unifying video-language understanding as masked language modeling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. 2023a. [Video-llava: Learning united visual representation by alignment before projection](#). *Preprint*, arXiv:2311.10122.
- Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, Ce Liu, and Lijuan Wang. 2023b. [Mm-vid: Advancing video understanding with gpt-4v\(ision\)](#). *Preprint*, arXiv:2310.19773.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. [Visual instruction tuning](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Haoyu Lu, Mingyu Ding, Nanyi Fei, Yuqi Huo, and Zhiwu Lu. 2022. [LGDN: Language-guided denoising network for video-language modeling](#). In *Advances in Neural Information Processing Systems*.

- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*.
- Tegan Maharaj, Nicolas Ballas, Anna Rohrbach, Aaron Courville, and Christopher Pal. 2017. [A dataset and exploration of models for understanding video data through fill-in-the-blank question-answering](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7359–7368.
- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2023. [Egoschema: A diagnostic benchmark for very long-form video language understanding](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Juhong Min, Shyamal Buch, Arsha Nagrani, Minsu Cho, and Cordelia Schmid. 2024. Morevqa: Exploring modular reasoning models for video question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. 2024. Compositional chain of thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Twarek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Pinelopi Papalampidi, Skanda Koppula, Shreya Pathak, Justin Chiu, Joe Heyward, Viorica Patraucean, Jiajun Shen, Antoine Miech, Andrew Zisserman, and Aida Nematzdeh. 2024. A simple recipe for contrastively pre-training video-first encoders beyond 16 frames. In *Proceedings of the IEEE/CVF Conference on Com-*

- puter Vision and Pattern Recognition (CVPR), pages 14386–14397.
- Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens Contente, Larisa Markeeva, Dylan Sunil Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, Carl Doersch, Tatiana Matejovicova, Yury Sulsky, Antoine Miech, Alexandre Fréchet, Hanna Klimczak, Raphael Koster, Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osindero, Dima Damen, Andrew Zisserman, and Joao Carreira. 2023. [Perception test: A diagnostic benchmark for multimodal video models](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Tianwen Qian, Ran Cui, Jingjing Chen, Pai Peng, Xiaowei Guo, and Yu-Gang Jiang. 2023. [Locate before answering: Answer guided question localization for video question answering](#). *Trans. Multi.*, 26:4554–4563.
- Hung-Ting Su, Yulei Niu, Xudong Lin, Winston H. Hsu, and Shih-Fu Chang. 2023. [Language models are causal knowledge extractors for zero-shot video question answering](#). In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 4951–4960.
- Sanjay Subramanian, Medhini Narasimhan, Kushal Khangaonkar, Kevin Yang, Arsha Nagrani, Cordelia Schmid, Andy Zeng, Trevor Darrell, and Dan Klein. 2023. Modular visual question answering via code generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, Toronto, Canada. Association for Computational Linguistics.
- Yuchong Sun, Hongwei Xue, Ruihua Song, Bei Liu, Huan Yang, and Jianlong Fu. 2022. [Long-form video-language pre-training with multimodal temporal contrastive learning](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 38032–38045. Curran Associates, Inc.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. [Vipergpt: Visual inference via python execution for reasoning](#). *Proceedings of IEEE International Conference on Computer Vision (ICCV)*.
- Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. 2024. [Internvid: A large-scale video-text dataset for multimodal understanding and generation](#). In *The Twelfth International Conference on Learning Representations*.
- Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Derek Hoiem, Shih-Fu Chang, Mohit Bansal, and Heng Ji. 2022. [Language models with image descriptors are strong few-shot video-language learners](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 8483–8497. Curran Associates, Inc.
- Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. 2021. [Star: A benchmark for situated reasoning in real-world videos](#). In *Thirty-fifth Conference on Neural Information Processing Systems (NeurIPS)*.
- Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. 2019. [Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. [Next-qa: Next phase of question-answering to explaining temporal actions](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9777–9786.
- Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 2023. [CLIP-vip: Adapting pre-trained image-text model to video-language alignment](#). In *The Eleventh International Conference on Learning Representations*.
- Hongyang Xue, Zhou Zhao, and Deng Cai. 2017. [Unifying the video and question attentions for open-ended video question answering](#). *IEEE Transactions on Image Processing*, 26(12):5656–5666.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2021. [Just ask: Learning to answer questions from millions of narrated videos](#). In *ICCV*.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. [Zero-shot video question answering via frozen bidirectional language models](#). In *NeurIPS*.
- Antoine Yang, Arsha Nagrani, Paul Hongsuck Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptev, Josef Sivic, and Cordelia Schmid. 2023. [Vid2seq: Large-scale pretraining of a visual language model for dense video captioning](#). In *CVPR*.
- Qinghao Ye, Guohai Xu, Ming Yan, Haiyang Xu, Qi Qian, Ji Zhang, and Fei Huang. 2023. [Hitea: Hierarchical temporal-aware video-language pre-training](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15405–15416.
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. [Coca: Contrastive captioners are image-text foundation models](#). *Transactions on Machine Learning Research*.
- Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. 2023. [Self-chained image-language model for video localization and question answering](#). In *NeurIPS*.

- Andy Zeng, Maria Attarian, brian ichter, Krzysztof Marcin Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael S Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. 2023. [Socratic models: Composing zero-shot multimodal reasoning with language](#). In *The Eleventh International Conference on Learning Representations*.
- Kuo-Hao Zeng, Tseng-Hung Chen, Ching-Yao Chuang, Yuan-Hong Liao, Juan Carlos Niebles, and Min Sun. 2017. [Leveraging video descriptions to learn video question answering](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1).
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. [Sigmoid loss for language image pre-training](#). In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11941–11952.
- Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. 2023a. [A simple llm framework for long-range video question-answering](#). *Preprint*, arXiv:2312.17235.
- Chaoyi Zhang, Kevin Lin, Zhengyuan Yang, Jianfeng Wang, Linjie Li, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2024. [Mm-narrator: Narrating long-form videos with multimodal in-context learning](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hang Zhang, Xin Li, and Lidong Bing. 2023b. [Video-LLaMA: An instruction-tuned audio-visual language model for video understanding](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 543–553, Singapore. Association for Computational Linguistics.
- Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. 2023. [Learning video representations from large language models](#). In *CVPR*.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Jeff Huang, Chuyue Sun, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2023. [Efficiently programming large language models using sglang](#). *Preprint*, arXiv:2312.07104.
- Linchao Zhu, Zhongwen Xu, Yi Yang, and Alexander G. Hauptmann. 2017. [Uncovering the temporal context for video question answering](#). *International Journal of Computer Vision*, 124(3):409–421.

Supplementary Material for “TravelER”

Here we provide additional information about our experimental results, qualitative examples, implementation details, and datasets. Specifically, Section A provides more experiment results, Section B provides additional implementation details, Section C provides additional related work, and Section D provides qualitative visualizations to illustrate our approach.

A Additional Experiment Results

We begin by presenting several additional ablations in Section A.1 that further demonstrate the benefits of our TravelER approach. We also present additional results in Section A.2.

A.1 Additional Ablations

In what follows, we provide additional ablations that further illustrate the benefits of TravelER. For all ablations, we compare the ablated experiment with the corresponding best-performing TravelER results on a random sample of 1000 examples from the training set of the NExT-QA dataset. We use GPT-3.5 as the LLM and LLaVA-1.6 as the LMM. **LMM response length.** The LMM in our framework is crucial because it allows us to capture more relevant and question-specific details from visual input. However, if the LMM’s responses are too long, the memory bank will become too large, whereas if the LMM’s responses are too short, insufficient information will be captured. Thus, we conduct an experiment to determine the optimal LMM response length, and display our results in Fig 4. We find that limiting the LMM response to 150 tokens yields the most optimal performance, while accuracy decreases by -2.2% and -1.7% if the response is limited to 75 tokens and 300 tokens respectively. This supports the fact that there is a tradeoff between not collecting enough information for short response lengths and collecting too much information as the LMM response size increases.

Prompt analysis. In each module of our framework, we use a task prompt to provide instructions to our agents (LLMs or LMMs). The construction of these prompts plays a large role in how instructions are executed. Currently, we use the question Q as input into all prompts $(P_T^{(1)}), (P_R^{(1)}), (P_E^{(1)}), (P_V^{(1)})$. However, we use the choices C as input only for the Planner and Evaluator prompts since the Planner needs the choices to tailor its plan, and the Evaluator needs

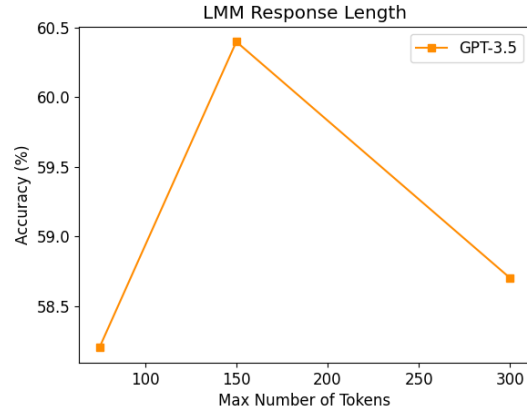


Figure 4: Comparison of different LMM Response Length (75, 150, 300 max tokens). 150 is optimal.

the choices to answer the question. We experiment with adding the choices C to the Retriever and Extractor prompts, and find that this degrades performance by -1.6%. This may be because the incorrect choices mislead the Retriever into searching for non-existent events or the Extractor into asking irrelevant questions.

A.2 Additional Results

Results on Causal-VidQA. We also report zero-shot results on Causal-VidQA (Li et al., 2022), a dataset designed to focus on causal related questions to facilitate deeper video understanding towards video reasoning. We see that our method outperforms the best zero-shot method by 16.9%. **Additional Baselines.** We also compare our approach with the concurrent work MoReVQA (Min et al., 2024). It can be seen that MoReVQA performs better on NExT-QA by 1%, while our method performs better on the more difficult long-form video understanding EgoSchema benchmark by 1.6%. We also note some key differences. Firstly, our method is able to re-plan based on feedback from previous iterations while MoReVQA uses a single forward pass through the stages and cannot modify its approach during the process. Secondly, we use an iterative approach while MoReVQA does not, which allows our method to adjust and improve using feedback from previous iterations.

B Additional Implementation Details

To run our models on larger benchmarks, we use 8 NVIDIA RTX 6000 GPUs and split the dataset across multiple processes. In addition, we use the

Table 5: Zero-shot (ZS) results on Causal-VidQA.

Model	Acc. (%)
Just-Ask (Yang et al., 2021)	27.1
CaKE + CoMem (Su et al., 2023)	30.0
CaKE + HGA (Su et al., 2023)	30.6
TravelER (ours)	47.5 (+16.9)

SGLang (Zheng et al., 2023) package, which provides a variety of performance optimizations for our LMMs and enables us to perform batched inference for models that do not natively support doing so. We serve our LMM on a single GPU and implement a queue that is shared across all runs. This allows individual runs to asynchronously call the LMM using an API request instead of creating a new instance of the LMM for each run. Typically, we have 4-5 processes sharing the same LMM.

Next, we use all default parameters. For our LMMs, we experiment with modifying the maximum output token length as an ablation, but use all default parameters otherwise. For the Llama 3 ablation in 3a, we serve Llama 3 70B on 4 NVIDIA A100s using vLLM (Kwon et al., 2023). We report results from a single run for all experiments.

B.1 Prompts

Our prompts are shown in Table 6, Table 7, Table 8, and 9. The black text is the base prompt template, and we replace the blue text with the corresponding information from the relevant video. The generated outputs are in the orange text.

B.2 NExT-QA

Dataset. NExT-QA is a challenging dataset that tests causal action reasoning and temporal understanding. It contains 5,440 videos with an average length of 44s. Compared to earlier VideoQA benchmarks (Papalampidi et al., 2024; Xue et al., 2017; Zeng et al., 2017; Zhu et al., 2017), NExT-QA requires going beyond simple recognition of objects and actions to answer the questions correctly. Each question requires selecting the best option out of 5 choices, often with very similar degrees of plausibility. Additionally, each question is categorized into either a Temporal, Causal, or Descriptive type. Temporal questions often ask what happens during, before, or after an event or action, while causal questions require advanced reasoning and inference about why an event or action occurs. Following the trend of works before us, we evaluate our method

on the 5,000 questions in the NExT-QA validation set, which consist of 500 different videos. This dataset is in English.

Inference Details. For NExT-QA, we use LLaVA-1.6 (Vicuna 13B) as the LMM and GPT-4 (gpt-4-1106-preview) as the LLM. We used the longer, comprehensive prompts, with no answer choices included in the Extractor prompt. We also initialize the memory bank with 5 frames, and use the multi-frame Retriever with 5 frames.

B.3 Perception Test

Dataset. In comparison with earlier VideoQA datasets (Maharaj et al., 2017; Zhu et al., 2017; Zeng et al., 2017) that focus on computational tasks such as classification, detection, or tracking, Perception Test is a dataset that focuses on skills such as memory, abstraction, physics, and semantics. Moreover, it is designed to test the transfer capabilities of different models and intended to be approached in a few-shot or zero-shot manner. The dataset consists of 11.6k real-world videos with an average length of 23 seconds, and 38K multiple choice QA questions. This dataset is in English.

Inference Details. For Perception Test, we use LLaVA-1.6 (Vicuna 13B) as the LMM and GPT-3.5 (gpt-3.5-turbo-0125) as the LLM. We used shorter, simplified prompts, with no answer choices included in the Extractor prompt. We also initialize the memory bank with 5 frames, and use the multi-frame Retriever with 5 frames.

B.4 EgoSchema

Dataset. EgoSchema is a dataset that tests long-form video reasoning and is intended to be answered in a few-shot or zero-shot manner. It introduces the idea of certificate lengths, which are the minimum number of seconds it takes to be able to answer the question correctly. It has 5k questions in the full set, and a 500 question subset is used as the validation set. We use the full test set for evaluation. This dataset is in English.

Inference Details. For EgoSchema, we use GPT-4 (gpt-4-1106-preview) as the LLM and LaViLa as the LMM. To prevent data leakage, we use a retrained version of LaViLa that does not use overlapping videos between Ego4D and EgoSchema. For window size in the Retriever, we take the selected frame, along with the frame 1 second before and 1-2 seconds after the selected frame.

<p>User: Create the best plan to gather information to answer the question. QUESTION: QUESTION CHOICES: CHOICES You are provided information collected from individual frames of the video and is represented as a dictionary keyed by the timestamps of the frames. INFORMATION: INFO You are also given an explanation for why you aren't able to definitively answer the question with the current information. EXPLANATION: EXPLANATION Follow these rules: 1. You only have access to individual frames of the video, with no audio. You can go to a certain timestamp, search for actions or settings, and describe or ask questions about individual frames. 2. Make sure that you have viewed the relevant frames. Make your plan as simple and straightforward as possible, and no longer than 5 steps long. Return your plan as a numbered list, after PLAN. Do not include any other response or explanation. Let's think step-by-step.</p> <hr/> <p>Assistant: Output: PLAN</p>

Table 6: Planner prompt P_T :

<p>User: You are given the following information about a LENGTH second video, with information from individual frames at different timestamps. INFORMATION: INFO PLAN: PLAN Currently, you are viewing second CURR. Choose the timestamp, in seconds, of the next frame to view. When choosing the next frame to view, remember that you are trying to collect information to answer this multiple choice question: QUESTION Think of what information you need, and consider what information you already have. Use the temporal nature of the video and your past information to choose the next frame. Do not choose a frame you already have information about, and make sure that the frame you choose is at least WINDOW SIZE seconds apart from the second you are currently viewing. Return your answer as a single Python float representing the second you want to view. Don't provide any other response or explanation.</p> <hr/> <p>Assistant: Output: TIMESTAMP</p>

Table 7: Retriever prompt P_R :

User:

You are given the following information about a **LENGTH** second video, with information from individual frames at different timestamps.

INFORMATION: **INFO**

Currently, you are viewing second **CURRENT TIMESTAMP**, which has the caption: **FRAME CAPTION**

Form up to 3 questions about this frame to best help answer the multiple-choice question: **QUESTION**.

Follow these rules:

1. Use the given information to decide what further visual information you need to answer the question.
2. Since you are asking questions about a single frame, you cannot ask about other frames, reference past or future events, or ask about specific timestamps.

Return your questions as a Python list of strings (in double quotes) and don't include any numbered lists, backticks, or language hints. Follow Python syntax. Make sure you have followed the steps. Don't provide any other response or explanation.

Assistant: Output: **QUESTIONS**

Table 8: Extractor prompt P_E :

User:

Evaluate if there is enough information to answer a multiple-choice question about a video and if the plan has been completed.

If there is enough information to choose the correct answer with complete certainty and the plan has been followed, return the index of the choice after a brief explanation. Otherwise, return None after a brief explanation of why you can't narrow down to a single answer choice. Be strict and don't guess.

INFORMATION: **INFO**

PLAN: **PLAN**

QUESTION: **QUESTION**

CHOICES: **CHOICES**

Give a brief explanation. Then, include your final answer after the words "Final Answer:" in your response at the end. Do not include anything other than the answer as an integer or None after "Final Answer:".

Let's think step by step.

Assistant: Output: **ANSWER**

Table 9: Evaluator prompt P_V :

B.5 STAR

Dataset. STAR is a dataset that tests reasoning in real-world video situations. It consists of 22K video clips, with 60K situated reasoning questions, with 4 possible choices each. Questions are broadly divided into 4 main categories: interaction, sequence, prediction, and feasibility. This dataset is in English.

Inference Details. For STAR, we use LLaVA-1.6 (Vicuna 13B) as the LMM and GPT-3.5 (gpt-3.5-turbo-0125) as the LLM. We used shorter, simplified prompts, with no answer choices included in the Extractor prompt. We also initialize the memory bank with 5 frames, and use the multi-frame Retriever with 3 frames, since the videos are shorter.

C Additional Related Work

Modular Vision Frameworks. There has been a long history of work (Andreas et al., 2016; Herzig et al., 2018; Krishna et al., 2018; Baradel et al., 2018; Battaglia et al., 2018; Herzig et al., 2020, 2022, 2019, 2023; Bar et al., 2021; Avraham et al., 2022; Jerbi et al., 2020; Herzig et al., 2024; Mitra et al., 2024) that attempts to combine deep neural networks with modularity. Recently, works like VisProg (Gupta and Kembhavi, 2023), CodeVQA (Subramanian et al., 2023), RVP (Ge et al., 2023), and ViperGPT (Surís et al., 2023) have leveraged the improved coding capabilities of LMMs to generate code to compose different submodules together to answer visual questions. In addition, ProViQ (Choudhury et al., 2023) extends ViperGPT’s work in the video domain by adding more modules for VideoQA. Similarly, we leverage the strong power of LMMs in a modular approach. However, while these approaches have shown promising results, they are limited to single-shot planning when generating code, resulting in a fixed plan that cannot adapt. In contrast to these works, our approach has the advantage of being able to iteratively replan based on new information collected.

D Qualitative Visualizations

We present further qualitative success and failure cases of our TravelER framework. For each dataset, we display qualitative visualizations for 2 successes and 2 failure cases. For the success cases, we show expanded visualizations in Figures Table 5, Table 6, Table 7, and abridged versions in 8 that demonstrate the benefits of our question-answering approach compared to regular caption-

ing. For failures, we also present abridged versions for each dataset in Figures Table 9, 10, and 11. Finally, we present some additional success and failure cases in 12. For the visualizations, we display 3 iterations of the question answering process with 2 Extractor QA pairs each, compressing the output text by displaying the most important parts for brevity. For the abridged versions, we display the video on top, and the traversal order using the numbered orange circles. In the row beneath, we display the frames in the order they are selected, and display corresponding Extractor question-answer output in yellow and captions in gray.

E Licenses and Privacy

The license, PII, and consent details of each dataset are in the respective papers. In addition, we wish to emphasize that the datasets we use do not contain any harmful or offensive content, as many other papers in the field also use them. Thus, we do not anticipate a specific negative impact, but, as with any Machine Learning method, we recommend to exercise caution.

Successes

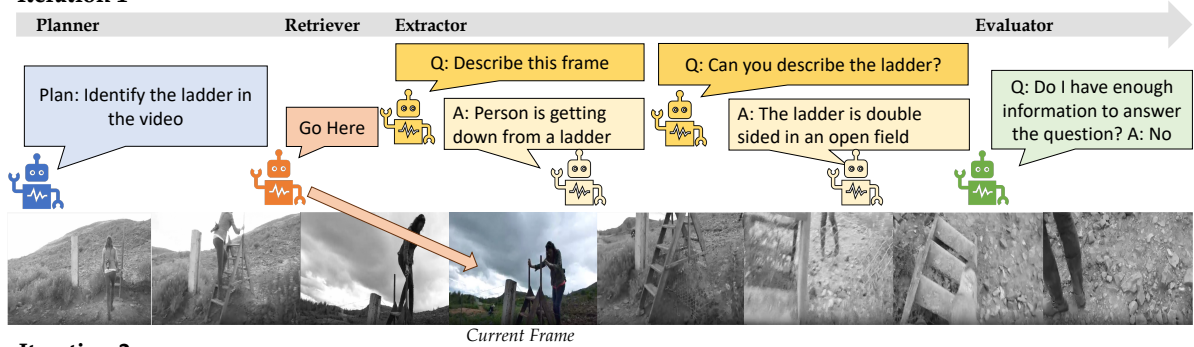
NeXT-QA

Question and Choices

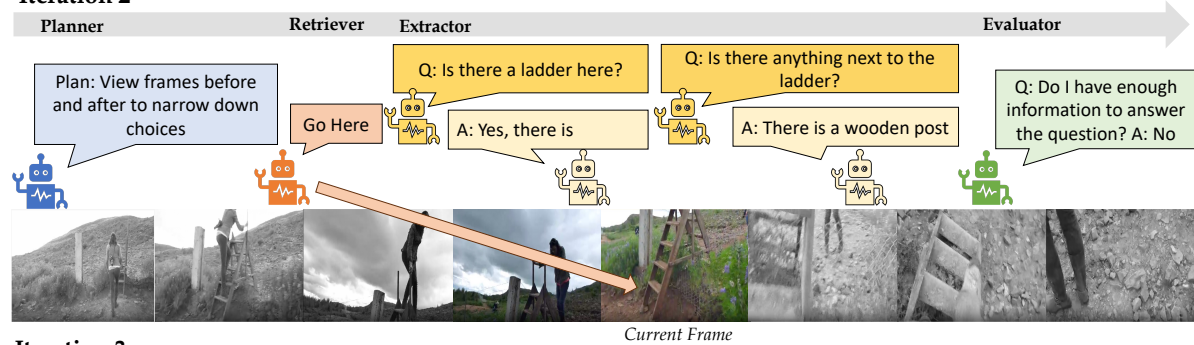
Q: Why is there a ladder in the middle of the land?

A: (A) for child to play in (B) supervising boy (C) tree trimming (D) **cross over fence** (E) to play

Iteration 1



Iteration 2



Iteration 3

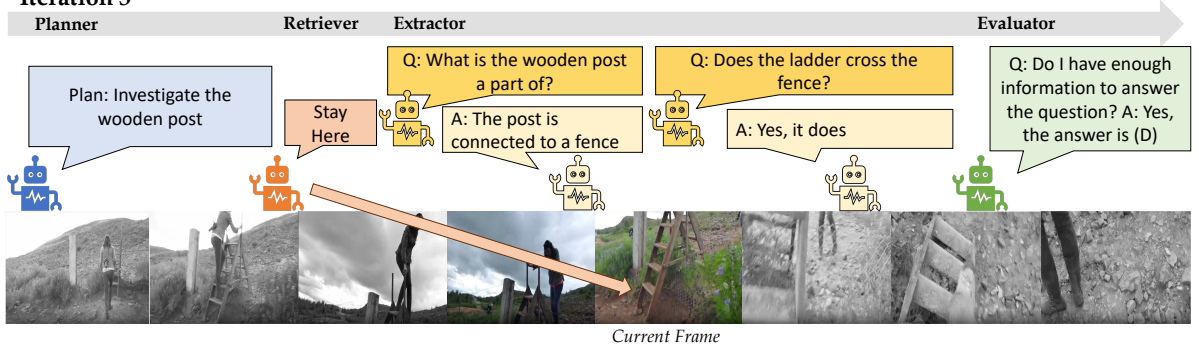


Figure 5: **NeXT-QA Success Predictions.** We can see that our framework can adapt to new information collected in past iterations. For example, in Iteration 3, our Planner module is able to use information about the wooden post from a previous iteration and ask further questions to identify the correct answer.

Successes

STAR

Question and Choices

Q: What happened before the person took the blanket?

A: (A) Opened the box. (B) Opened the door. (C) **Ate the sandwich.** (D) Threw the towel.



Figure 6: **STAR Success Predictions.** Here, we can see that our method does not require viewing frames sequentially. For example, we view the beginning of the video in Iteration 1, the middle of the video in Iterations 2 and 3, and return to the beginning in Iteration 4. Moreover, our method can collect information and double-check ambiguous information across different timestamps. For example, in Iteration 2, we are told the man is taking a blanket, and then we can view a different frame to confirm that he is indeed holding a towel in Iteration 3.

Successes

Perception Test

Question and Choices

Q: What stops the motion of the object placed on the slanted plane after being released by the person the first time?
 A: (A) I don't know. (B) **The person or collision with another object.** (C) High friction with the surface.

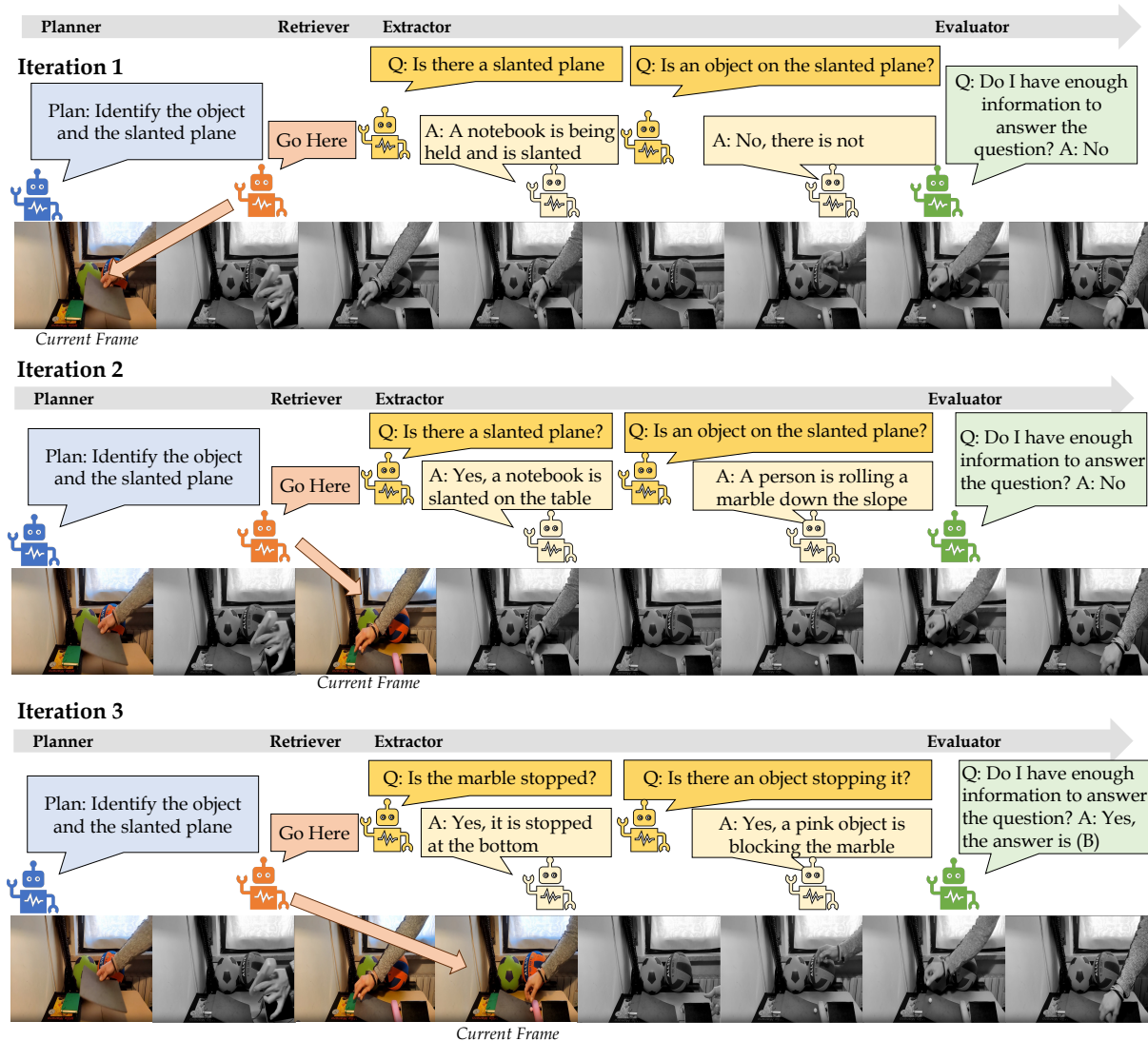


Figure 7: **Perception Test Success Predictions.** We display some success cases for the challenging Perception Test dataset. Here, our method is able to infer which objects the question refers to through our question-asking approach, even though the question does not explicitly describe them.

Successes

NeXT-QA

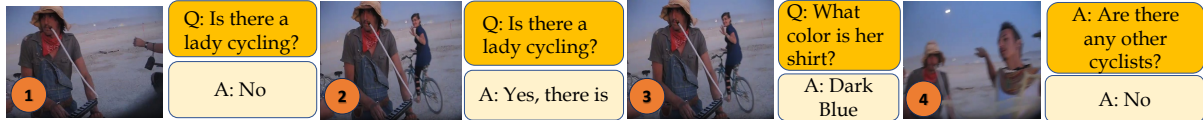
Question and Choices

Q: What is the color of the shirt of the lady cycling?

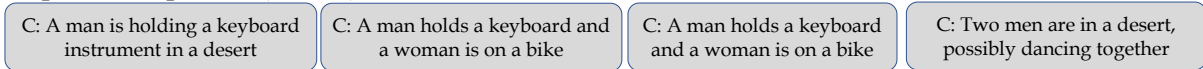
A: (A) black (B) light blue (C) white and red (D) **blue** (E) purple



TravelER Output (Ours)



Caption Comparison (Others)

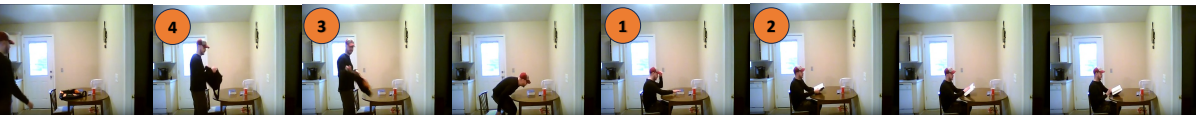


STAR

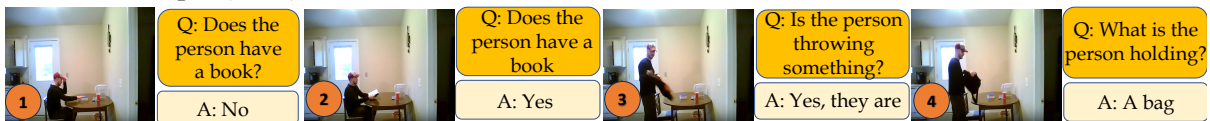
Question and Choices

Q: What object did the person throw before they took the book?

A: (A) the towel (B) the laptop (C) **the bag** (D) the blanket



TravelER Output (Ours)



Caption Comparison (Others)

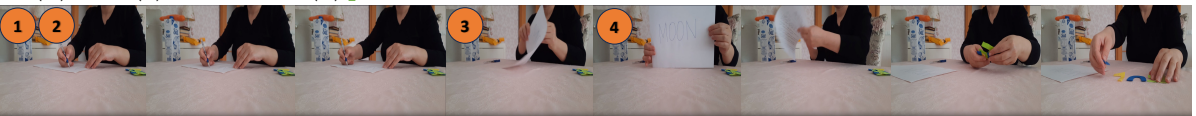


Perception Test

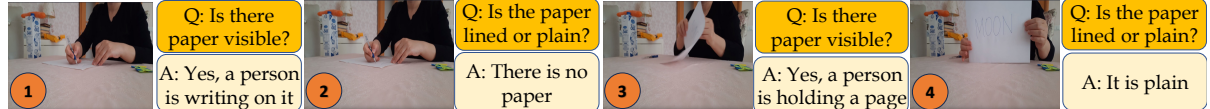
Question and Choices

Q: Is the paper lined or plain?

A: (A) lined (B) I don't know (C) **plain**



TravelER Output (Ours)



Caption Comparison (Others)

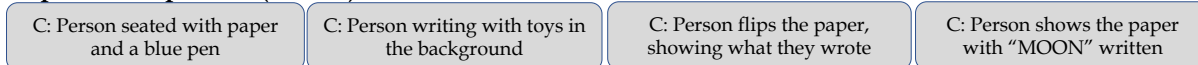


Figure 8: **Comparison with Captioning Approaches.** For each example, we display the videos on top, and the traversal order using the numbered orange circles. In the rows beneath, we display the frames in the order they are selected, and display corresponding Extractor question-answer output in yellow and captions in gray. We display compressed versions of GPT-4V-generated captions for a visual comparison. By asking specific questions, we can extract more detailed and relevant information than a general description.

Failures

NeXT-QA

Example 1

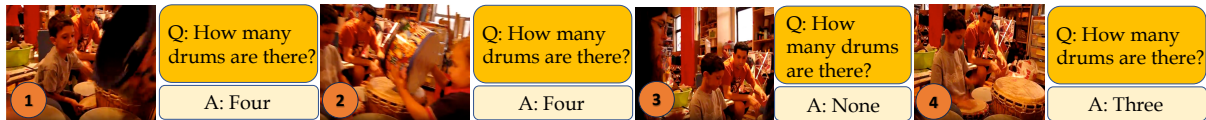
Question and Choices

Q: How many drums in total are there for the people to play?

A: (A) **four** (B) two (C) **six** (D) three (E) eleven



TravelER Output (Ours)



Caption Comparison (Others)

C: A child holds a drum set and looks forward

C: A child is focused, while another person plays the drums in the foreground

C: A woman appears to be instructing the people playing the drums

C: A child is playing a drum with concentration, with another person behind them

Example 2

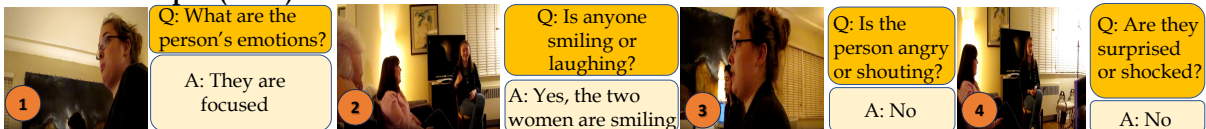
Question and Choices

Q: How is the atmosphere of this group discussion?

A: (A) **happy** (B) **serious** (C) sad (D) angry (E) shocking



TravelER Output (Ours)



Caption Comparison (Others)

C: A person looking towards something off-camera

C: Two people smiling and engaged in conversation

C: A person looking towards something off-camera

C: People gathered in a cozy living room

Figure 9: **NeXT-QA Failure Predictions.** Here, we display some failure cases for NeXT-QA. Like before, we display the video on top, and the traversal order using the numbered orange circles. In the row beneath, we display the frames in the order they are selected, and display corresponding Extractor question-answer output in yellow and captions in gray. We can see that conflicting information or false positives can mislead our approach. We also observe that counting can be a challenge for certain LMMs, but this can be mitigated in the future by swapping in stronger LMMs.

Failures

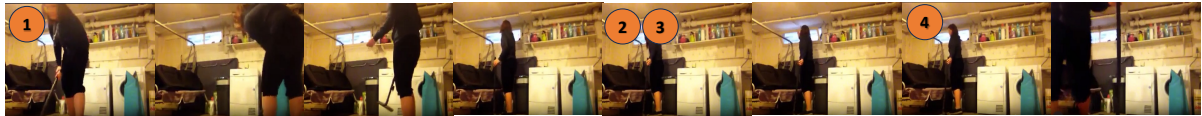
STAR

Example 1

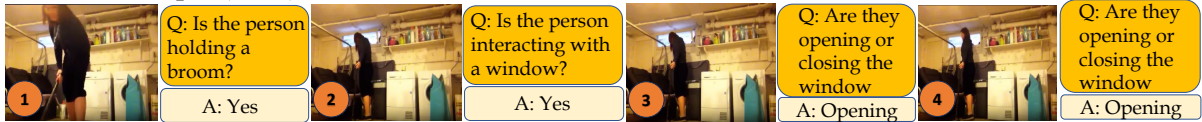
Answer

Q: What did the person do to the window after holding the broom?

A: (A) **Closed**. (B) **Opened**. (C) Washed. (D) Lied on.

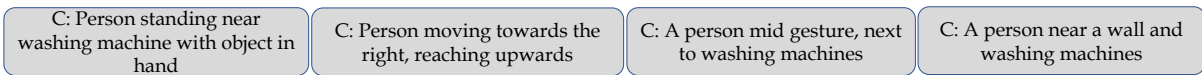


TravelER Output (Ours)



STAR

Caption Comparison (Others)

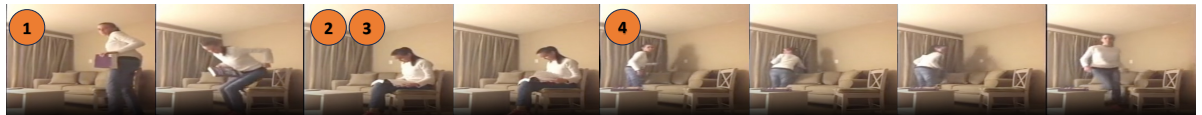


Example 2

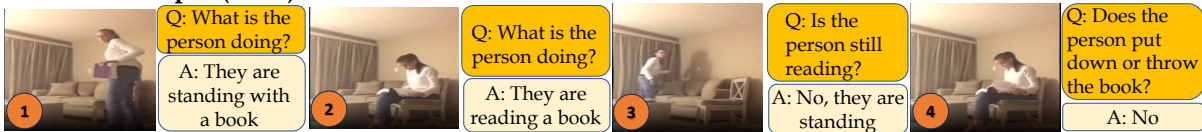
Answer

Q: What did the person do with the book?

A: (A) **Took**. (B) **Opened**. (C) Threw. (D) Put down.



TravelER Output (Ours)



Caption Comparison (Others)

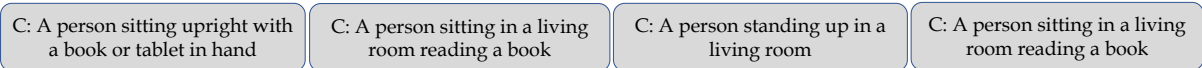


Figure 10: **STAR Failure Cases**. Here, we display some failure cases for the STAR dataset, using the same abridged representation described previously. We see that a limitation of a frame-wise approach is that it may be difficult to capture very temporal actions. For example, in Example 1, it is difficult to understand if the woman is opening or closing the window.

Failures

Perception Test

Example 1

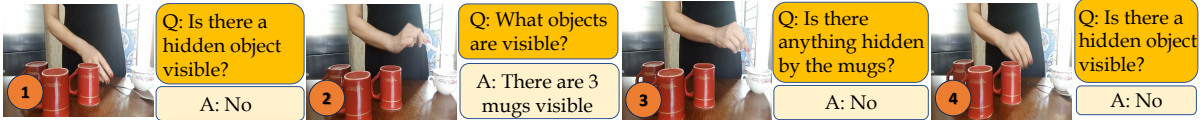
Question and Choices

Q: The person uses multiple similar objects to play an occlusion game. Where is the hidden object at the end of the game from the person's point of view?

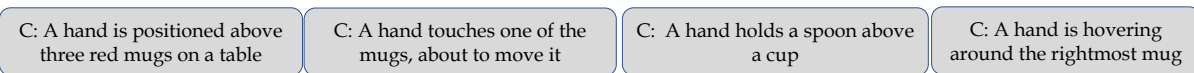
A: (A) **Under the third object from the left.** (B) **Under the second object from the left.** (C) Under the first object from the left.



TravelER Output (Ours)



Caption Comparison (Others)



Example 2

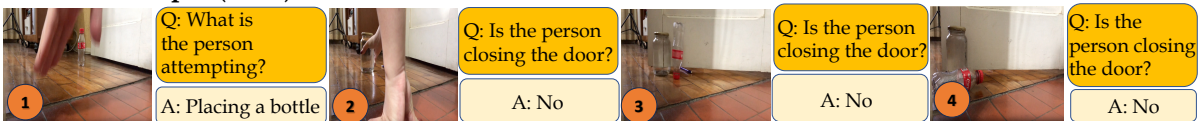
Question and Choices

Q: Can the person close the door during the last attempt?

A: (A) No, the person cannot close the door. (B) **Yes, the person can close the door.** (C) **The person didn't attempt to close the door in this video.**



TravelER Output (Ours)



Caption Comparison (Others)

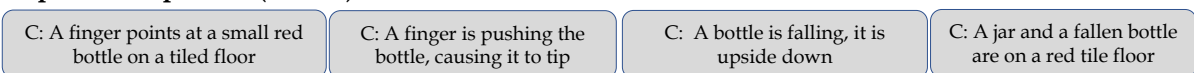






Figure 11: **Perception Test Failure Cases.** We display some qualitative visualizations for Perception Test failure cases using the abridged representation discussed previously. We see that for some cases where objects of interest are occluded or not in the frame, our method might have difficulties extracting useful information.

Successes

NEXT-QA

Q: Why are the men wearing hats?

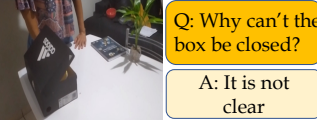
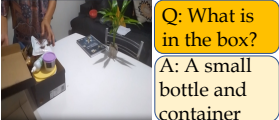
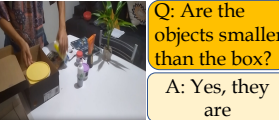
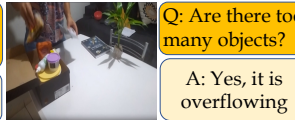
A: (A) Sunny (B) Fashion (C) Costume of play (D) Uniform costume (E) Road repairing

	Q: What is the weather like? A: It is sunny		Q: Are the men wearing uniforms? A: No, they wear outdoor gear		Q: What are the men doing? A: They are walking		Q: Are there shadows that show sunny weather? A: Yes, there are
C: A group of hikers with backpacks outdoors		C: A group of hikers, with one person resting an arm on another		C: Friends hiking together, with one's arm on another		C: A group of hikers, walking together along a road	

Perception Test

Q: Why can't the box be closed at the end?

A: (A) The person put one or more objects that are too large to fit in the box, even if placed by themselves. (B) Each of the objects would fit in the box individually, but the person put too many objects. (C) The lid does not match the box.





	Q: Why can't the box be closed? A: It is not clear		Q: What is in the box? A: A small bottle and container		Q: Are the objects smaller than the box? A: Yes, they are		Q: Are there too many objects? A: Yes, it is overflowing
C: A person standing behind a table, presenting a shoebox with a GoPro logo on it.		C: A person is unpacking items on a table. There are objects like a plant vase and boxes		C: A person organizes items on a table with objects like a box with a yellow lid, and a potted plant.		C: A person appears to be sorting through items on a table, with a stack of boxes and containers	

Failures

NEXT-QA

Q: What does the lady do after taking the microphone from the man in red cape at the start?

A: (A) put 2 fingers up (B) talk to the audience (C) turn behind her (D) eat (E) put hands together

	Q: Does the lady have the mic? A: No, the man has it		Q: Does the lady have the mic? A: Yes		Q: What is the lady doing? A: She is holding it		Q: What is she doing? A: She is high-fiving
C: A karaoke night with participants on stage, one draped in a Swiss flag.		C: An eclectic performance at a colorful music event.		C: Band performing on stage with vocalist as other members look on.		C: A high-spirited karaoke night where a performer gets a high-five on stage.	

Perception Test

Q: The person interacts with a lighting device among other objects. Is the lighting device on at any point?

A: (A) I don't know (B) yes (C) no

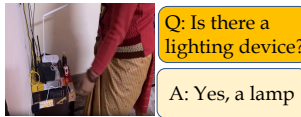
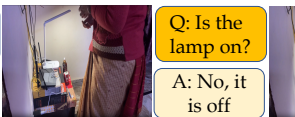
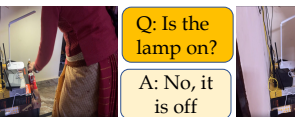
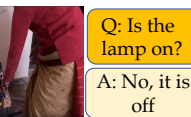
	Q: Is there a lighting device? A: Yes, a lamp		Q: Is the lamp on? A: No, it is off		Q: Is the lamp on? A: No, it is off		Q: Is the lamp on? A: No, it is off
C: A person is reaching out towards an electronic device with wires and routers.		C: A person is standing next to a shelf with household items.		C: An individual is standing by a shelf with electronic devices and a router.		C: A person is standing next to a shelf with household items.	

Figure 12: An additional visualization of predictions. We show more qualitative visualizations of our method on NEXT-QA and Perception Test using our abridged representation, with successes on top and failures on the bottom. We compare our generated question-answer pairs for each frame (in yellow) with captions (labeled C in gray) generated from the same frame. We see that our method is able to extract more fine-grained and relevant information compared to simple captioning.