

A Simple LLM Framework for Long-Range Video Question-Answering

Ce Zhang* Taixi Lu* Md Mohaiminul Islam Ziyang Wang Shoubin Yu
 Mohit Bansal Gedas Bertasius

Department of Computer Science, UNC Chapel Hill

{cezhang, mmiemon, ziyangw, shoubin, mbansal, gedas}@cs.unc.edu, taixi@email.unc.edu

Abstract

We present LLoVi, a language-based framework for long-range video question-answering (LVQA). Unlike prior long-range video understanding methods, which are often costly and require specialized long-range video modeling design (e.g., memory queues, state-space layers, etc.), our approach uses a frame/clip-level visual captioner (e.g., BLIP2, LaViLa, LLaVA) coupled with a Large Language Model (GPT-3.5, GPT-4) leading to a simple yet surprisingly effective LVQA framework. Specifically, we decompose short and long-range modeling aspects of LVQA into two stages. First, we use a short-term visual captioner to generate textual descriptions of short video clips (0.5-8s in length) densely sampled from a long input video. Afterward, an LLM aggregates the densely extracted short-term captions to perform long-range temporal reasoning needed to understand the whole video and answer a question. To analyze what makes our simple framework so effective, we thoroughly evaluate various components of our system. Our empirical analysis reveals that the choice of the visual captioner and LLM is critical for good LVQA performance. Furthermore, we show that a specialized prompt that asks the LLM first to summarize the noisy short-term visual captions and then answer a given input question leads to a significant LVQA performance boost. On EgoSchema, which is best known as a very long-form video question-answering benchmark, our method achieves 50.3% accuracy, outperforming the previous best-performing approach by 18.1% (absolute gain). In addition, our approach outperforms the previous state-of-the-art by 4.1% and 3.1% on NeXT-QA and IntentQA. We also extend LLoVi to grounded LVQA and show that it outperforms all prior methods on the NeXT-GQA dataset. We will release our code at <https://github.com/CeeZh/LLoVi>.

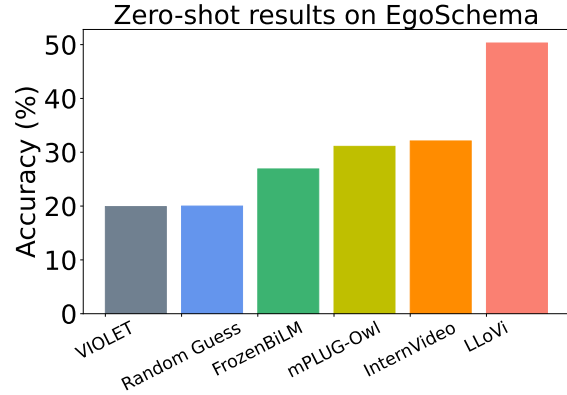


Figure 1. Comparison of long-range video question-answering (LVQA) accuracy on EgoSchema, the new LVQA benchmark that addresses many limitations of prior long-range video benchmarks. Our simple LLM framework, named LLoVi, outperforms all prior LVQA approaches by a significant margin (+18.1%).

1. Introduction

Recent years have witnessed remarkable progress in short video understanding (5-15s in length) [11, 61, 63, 73, 77]. However, extending these models to long videos (e.g., several minutes or hours in length) is not trivial due to the need for sophisticated long-range temporal reasoning capabilities. Most existing long-range video models rely on costly and complex long-range temporal modeling schemes, which include memory queues [4, 23, 24, 67], long-range feature banks [5, 66, 82], space-time graphs [18, 60], state-space layers [19, 20, 56] and other complex long-range modeling modules [1, 17, 76].

Recently, Large Language Models (LLMs) have shown impressive capability for long-range reasoning on a wide range of tasks such as document understanding [14, 49, 58] and long-horizon planning [15, 35, 47]. Motivated by these results in the natural language and decision-making domain, we explore using LLMs for long-range video question answering (LVQA). Specifically, we propose LLoVi, a simple language-based framework for long-range video understanding. Unlike prior long-range video models, our ap-

*The first two authors contribute equally.

proach does not require specialized long-range video modules (e.g., memory queues, state-space layers, spatiotemporal graphs, etc.) but instead uses a short-term visual captioner coupled with an LLM, thus exploiting the long-range temporal reasoning ability of LLMs. Our simple two-stage approach tackles the LVQA task by decomposing it into short and long-range modeling subproblems:

1. First, given a long video input, we segment it into multiple short clips and convert them into one-sentence textual descriptions using a pre-trained frame/clip-level visual captioner (e.g., BLIP2, LaViLa).
2. Afterward, we concatenate the temporally ordered captions from Step 1 and feed them into an LLM (e.g., GPT-3.5, GPT-4) with proper prompts to perform long-range video reasoning for LVQA.

Recent works have explored incorporating visual captioning models into LLMs for video understanding. Several methods proposed using LLMs for various short-term video understanding tasks. For example, Video ChatCaptioner [3] generates enriched video descriptions by learning to answer questions from ChatGPT. Similarly, ChatVideo [55] stores the detected tracklets in a database, which are then used for interacting with the user. Additionally, Socratic Models [81], VidIL [62], and VideoChat [31] all use pretrained visual models to extract low-level video concepts (e.g., objects, actions, scenes, etc.) and then leverage LLM to perform various short-term video understanding tasks. Beyond short-term video understanding tasks, we note that the concurrent work in [10] applies an LLM-based framework for long-range video understanding. However, their analysis is limited to movie-based datasets that rely heavily on non-visual inputs such as speech and subtitles, thus requiring limited visual analysis [37]. Compared to this concurrent work, we conduct our experiments on EgoSchema, the newly introduced LVQA benchmark that addresses many limitations of prior long-range video benchmarks. Furthermore, unlike prior LLM-based video understanding frameworks, many of which perform only qualitative analysis of their models [3, 31, 33, 55], we conduct a thorough empirical study investigating the effectiveness of our framework.

Specifically, we investigate (i) the selection of a visual captioner, (ii) the choice of an LLM, (iii) the LLM prompt design, (iv) few-shot in-context learning, (v) optimal video processing configurations (i.e., clip length, sampling rate, etc.), and (vi) the generalization of our framework to other datasets and tasks. Our key empirical findings include:

- A multi-round prompt that first asks the LLM to summarize the noisy/redundant short-term visual captions and then answer a given question based on the summarized captions leads to the most significant boost in performance (+5.8%) among the prompts we have tried (e.g., zero-shot CoT, Self-Consistency).
- GPT-3.5 provides the best tradeoff between accuracy,

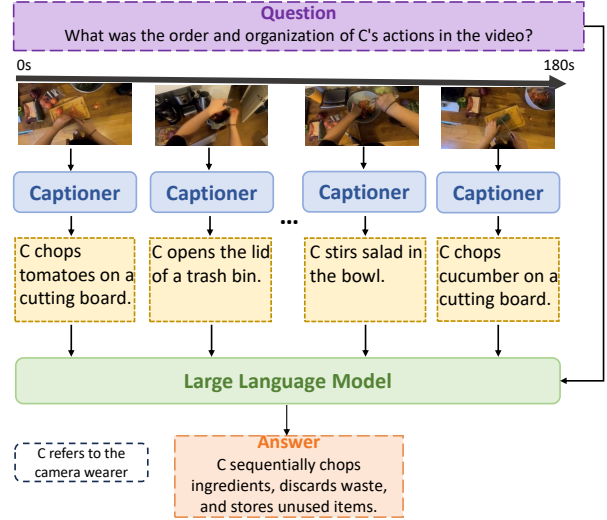


Figure 2. An illustration of LLoVi, our simple LLM framework for long-range video question-answering (LVQA). We use Large Language Models (LLMs) like GPT-3.5 and GPT-4 for their long-range modeling capabilities. Our method involves two stages: first, we use short-term visual captioners (e.g., LaViLa, BLIP2) to generate textual descriptions for brief video clips (0.5s-8s). Then, an LLM aggregates these dense, short-term captions for long-range reasoning required for LVQA. This simple approach yields impressive results, demonstrating LLMs’ effectiveness in LVQA.

- computational cost, and sufficiently large context length.
- Using LaViLa visual captioner [83] leads to best results (**51.8%**) followed by BLIP-2 [28] (**46.7%**) and EgoVLP [43] (**46.6%**).
- Few-shot in-context learning leads to a large improvement on both the variant of our model with the simplest prompt (+4.7%) and our best-performing variant with a multi-round prompt (+4.1%).
- Densely extracting visual captions from consecutive 1-second video clips of the long video input leads to the strongest performance on EgoSchema.
- Our final method, built using the above-listed empirical insights, achieves **50.3%** zero-shot LVQA accuracy on the full test set of EgoSchema, outperforming the previous best approach by **18.1%** as shown in Figure 1.
- Our LLoVi also outperforms prior approaches on NeXT-QA and IntentQA by **4.1%** and **3.1%**, and it also achieves state-of-the-art performance on NeXT-GQA, a recently introduced grounded LVQA benchmark.

We hope that our simple, training-free method will encourage new ideas and a simpler model design in LVQA. We will release our code to enable the community to build on our work.

2. Related Work

Long-range Video Understanding. Modeling long-range videos (e.g., several minutes or longer) typically requires

models with sophisticated temporal modeling capabilities, often leading to complex model design. LF-VILA [50] proposes a Temporal Window Attention (HTWA) mechanism to capture long-range dependency in long-form video. MeMViT [67] and MovieChat [48] adopt a memory-based design to store information from previously processed video segments. Several prior methods use space-time graphs [18, 60] or relational space-time modules [76] to capture spatiotemporal dependencies in long videos. Lastly, the recently introduced S4ND [40], ViS4mer [19] and S5 [56] use Structured State-Space Sequence (S4) [13] layers to capture long-range dependencies in the video. Unlike these prior approaches, we *do not* propose any complex long-range temporal modeling modules but instead develop a simple and strong LLM-based framework for zero-shot LVQA.

LLMs for Video Understanding. The recent surge in large language models (LLMs) [2, 9, 41, 46, 53, 54] has inspired many LLM-based applications in video understanding. Models like Socratic Models [81] and VideoChat [31] integrate pretrained visual models with LLMs for extracting visual concepts and applying them to video tasks. Video ChatCaptioner [3] and ChatVideo [55] leverage LLMs for video representation and dialog-based user interaction, respectively. VidIL [62] employs LLMs for adapting image-level models to video tasks using few-shot learning. Beyond short-term understanding, studies [10, 33] have explored LLMs for long-range video modeling. The work in [33] uses GPT-4 for video summarization but lacks quantitative evaluation. Meanwhile, [10] focuses on movie datasets, requiring limited visual analysis [37]. In contrast, we conduct our experiments on the EgoSchema benchmark for long-range LVQA and provide an extensive empirical analysis of various design choices of our simple approach.

Video Question Answering. Unlike image question-answering, video question-answering (VidQA) presents unique challenges, requiring both spatial and temporal reasoning. Most existing VidQA methods, either using pretraining-finetuning paradigms [6, 26, 79], zero-shot [34, 51, 74, 79], or few-shot learning [62], focus on short-term video analysis (5-30s). To overcome the limitations of short-term VidQA, new benchmarks have been proposed covering longer video durations: NextQA [68] averages 44s, while ActivityNet-QA [80], TVQA [25], How2QA [72], MovieQA [52], and DramaQA [7] range from 100s to several minutes. Despite these lengths, works in [21, 37, 75] found that many benchmarks can be solved by analyzing only brief clips, not requiring extensive video modeling, and exhibited language biases, being solvable using pure text-only methods that ignore visual content. To address these issues, the EgoSchema benchmark [37] was recently introduced, requiring at least 100 seconds of video analysis and not exhibiting any language biases. Thus, our work focuses on the LVQA analysis on EgoSchema.

3. Method

Recently, LLMs have been shown to excel on a wide range of long-range modeling tasks [14, 15, 35, 47, 49, 58]. Motivated by these long-range modeling capabilities, we propose to tackle the long-range video question-answering (LVQA) task by decomposing it into two subtasks: 1) short-term video clip captioning, and 2) long-range text-based video understanding. Our motivation for such decomposition is to separate the short-term and long-range video modeling so that we could leverage the strong existing short-term visual captioners (e.g., LaViLa, BLIP2) for the first sub-task, and powerful zero-shot LLMs (e.g., GPT-3.5, GPT-4, LLaMA, T5) for the second sub-task, thus, exploiting the power of LLMs for long-range modeling.

Our decomposed LVQA framework, named LLoVi (Language-based Long-range Video Question-answering), brings important advantages. First, our approach is simple as it does not rely on complex/specialized long-range video modeling operators (e.g., memory queues, state-space layers, space-time graphs, etc.). Second, our framework is training-free, which makes it easy to apply it to LVQA in zero-shot settings. Third, our method is agnostic to the exact choices of a visual captioner and LLM, and thus, it can benefit from future improvements in visual captioning/LLM model design. Figure 2 presents a detailed illustration of our high-level approach. Below, we provide details about each component of our framework.

3.1. Short-term Video Clip Captioning

Given a long untrimmed video input V , we first segment it into N_v non-overlapping short video clips $v = \{v_m\}_{m=1}^{N_v}$, where $v_m \in \mathbb{R}^{T_v \times H \times W \times 3}$ and T_v, H, W are the number of frames, height and width of a short video clip respectively. Afterward, we feed each video clip v_m into a pretrained short-term visual captioner ϕ , which produces textual captions $c_m = \phi(v_m)$, where $c_m = (w_1, \dots, w_{L_m})$ and w_i represents the i -th word in caption c_m of length L_m . Note that our model is not restricted to any specific visual captioning model. Our experimental section demonstrates that we can incorporate various video (LaViLa [83], EgoVLP [43], VideoBLIP [78]) and image (BLIP-2 [27]) captioning models. Next, we describe how our extracted short-term captions are processed by an LLM.

3.2. Long-range Reasoning with an LLM

We want to leverage foundational LLMs for holistic long-range video understanding rather than relying on complex long-range temporal modeling schemes (e.g., memory queues, state-space layers, space-time graphs, etc.) used by prior methods [61, 67].

Formally, given short-term visual captions $\{c_m\}_{m=1}^{N_v}$ for all N_v short video clips, we first concatenate the clip captions into the full video captions $C = [c_1, \dots, c_{N_v}]$ in the

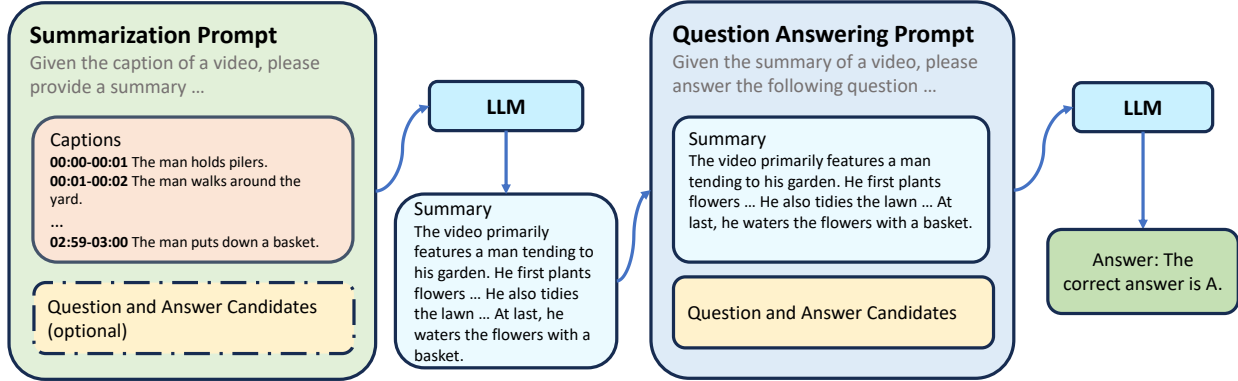


Figure 3. An illustration of the multi-round summarization-based prompt. Our empirical analysis shows that many modern LLMs may struggle when provided with long ($>1K$ words), noisy, and redundant/irrelevant caption sequences. To address this issue, we study a more specialized LLM prompt that asks an LLM first to summarize the noisy short-term visual captions (first round of prompting) and then answer a given question about the video (second round of prompting). Our results indicate that such a multi-round prompting strategy significantly boosts LVQA performance (+5.8%) compared to standard prompting techniques.

same order as the captions appear in the original video. Afterward, the concatenated video captions C are fed into an LLM for long-range video reasoning. Specifically, given the concatenated video captions C , the question Q , and the answer candidates A , we prompt the LLM to select the correct answer using the following prompt template: “Please provide a single-letter answer (A, B, C, D, E) to the following multiple-choice question $\{Q\}$. You are given language descriptions of a video. Here are the descriptions: $\{C\}$. Here are the choices $\{A\}$.”. The full prompt is included in Supplementary Materials.

Our experiments in Section 4.3 suggest that this simple approach works surprisingly well for LVQA. However, we also discovered that many modern LLMs (e.g., GPT-3.5, LLaMA) may struggle when provided with long ($>1K$ words), noisy, and potentially redundant/irrelevant caption sequences. To address these issues, we investigate more specialized LLM prompts that ask an LLM first to summarize the noisy short-term visual captions (first round of prompting) and then answer a given question about the video (second round of prompting). Specifically, we formulate such a multi-round prompt as follows: given the video captions C , the question Q , and the answer candidates A , instead of directly feeding the $\{C, Q, A\}$ triplet into LLM for LVQA, we first ask the LLM to provide a summary of the captions in the first round, which we denote as S using the following prompt template: “You are given language descriptions of a video: $\{C\}$. Please give me a $\{N_w\}$ word summary.” N_w denotes the desired number of words in the summary S . Afterward, during the second round of prompting, instead of using the captions C , we use the summary S as input for the LLM to select one of the answer candidates. Conceptually, this may be beneficial, as the LLM-generated summary S filters out potentially irrelevant/noisy informa-

tion from the initial set of captions C , making LLM inputs for the subsequent QA process more succinct and cleaner. A detailed illustration of our multi-round prompt is shown in Figure 3.

3.3. Implementation Details

For the experiments on EgoSchema, we use LaViLa [83] as our short-term captioner. We note that LaViLa is pretrained on Ego4D, the same dataset that the authors of EgoSchema use to build their benchmark. Thus, to avoid any data intersection between training and testing sets, we retrain our variant of the LaViLa model on a subset of 6K Ego4D videos that do not include any EgoSchema videos. To process video inputs, we segment each video into multiple 1s clips with a stride of 1s, resulting in a list of consecutive clips that cover the entire video. We use GPT-3.5 as the LLM for long-range reasoning on EgoSchema. For NeXT-QA, IntentQA and NeXT-GQA, we use LLaVA-1.5 [36] as the visual captioner and GPT-4 as the LLM. We down-sample the videos to 0.5 FPS and prompt LLaVA to generate captions that contain roughly 30 words for each frame. We use LaViLa as our visual captioner on EgoSchema because compared to other captioners (e.g. BLIP-2, LLaVA), LaViLa works better on first-person view videos. For the third-person view video datasets, we use LLaVA because we found its performance to be better than LaViLa. We also observe that GPT-4 generally performs better than other language models. However, on EgoSchema, we found the output of GPT-4 inconclusive for many questions (e.g. an output message that the provided information is insufficient to answer the given question). In contrast, we did not experience such issues with GPT-3.5. We also did not observe these issues with GPT-4 on other datasets except EgoSchema. We provide more implementation details in

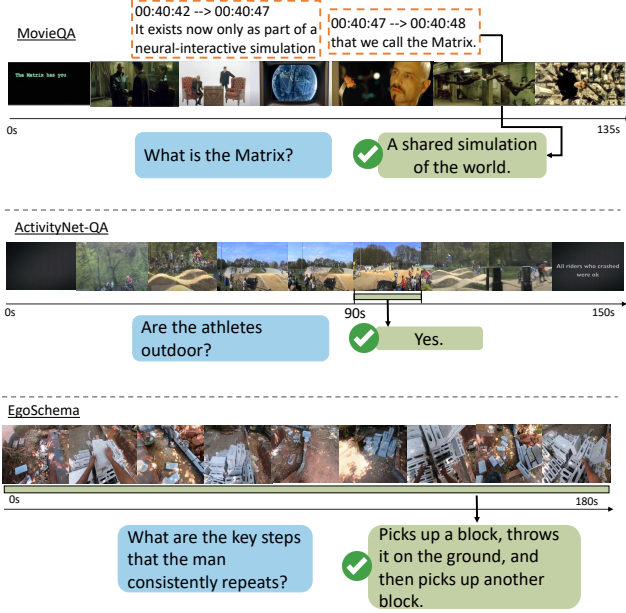


Figure 4. **An illustration of prior LVQA dataset/benchmark limitations.** **Top:** An example from MovieQA [52]. The model can use the provided subtitle information to answer a question while completely ignoring visual cues in a video. **Middle:** An example from ActivityNet-QA Dataset [80]. Despite long video inputs, the model only needs to analyze a short 1s video clip to answer the question. **Bottom:** An example from EgoSchema Dataset [37]. The model must analyze visual cues from the whole 3-minute video to answer a given question without relying on any additional textual inputs (e.g., speech, subtitles).

the Supplementary Material.

4. Experiments

4.1. Datasets and Metrics

Unlike short-term video question-answering, long-range video question-answering (LVQA) lacks robust and universally agreed-upon benchmarks. Many prior movie-based long-range video understanding benchmarks [8, 52] exhibit significant language biases, i.e., pure text-based approaches can achieve excellent performance by using information from subtitles/speech while ignoring the video content entirely [21, 75] (See the example in the top row of Figure 4). Furthermore, while many existing long-range video understanding benchmarks [38, 65, 80] contain long video inputs, prior work [37] have shown that all of these benchmarks can be solved by analyzing only several second-long video clips within the longer video inputs, without requiring any long-range video modeling capabilities (See the example in the middle row of Figure 4).

To address these limitations, recent work introduced **EgoSchema** [37], a new long-range video question-answering benchmark, consisting of 5K multiple choice

question-answer pairs, spanning 250 hours of video, and covering a wide range of human activities. Unlike prior benchmarks in this area, the authors in [37] have manually verified that to answer a given question correctly, one must watch at least 100 seconds of video content, which is orders of magnitude longer than any existing benchmark (See the example in the bottom row of Figure 4). Furthermore, unlike prior LVQA benchmarks, EgoSchema has no textual inputs such as speech, subtitles, or storylines. This means the questions must be answered based only on video inputs, preventing language-based biases. The EgoSchema benchmark consists of 5,000 questions, each requiring the correct answer to be selected between 5 given options based on a three-minute-long video clip. The entire dataset is designed for zero-shot evaluation and has no training set. For validation, the authors released a subset of 500 questions with ground truth answers (EgoSchema Subset). By default, our experiments are conducted on the EgoSchema Subset. The metric we use is QA accuracy, i.e., the percentage of correctly answered questions among all questions.

In addition to **EgoSchema**, we also perform zero-shot LVQA experiments on three other LVQA benchmarks:

- **NExT-QA** [68] contains 5,440 videos with an average duration of 44s and 48K multi-choice questions and 52K open-ended questions. There are 3 different question types: Temporal, Causal, and Descriptive. Following common practice, we perform zero-shot evaluation on the validation set, which contains 570 videos and 5K multiple-choice questions.
- **IntentQA** [30] contains 4,303 videos and 16K multiple-choice question-answer pairs focused on reasoning about people’s intent in the video. We perform a zero-shot evaluation on the test set containing 2K questions.
- **NExT-GQA** [71] is an extension of NExT-QA with 10.5K temporal grounding annotations associated with the original QA pairs. The dataset was introduced to study whether the existing LVQA models can temporally localize video segments needed to answer a given question. We evaluate all methods on the test split, which contains 990 videos with 5,553 questions, each accompanied by a temporal grounding label. The metrics we used include: 1) Intersection over Prediction (IoP) [71], which measures whether the predicted temporal window lies inside the ground truth temporal segment, 2) temporal Intersection over Union (IoU), and 3) Acc@GQA, which depicts the percentage of accurately answered and grounded predictions. For IoP and IoU, we report the mean values and values with the overlap thresholds of 0.5.

To study the most important factors in our framework design, we first conduct an empirical study on the released subset of EgoSchema. Afterward, we present our main results on the full EgoSchema test set. Lastly, we present our results on other datasets and tasks.

Captioner	Caption Type	Ego4D Pre-training	Acc. (%)
VideoBLIP	clip-level	✓	40.0
EgoVLP	clip-level	✓	46.6
BLIP-2	frame-level	✗	46.7
LaViLa	clip-level	✓	51.8
Oracle	clip-level	-	65.8

Table 1. **Accuracy of our framework with different visual captioners.** LaViLa visual captioner achieves the best results, outperforming other clip-level (e.g., EgoVLP, VideoBLIP) and image-level (e.g., BLIP-2) captioners. We also observe that the Oracle baseline using ground truth captions greatly outperforms all other variants, suggesting that our framework can benefit from the future development of visual captioners.

4.2. Empirical Study on EgoSchema

Before presenting our main results, we first study the effectiveness of different components within our LLoVi framework, including (i) the visual captioner, (ii) the LLM, (iii) the LLM prompt design, and (iv) few-shot in-context learning. The experiments are conducted on the EgoSchema Subset with 500 multi-choice questions. We discuss our empirical findings below. We also include additional experiments in the supplementary material.

4.2.1 Visual Captioning Model

In Table 1, we study the effectiveness of various clip-level video captioners, including LaViLa [83], EgoVLP [43], and VideoBLIP [78]. In addition to video captioners, we also try the state-of-the-art image captioner, BLIP-2 [28]. Lastly, to study the upper bound of our visual captioning results, we include the ground truth Oracle captioning baseline obtained from the Ego4D dataset. All baselines in Table 1 use similar experimental settings, including the same LLM model, i.e., GPT-3.5. The results are reported as LVQA accuracy on the EgoSchema Subset.

Based on the results in Table 1, we observe that LaViLa is the best captioning model, outperforming BLIP-2, EgoVLP, and VideoBLIP. We also observe that despite not being pre-trained on Ego4D [12], BLIP-2 performs reasonably well (**46.7%**) and even outperforms other strong Ego4D-pretrained baselines, EgoVLP and VideoBLIP. Lastly, the Oracle baseline with ground truth captions outperforms LaViLa captions by a large margin (**14.0%**). This demonstrates that our framework can benefit from future improvements in visual captioning models.

4.2.2 Large Language Model

In Table 2, we analyze the performance of our LLoVi framework using different LLMs while fixing the visual captioner

LLM	Model Size	Acc. (%)
Llama2-7B	7B	34.0
Llama2-13B	13B	40.4
Llama2-70B	70B	50.6
GPT-3.5	175B	51.8
GPT-4	N/A	58.3

Table 2. **Accuracy of our framework with different LLMs.** GPT-4 achieves the best accuracy, suggesting that stronger LLMs perform better in LVQA. However, due to the best accuracy and cost tradeoff, we use GPT3.5 for our remaining experiments.

to be LaViLa. Based on these results, we observe that GPT-4 achieves the best performance (**58.3%**), followed by GPT-3.5 (**51.8%**). These results suggest that stronger LLMs (GPT-4) are better at long-range modeling, as indicated by a significant margin in LVQA accuracy between GPT-4 and all other LLMs (**>6.5%**). We also note that the Llama2 performs reasonably well with its 70B variant (**50.6%**), but its performance drastically degrades with smaller capacity LLMs (i.e., Llama2-7B, Llama2-13B). Due to the tradeoff between accuracy and cost, we use GPT-3.5 for all of our remaining experiments unless noted otherwise.

4.2.3 LLM Prompt Analysis

In this section, we (1) analyze several variants of our summarization-based prompt (described in Section 3), and (2) experiment with other commonly used prompt designs, including Zero-shot Chain-of-Thought (Zero-shot CoT) [64], Plan-and-Solve [57], and Self-Consistency [59]. Below, we present a detailed analysis of these results.

A Multi-round Summarization-based Prompt. As discussed in Section 3, we found that a specialized multi-round prompt that first asks the LLM to summarize the noisy short-term captions, and then answer the question using the LLM-generated summary performs better than our base approach that uses caption inputs directly for LVQA. Given a concatenated set of captions C , an input question Q , and a set of candidate answers A , there are several input combinations that we can use to obtain the summary S . Thus, here, we investigate three distinct variants of obtaining summaries S :

1. $(C) \rightarrow S$: the LLM uses caption-only inputs C to obtain summaries S in the first round of prompting.
2. $(C, Q) \rightarrow S$: the LLM uses captions C and a question Q as inputs for generating summaries S . Having additional question inputs is beneficial as it allows the LLM to generate a summary S specifically tailored for answering an input question Q .
3. $(C, Q, A) \rightarrow S$: the LLM takes captions C , a question Q , and the answer candidates A as its inputs to produce summaries S . Like above, having additional answer can-

Variant	Standard	(C) \rightarrow S	(C, Q) \rightarrow S	(C, Q, A) \rightarrow S
Acc. (%)	51.8	53.6	57.6	55.9

Table 3. **Different variants of our multi-round summarization-based prompt.** Our results indicate that the (C, Q) \rightarrow S variant that takes concatenated captions C and a question Q for generating a summary S works the best, significantly outperforming (+5.8%) the standard prompt. This confirms our hypothesis that additional inputs in the form of a question Q enable the LLM to generate a summary S tailored to a given question Q .

Number of words	50	100	300	500	700
Acc. (%)	55.6	57.4	55.8	57.6	55.0

Table 4. **Number of words in a generated summary.** We study the optimal number of words in an LLM-generated summary. These results suggest that the optimal LVQA performance is obtained when using 500-word summaries.

didate inputs enables the LLM to generate a summary S most tailored to particular question-answer pairs.

In Table 3, we explore the effectiveness of these three prompt variants. Our results show that all three variants significantly outperform our simple, yet already strong baseline that uses a standard LVQA prompt (described in Section 3). Specifically, we note that the variant (C) \rightarrow S that uses caption-only inputs to obtain the summaries outperforms the standard baseline by 1.8%. Furthermore, we observe that incorporating a given question as an input (i.e., the (C, Q) \rightarrow S variant) leads to the best performance (57.6%) with a significant 5.8% boost over the standard LVQA prompt baseline. This confirms our earlier intuition that having additional question Q inputs enables the LLM to generate a summary S specifically tailored for answering that question, thus leading to a big boost in LVQA performance. Lastly, we observe that adding answer candidates A as additional inputs (i.e., the (C, Q, A) \rightarrow S variant) leads to a drop in performance (-1.7%) compared with the (C, Q) \rightarrow S variant. We conjecture that this might be because the wrong answers in the candidate set A may mislead the LLM, leading to a suboptimal summary S .

We also investigate the optimal length of the generated summary S , and present these results in Table 4. Specifically, for these experiments, we ask the LLM to generate a summary S using a different number of words (as part of our prompt). We use the best performing (C, Q) \rightarrow S variant for these experiments. Our results indicate that using a very small number of words (e.g., 50) leads to a drop in performance, indicating that compressing the caption information too much hurts the subsequent LVQA performance. Similarly, generating summaries that are quite long (e.g., 700 words) also leads to worse results, suggesting that the filtering of the potentially noisy/redundant information in the captions is important for good LVQA performance. The best performance is obtained using 500-word summaries.

Prompting Technique	Acc. (%)
<i>Zero-shot</i>	
Standard	51.8
Zero-shot Chain-of-Thought [64]	53.2
Plan-and-Solve [57]	54.2
Self-Consistency [59]	55.4
Ours	57.6
<i>Few-shot</i>	
Standard	56.5
Ours	61.7

Table 5. **Comparison with commonly used prompting techniques.** The “Standard” means a standard LVQA prompt is used (see Section 3). We demonstrate that our framework benefits from more sophisticated prompting techniques. Our multi-round summarization-based prompt achieves the best performance in both zero-shot and few-shot learning settings.

Comparison with Commonly Used Prompts. Next, in Table 5, we compare our multi-round summarization-based prompt with other commonly used prompts such as Zero-shot Chain-of-Thought [64], Plan-and-Solve [57], and Self-Consistency [59]. From these results, we observe that all of these prompts outperform the base variant of our model that uses a standard prompt. In particular, among these commonly used prompts, the self-consistency prompting technique achieves the best results (55.4%). Nevertheless, our multi-round summarization-based prompt achieves the best performance (57.6%).

4.2.4 Few-shot In-Context Learning

In-context learning with LLMs has shown strong few-shot performance in many NLP tasks [2, 64]. In Table 5, we evaluate the few-shot in-context learning capabilities of our LLoVi framework. Our results show that our LLoVi framework greatly benefits from few-shot in-context learning. Specifically, the few-shot in-context learning leads to a 4.7% boost on the variant of our framework that uses a standard prompt and 4.1% boost on our advanced framework using a multi-round summarization-based prompt. We used 6 few-shot examples for our experiments as we found this configuration to produce the best performance.

4.3. Main Results on EgoSchema

In Table 6, we evaluate our best-performing LLoVi framework, developed using our empirical insights on the full EgoSchema test set containing 5K video samples. We compare our approach with prior state-of-the-art methods including InternVideo [61], mPLUG-Owl [77], VIO-LET [11], FrozenBiLM [73]. Based on these results, we observe that the best-performing zero-shot variant of our LLoVi framework achieves 50.3% accuracy, outperforming

Model	Acc. (%)
<i>Zero-shot</i>	
VIOLET [11]	19.9
FrozenBiLM [73]	26.9
mPLUG-Owl [77]	31.1
InternVideo [61]	32.1
LLoVi (Ours)	50.3
<i>Few-shot</i>	
LLoVi (Ours)	52.5

Table 6. **Results on the full set of EgoSchema.** The best-performing zero-shot variant of our LLoVi framework achieves **50.3%** accuracy, outperforming the previous best-performing InternVideo model by **18.2%**. For fair comparisons, we gray out our best few-shot variant.

Model	Causal (%)	Temporal (%)	Descriptive (%)	All (%)
VFC [39]	45.4	51.6	64.1	51.5
InternVideo [61]	43.4	48.0	65.1	49.1
ViperGPT [51]	-	-	-	60.0
SeViLA [79]	61.3	61.5	75.6	63.6
LLoVi (ours)	69.5	61.0	75.6	67.7

Table 7. **Zero-shot results on NeXT-QA.** LLoVi achieves **67.7%** accuracy, outperforming previous best-performing model SeViLA by **4.1%**. Notably, LLoVi excels at causal reasoning outperforming SeViLA by **8.2%** in the causal question category.

the previous best-performing InternVideo model by a large margin (**18.2%**). Additionally, we show that using few-shot in-context learning, our best variant gets further improvement. These results validate our design choice of leveraging the long-range modeling capabilities of LLMs for the LVQA task. Furthermore, since our proposed LLoVi framework is agnostic to the visual captioning model and an LLM it uses, we believe that in the future, we could further improve these results by leveraging more powerful visual captioners and LLMs.

4.4. Results on Other Datasets

Next, we demonstrate that our simple framework generalizes well to other LVQA benchmarks.

NeXT-QA. In Table 7, we evaluate LLoVi on the NeXT-QA [68] validation set in a zero-shot setting. We compare our approach with prior methods: VFC [39], InternVideo [61], ViperGPT [51], and SeViLA [79]. We observe that LLoVi outperforms the previous best-performing method, SeViLA by **4.1%**. Notably, in the Causal category, LLoVi achieves **8.2%** improvement. We conjecture this improvement comes from the simple 2-stage design of our LLoVi framework: captioning followed by LLM reasoning. By captioning the video, we are able to directly leverage the reasoning ability of the powerful LLMs and thus achieve good causal reasoning performance.

IntentQA. In Table 8, we evaluate our method on the IntentQA [30] test set. In our comparisons, we in-

Model	Acc. (%)
<i>Supervised</i>	
HQGA [69]	47.7
VGT [70]	51.3
BlindGPT [42]	51.6
CaVIR [29]	57.6
<i>Zero-shot</i>	
SeViLA [79]	60.9
LLoVi (ours)	64.0

Table 8. **Results on IntentQA.** Our zero-shot framework outperforms previous supervised methods by a large margin (**6.4%**). LLoVi also outperforms the recent zero-shot method, SeViLA, by **3.1%**.

Model	mIoP	IoP@0.5	mIoU	IoU@0.5	Acc@GQA
<i>Weakly-Supervised</i>					
IGV [32]	21.4	18.9	14.0	9.6	10.2
Temp[CLIP](NG+) [71]	25.7	25.5	12.1	8.9	16.0
FrozenBiLM (NG+) [71]	24.2	23.7	9.6	6.1	17.5
SeViLA [79]	29.5	22.9	21.7	13.8	16.6
<i>Zero-shot</i>					
LLoVi (ours)	37.3	36.9	20.0	15.3	24.3

Table 9. **Grounded LVQA results on NeXT-GQA.** We extend LLoVi to the grounded LVQA task and show that it outperforms prior weakly-supervised approaches on all evaluation metrics. For a fair comparison, we de-emphasize the models that were pre-trained on video-language grounding datasets.

clude several supervised methods (HQGA [69], VGT [70], BlindGPT [42], CaVIR [29]) and one recent zero-shot approach, SeViLA. From the results in Table 8, we observe that our method greatly outperforms all prior approaches, both in the fully supervised and zero-shot settings.

NeXT-GQA. In Table 9, we extend our framework to the grounded LVQA task and evaluate it on the NeXT-GQA [71] test set. We compare LLoVi with the weakly-supervised methods: IGV [32], Temp[CLIP](NG+) [71], FrozenBiLM(NG+) [71] and SeViLA [79]. These baselines are first trained on NeXT-GQA to maximize the QA accuracy, and then use ad-hoc methods [71] to estimate a relevant video segment for question-answering. Although LLoVi is not trained on NeXT-GQA, it still outperforms these weakly-supervised methods by a large margin on all evaluation metrics. These results demonstrate that in addition to LVQA, our framework can also be used to temporally ground its predictions for more explainable long-range reasoning.

5. Conclusion

In this work, we present a simple, yet highly effective LLM-based framework for long-range video question-answering (LVQA). We thoroughly evaluate various design choices of our approach and use our empirical insights to develop

a method that improves upon prior LVQA approaches by a significant margin on the newly introduced EgoSchema benchmark. We also demonstrate that our framework generalizes to other LVQA benchmarks such as NeXT-QA, IntentQA, and it can be extended to grounded LVQA tasks. While our paper does not provide any major technical contributions, we hope that our simple LVQA framework will help inspire new ideas and simplify model design in long-range video understanding.

References

- [1] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, page 4, 2021. 1, 2
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 3, 7
- [3] Jun Chen, Deyao Zhu, Kilichbek Haydarov, Xiang Li, and Mohamed Elhoseiny. Video chatcaptioner: Towards the enriched spatiotemporal descriptions. *arXiv preprint arXiv:2304.04227*, 2023. 2, 3
- [4] Yihong Chen, Yue Cao, Han Hu, and Liwei Wang. Memory enhanced global-local aggregation for video object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10337–10346, 2020. 1
- [5] Feng Cheng and Gedas Bertasius. Tallformer: Temporal action localization with a long-memory transformer. In *European Conference on Computer Vision*, pages 503–521. Springer, 2022. 1
- [6] Feng Cheng, Xizi Wang, Jie Lei, David Crandall, Mohit Bansal, and Gedas Bertasius. Vindlu: A recipe for effective video-and-language pretraining. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [7] Seongho Choi, Kyoung-Woon On, Yu-Jung Heo, Ahjeong Seo, Youwon Jang, Min Su Lee, and Byoung-Tak Zhang. Dramaqa: Character-centered video story understanding with hierarchical QA. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 1166–1174. AAAI Press, 2021. 3
- [8] Seongho Choi, Kyoung-Woon On, Yu-Jung Heo, Ahjeong Seo, Youwon Jang, Minsu Lee, and Byoung-Tak Zhang. Dramaqa: Character-centered video story understanding with hierarchical qa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1166–1174, 2021. 5
- [9] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 3
- [10] Jiwan Chung and Youngjae Yu. Long story short: a summarize-then-search method for long video question answering. In *BMVC*, 2023. 2, 3
- [11] Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. VIOLET: End-to-End Video-Language Transformers with Masked Visual-token Modeling. In *arXiv:2111.1268*, 2021. 1, 7, 8
- [12] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu 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. 6
- [13] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021. 3
- [14] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*, 2023. 1, 3
- [15] Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023. 1, 3
- [16] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019. 2
- [17] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Timeception for complex action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 254–263, 2019. 1
- [18] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Videograph: Recognizing minutes-long human activities in videos. *arXiv preprint arXiv:1905.05143*, 2019. 1, 3
- [19] Md Mohaiminul Islam and Gedas Bertasius. Long movie clip classification with state-space video models. In *European Conference on Computer Vision*, pages 87–104. Springer, 2022. 1, 3
- [20] Md Mohaiminul Islam, Mahmudul Hasan, Kishan Shamsundar Athrey, Tony Braskich, and Gedas Bertasius. Efficient movie scene detection using state-space transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18749–18758, 2023. 1
- [21] Bhavan Jasani, Rohit Girdhar, and Deva Ramanan. Are we asking the right questions in movieqa? In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3, 5
- [22] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 2
- [23] Sangho Lee, Jinyoung Sung, Youngjae Yu, and Gunhee Kim. A memory network approach for story-based temporal summarization of 360 videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1410–1419, 2018. 1

- [24] Sangmin Lee, Hak Gu Kim, Dae Hwi Choi, Hyung-Il Kim, and Yong Man Ro. Video prediction recalling long-term motion context via memory alignment learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3054–3063, 2021. 1
- [25] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. Tvqa: Localized, compositional video question answering. In *EMNLP*, 2018. 3
- [26] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, 2021. 3
- [27] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 3
- [28] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2, 6, 1
- [29] Jiapeng Li, Ping Wei, Wenjuan Han, and Lifeng Fan. Inten-tqa: Context-aware video intent reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11963–11974, 2023. 8
- [30] Jiapeng Li, Ping Wei, Wenjuan Han, and Lifeng Fan. Inten-tqa: Context-aware video intent reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11963–11974, 2023. 5, 8
- [31] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding, 2023. 2, 3
- [32] Yicong Li, Xiang Wang, Junbin Xiao, Wei Ji, and Tat-Seng Chua. Invariant grounding for video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2928–2937, 2022. 8
- [33] 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. Mm-vid: Advancing video understanding with gpt-4v(ision). *arXiv preprint arXiv:2310.19773*, 2023. 2, 3
- [34] Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, et al. Mm-vid: Advancing video understanding with gpt-4v (ision). *arXiv preprint arXiv:2310.19773*, 2023. 3
- [35] Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. Llm+ p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*, 2023. 1, 3
- [36] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 4
- [37] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *arXiv preprint arXiv:2308.09126*, 2023. 2, 3, 5
- [38] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2630–2640, 2019. 5
- [39] Liliane Momeni, Mathilde Caron, Arsha Nagrani, Andrew Zisserman, and Cordelia Schmid. Verbs in action: Improving verb understanding in video-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15579–15591, 2023. 8
- [40] Eric Nguyen, Karan Goel, Albert Gu, Gordon Downs, Preeti Shah, Tri Dao, Stephen Baccus, and Christopher Ré. S4nd: Modeling images and videos as multidimensional signals with state spaces. *Advances in neural information processing systems*, 35:2846–2861, 2022. 3
- [41] OpenAI. Gpt-4 technical report, 2023. 3
- [42] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022. 8
- [43] Kevin Qinghong Lin, Alex Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Zhongcong Xu, Difei Gao, Rongcheng Tu, Wenzhe Zhao, Weijie Kong, et al. Egocentric video-language pretraining. *arXiv e-prints*, pages arXiv–2206, 2022. 2, 3, 6, 1
- [44] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 2
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [46] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 3
- [47] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009, 2023. 1, 3
- [48] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. Moviechat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 3
- [49] Simeng Sun, Yang Liu, Shuohang Wang, Chenguang Zhu, and Mohit Iyyer. Pearl: Prompting large language models to plan and execute actions over long documents. *arXiv preprint arXiv:2305.14564*, 2023. 1, 3
- [50] Yuchong Sun, Hongwei Xue, Ruihua Song, Bei Liu, Huan Yang, and Jianlong Fu. Long-form video-language pre-training with multimodal temporal contrastive learning. *Ad-*

- vances in neural information processing systems, 35:38032–38045, 2022. 3
- [51] D     Sur  s, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2023. 3, 8
- [52] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhofen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4631–4640, 2016. 3, 5
- [53] Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Steven Zheng, et al. U12: Unifying language learning paradigms. In *The Eleventh International Conference on Learning Representations*, 2022. 3
- [54] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 3
- [55] Junke Wang, Dongdong Chen, Chong Luo, Xiyang Dai, Lu Yuan, Zuxuan Wu, and Yu-Gang Jiang. Chatvideo: A tracklet-centric multimodal and versatile video understanding system. *arXiv preprint arXiv:2304.14407*, 2023. 2, 3
- [56] Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. Selective structured state-spaces for long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6387–6397, 2023. 1, 3
- [57] Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2609–2634, Toronto, Canada, 2023. Association for Computational Linguistics. 6, 7
- [58] Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*, 2023. 1, 3
- [59] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *ICLR*, 2023. 6, 7, 2
- [60] Yang Wang, Gedas Bertasius, Tae-Hyun Oh, Abhinav Gupta, Minh Hoai, and Lorenzo Torresani. Supervoxel attention graphs for long-range video modeling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 155–166, 2021. 1, 3
- [61] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 1, 3, 7, 8
- [62] Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chengguang Zhu, Derek Hoiem, et al. Language models with image descriptors are strong few-shot video-language learners. *Advances in Neural Information Processing Systems*, 35: 8483–8497, 2022. 2, 3
- [63] Ziyang Wang, Yi-Lin Sung, Feng Cheng, Gedas Bertasius, and Mohit Bansal. Unified coarse-to-fine alignment for video-text retrieval. *arXiv preprint arXiv:2309.10091*, 2023. 1
- [64] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022. 6, 7
- [65] Chao-Yuan Wu and Philipp Krahenbuhl. Towards long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1884–1894, 2021. 5
- [66] Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krahenbuhl, and Ross Girshick. Long-term feature banks for detailed video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 284–293, 2019. 1
- [67] Chao-Yuan Wu, Yanghao Li, Karttikeya Mangalam, Haoqi Fan, Bo Xiong, Jitendra Malik, and Christoph 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. 1, 3
- [68] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng 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. 3, 5, 8
- [69] Junbin Xiao, Angela Yao, Zhiyuan Liu, Yicong Li, Wei Ji, and Tat-Seng Chua. Video as conditional graph hierarchy for multi-granular question answering. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI)*, pages 2804–2812, 2022. 8
- [70] Junbin Xiao, Pan Zhou, Tat-Seng Chua, and Shuicheng Yan. Video graph transformer for video question answering. In *European Conference on Computer Vision*, pages 39–58. Springer, 2022. 8
- [71] Junbin Xiao, Angela Yao, Yicong Li, and Tat Seng Chua. Can i trust your answer? visually grounded video question answering. *arXiv preprint arXiv:2309.01327*, 2023. 5, 8
- [72] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just ask: Learning to answer questions from millions of narrated videos. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1686–1697, 2021. 3
- [73] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Zero-shot video question answering via frozen bidirectional language models. In *NeurIPS*, 2022. 1, 7, 8
- [74] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Zero-shot video question answering via

- frozen bidirectional language models. *Advances in Neural Information Processing Systems*, 35:124–141, 2022. 3
- [75] Jianing Yang, Yuying Zhu, Yongxin Wang, Ruitao Yi, Amir Zadeh, and Louis-Philippe Morency. What gives the answer away? question answering bias analysis on video qa datasets. *arXiv preprint arXiv:2007.03626*, 2020. 3, 5
- [76] Xitong Yang, Fu-Jen Chu, Matt Feiszli, Raghav Goyal, Lorenzo Torresani, and Du Tran. Relational space-time query in long-form videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6398–6408, 2023. 1, 3
- [77] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023. 1, 7, 8
- [78] Keunwoo Peter Yu. VideoBLIP. 3, 6
- [79] Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. *NeurIPS*, 2023. 3, 8
- [80] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ing Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9127–9134, 2019. 3, 5
- [81] Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv*, 2022. 2, 3
- [82] Chuhan Zhang, Ankush Gupta, and Andrew Zisserman. Temporal query networks for fine-grained video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4486–4496, 2021. 1
- [83] Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit 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. 2, 3, 4, 6

Our appendix consists of Additional Analysis (Section A), Additional Implementation Details (Section B) and Qualitative Analysis (Section C).

A. Additional Analysis

In this section, we provide additional analysis on the EgoSchema Subset using the standard prompt.

A.1. Video Sampling Configurations

In Figure 5, we investigate the sensitivity of LVQA performance on EgoSchema with different video sampling configurations. Specifically, in Subfigure 5a, we experiment with 4 different clip lengths: 0.5s, 1s, 4s, and 8s. For each clip length, we use the stride that would be sufficient to cover the entire long video input. For these experiments, we use a LaViLa visual captioner and a GPT-3.5 LLM. Our results indicate that LVQA performance is the best when the sampled clip length is 1s. We observe that using an even shorter video clip length (i.e., 0.5s) produces many repetitive/redundant captions, which leads to 2% drop in LVQA performance. Furthermore, we also note that increasing the clip length to longer durations (e.g., 2s-8s) makes the accuracy lower since the extracted captions start to lack detailed visual information needed to answer the question. In addition, in Subfigure 5b, we fix the clip length to 1s and experiment with 4 different stride values: 1s, 2s, 4s, and 8s. Note that the 1s clips sampled using a 1s stride will cover the entire video without overlap. Our results suggest a gradual decrease in LVQA accuracy when increasing the stride from 1s to 8s. This indicates that having gaps in long video coverage leads to suboptimal LVQA performance. Thus, for the rest experiments, we use 1s video clips sampled with a 1s stride.

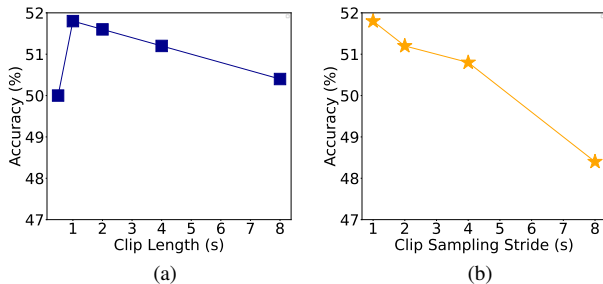


Figure 5. **The Analysis of Video Clip Sampling Strategy on EgoSchema.** These results show that sampling 1s video clips with a 1s stride leads to the best LVQA performance.

A.2. Accuracy on Different Question Types

To better understand the strengths and limitations of our LVQA framework, we manually categorize questions in the EgoSchema Subset into 5 categories: (1) Purpose/Goal

Question Category	Category Percentage(%)	Acc.(%)
Purpose/Goal Identification	49.2	54.9
Tools and Materials Usage	21.8	50.5
Key Action/Moment Detection	21.6	43.5
Action Sequence Analysis	18.2	52.7
Character Interaction	9.4	63.8

Table 10. **Accuracy on different question categories of EgoSchema.** We manually categorize each question in the EgoSchema Subset into 5 categories. Note that each question may belong to one or more categories. Our system performs the best on questions that involve character interaction analysis or human purpose/goal identification. This is encouraging as both of these questions typically require a very long-form video analysis.

Identification, (2) Tools and Materials Usage, (3) Key Action/Moment Detection, (4) Action Sequence Analysis, (5) Character Interaction (see Supplementary Materials for details). Note that some questions belong to more than one category. Based on this analysis, we observe that almost half of the questions relate to purpose/goal identification, which makes intuitive sense as inferring human goals/intent typically requires a very long video analysis. We also observe that a significant portion of the questions relate to tool usage, key action detection, and action sequence analysis. Lastly, the smallest fraction of the questions belong to character interaction analysis.

In Table 10, we break down our system’s performance according to each of the above-discussed question categories. Our results indicate that our system performs the best in the Character Interaction category (**63.8%**). One possible explanation is that the LaViLa model, which we use as our visual captioner, is explicitly pretrained to differentiate the camera wearer from other people, making it well-suited for understanding various interactions between characters in the video. We also observe that our framework performs much worse in the Key Action/Moment Detection category (**43.5%**). We conjecture that this might be caused by the limitations in the visual captioning model, i.e., if the key action fails to appear in any of the visual captions, the question will be almost impossible to answer. Lastly, we note that our model performs quite well on the remaining categories (**>50%**). It is especially encouraging to see strong results (**54.9%**) in the Purpose/Goal Identification category since inferring human intentions/goals from the video inherently requires very long-form video analysis.

B. Additional Implementation Details

B.1. Captioners

For most experiments on EgoSchema, we use LaViLa as the visual captioner. For other pre-trained visual captioners, we use off-the-shelf pre-trained models, e.g., BLIP2 [28], EgoVLIP [43].

LaViLa is trained on the Ego4D dataset. The original LaViLa train set has 7743 videos with 3.9M video-text pairs and the validation set has 828 videos with 1.3M video-text pairs. The EgoSchema dataset is cropped from Ego4D. Since EgoSchema is designed for zero-shot evaluation and the original LaViLa train set includes EgoSchema videos, we retrain LaViLa on Ego4D videos that do not have any overlap with EgoSchema videos to avoid unfair comparison with other methods. After removing the EgoSchema videos, the train set consists 6100 videos with 2.3M video-text pairs, and the validation set has 596 videos with 0.7M video-text pairs. We retrain LaViLa on this reduced train set to prevent data leakage. LaViLa training consists of two stages: 1) dual-encoder training and 2) narrator training. Below we provide more details.

Dual-encoder. We use TimeSformer [1] base model as the visual encoder and a 12-layer Transformer as the text encoder. The input to the visual encoder comprises 4 RGB frames of size 224×224 . We randomly sample 4 frames from the input video clip and use RandomResizedCrop for data augmentation. The video-language model follows a dual-encoder architecture as CLIP [45] and is trained contrastively. Following LaViLa [83], we use 1024 as batch size. We train at a 3×10^{-5} learning rate for 5 epochs on 32 NVIDIA RTX 3090 GPUs.

Narrator is a visually conditioned autoregressive Language Model. It consists of a visual encoder, a resampler module, and a text encoder. We use the visual encoder (TimeSformer [1] base model) from the pretrained dual-encoder (See the previous paragraph). The resampler module takes as input a variable number of video features from the visual encoder and produces a fixed number of visual tokens (i.e. 256). The text decoder is the pretrained GPT-2 [44] base model with a cross-attention layer inserted in each transformer block which attends to the visual tokens of the resampler module. We freeze the visual encoder and the text decoder, while only training the cross-attention layers of the decoder and the resampler module. Following the design in LaViLa [83], we use a batch size of 256 and a learning rate of 3×10^{-5} . We use AdamW optimizer [22] with $(\beta_1, \beta_2) = (0.9, 0.999)$ and weight decay 0.01. We train the model on 8 NVIDIA RTX 3090 GPUs for 5 epochs.

Narrating video clips. We use nucleus sampling [16] with $p = 0.95$ and return $K = 5$ candidate outputs. Then we take the narration with the largest confidence score as the final caption of the video clip.

For NExT-QA, IntentQA and NExT-GQA datasets, we use LLaVA1.5 as the visual captioner and GPT-4 as the LLM. Specifically, we use the `llava-1.5-7b-hf` variant with the prompt “*USER: <image>. Describe the image in 30 words. ASSISTANT:* ”.

B.2. LLMs

For most experiments on EgoSchema we use GPT-3.5 as the LLM. Specifically, we use the `gpt-3.5-turbo-0613` variant which has 4K context. When the context length is not enough, we use the `gpt-3.5-turbo-16k` variant. We use 0 as temperature for all experiments.

We use `Llama-2-7b-chat-hf`, `Llama-2-13b-chat-hf`, and `Llama-2-70b-chat-hf` variants as Llama2 models. For all Llama2 models, we use greedy sampling to generate the output.

For NExT-QA, IntentQA and NExT-GQA datasets, we use GPT-4 as the LLM with the variant `gpt-4-1106-preview`.

B.3. Prompting Techniques Implementation

Prompt Details. We provide detailed prompts for our standard prompt in Table 11, multi-round summarization-based prompt in Table 12, Zero-shot Chain of Thought in Table 13, and Plan-and-Solve prompting in Table 14. To implement Self-Consistency, we set the temperature of GPT3.5 to 0.7 and run Zero-shot Chain of Thought for 5 times following the design in Self-Consistency [59]. Each run of the model provides a result, and the final output is determined by a majority vote. The prompt for the grounded LVQA benchmark is shown in Table 15.

Output Processing. When answering multiple choice questions, GPT3.5 usually outputs complete sentences instead of a single-letter answer, i.e. A, B, C, D, or E. One way to obtain the single-character response is to perform post-processing on the output, which usually requires substantial engineering efforts. In our work, however, we observe that GPT3.5 is very sensitive to the starting sentences of the prompts. Therefore, we explicitly prompt it as in Table 11 to force GPT3.5 to generate a single character as response. In practice, we take out the first character of the output as the final answer.

User

Please provide a single-letter answer (A, B, C, D, E) to the following multiple-choice question, and your answer must be one of the letters (A, B, C, D, or E). You must not provide any other response or explanation. You are given some language descriptions of a first person view video. The video is 3 minute long. Each sentence describes a `clip.length` clip. Here are the descriptions: `Captions`

You are going to answer a multiple choice question based on the descriptions, and your answer should be a single letter chosen from the choices.

Here is the question: `Question`

Here are the choices. A: `Option-A`. B: `Option-B`. C: `Option-C`. D: `Option-D`. E: `Option-E`.

Assistant

`Answer`

Table 11. LLoVi with Standard Prompt on EgoSchema.

User

You are given some language descriptions of a first person view video. Each video is 3 minute long. Each sentence describes a `clip.length` clip. Here are the descriptions: `Captions`

Please give me a `num.words` words summary. When doing summarization, remember that your summary will be used to answer this multiple choice question: `Question`.

Assistant

`Summary`

User

Please provide a single-letter answer (A, B, C, D, E) to the following multiple-choice question, and your answer must be one of the letters (A, B, C, D, or E). You must not provide any other response or explanation. You are given some language descriptions of a first person view video. The video is 3 minute long. Here are the descriptions: `Summary`

You are going to answer a multiple choice question based on the descriptions, and your answer should be a single letter chosen from the choices.

Here is the question: `Question`

Here are the choices. A: `Option-A`. B: `Option-B`. C: `Option-C`. D: `Option-D`. E: `Option-E`.

Assistant

`Answer`

Table 12. LLoVi with Multi-round Summarization-based Prompt on EgoSchema. We show the variant (C, Q) \rightarrow S, where we feed the question without potential choices to the summarization stage. **Top:** caption summarization prompt. **Bottom:** question answering prompt. In the first stage, GPT3.5 outputs a question-guided summary. In the second stage, GPT3.5 takes the summary without the original captions, then answer the multiple choice question.

User

You are given some language descriptions of a first person view video. The video is 3 minute long. Each sentence describes a `clip.length` clip. Here are the descriptions: `Captions`

You are going to answer a multiple choice question based on the descriptions, and your answer should be a single letter chosen from the choices.

Here is the question: `Question`

Here are the choices. A: `Option-A`. B: `Option-B`. C: `Option-C`. D: `Option-D`. E: `Option-E`.

Before answering the question, let's think step by step.

Assistant

`Answer` and `Rationale`

User

Please provide a single-letter answer (A, B, C, D, E) to the multiple-choice question, and your answer must be one of the letters (A, B, C, D, or E). You must not provide any other response or explanation. Your response should only contain one letter.

Assistant

`Answer`

Table 13. LLoVi with Zero-shot Chain of Thought Prompting on EgoSchema.

User

You are given some language descriptions of a first person view video. The video is 3 minute long. Each sentence describes a `clip.length` clip. Here are the descriptions: `Captions`

You are going to answer a multiple choice question based on the descriptions, and your answer should be a single letter chosen from the choices.

Here is the question: `Question`

Here are the choices. A: `Option-A`. B: `Option-B`. C: `Option-C`. D: `Option-D`. E: `Option-E`.

To answer this question, let's first prepare relevant information and decompose it into 3 sub-questions. Then, let's answer the sub-questions one by one. Finally, let's answer the multiple choice question.

Assistant

`Sub-questions` and `Sub-answers`

User

Please provide a single-letter answer (A, B, C, D, E) to the multiple-choice question, and your answer must be one of the letters (A, B, C, D, or E). You must not provide any other response or explanation. Your response should only contain one letter.

Assistant

`Answer`

Table 14. LLoVi with Plan-and-Solve Prompting on EgoSchema.

User

I will provide video descriptions and one question about the video. The video is 1 FPS and the descriptions are the captions every 2 frames. Each caption starts with the frame number. To answer this question, what is the minimum frame interval to check? Follow this format: [frame_start_index, frame_end_index]. Do not provide any explanation.

Here are the descriptions: [Captions](#)

Here is the question: [Question](#)

Please follow the output format as follows: #Example1: [5, 19]. #Example2: [30, 60]. #Example3: [1, 10] and [50, 60]

Assistant

[Answer](#)

Table 15. LLoVi Prompt on NExT-GQA.





				
LaViLa	#C C drops the brick mould.	#O man X moves the cards.	#C C puts the cloth on the table.	#C C moves the dough in the tray.
BLIP2	A person is laying a brick in the dirt.	A child is playing a game of monopoly with a tray of paper plates.	A person is working on a tool.	Woman making dough in a kitchen.

Table 16. **Comparison between different captioners.** **Top:** frames from EgoSchema videos. **Middle:** captions generated by LaViLa. **Bottom:** captions generated by BLIP2. LaViLa captions are more concise than BLIP2 captions. LaViLa is better at differentiating the camera wearer and other people.

C. Qualitative Analysis

C.1. Captioners

In Table 16 we compare different captions generated by BLIP2 and LaViLa on EgoSchema. LaViLa captions are generally more concise than BLIP2 captions, focusing more on the actions while BLIP2 focuses more on describing the objects. We also observe that LaViLa is better at differentiating the camera wearer and other person. As shown in the second image in Table 16, LaViLa tends to focus more on the action of the other person when the camera wearer and other person both appear in the video.

C.2. LLoVi with Standard Prompt

We show two examples of our method with standard prompt, including a successful one and a failed one in Figure 6. Our method performs long-range modeling from short-term video captions through LLM to understand the video. In the success case demonstrated in Subfigure 6a, the captions describe the camera wearer’s action in a short period of time, such as the interaction with the tape measure and the wood. With the short-term captions, LLM understand the long video and answers the question correctly.

In the failure case shown in Subfigure 6b, although the video captioner identifies the object in the video correctly as a tablet, LLM understands the action of the camera wearer

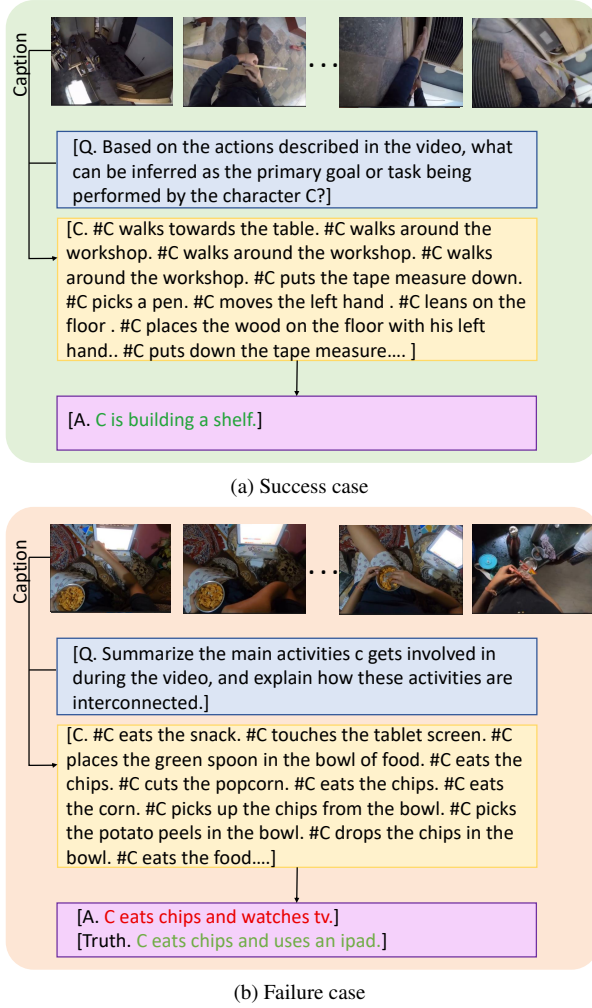


Figure 6. **Examples of our framework with a standard prompt on EgoSchema.** We show two examples, a successful one (a) and a failed one (b).

as watching TV rather than using an iPad. This might be caused by misguidance from the redundant captions that are not related to the question.

C.3. LLoVi with Multi-round Summarization-based Prompt

Figure 7 illustrates two EgoSchema questions that our framework with multi-round summarization-based prompt answers correctly. In Subfigure 7a, the question asks for the primary function of a tool that the video taker uses. However, shown in the first two images, the long video contains descriptions that are not related to the question, such as operating a machine and rolling a dough. As a result, the generated text captions would contain a large section that is not our direction of interest. By summarizing the captions with awareness to the question, LLM extracts key informa-

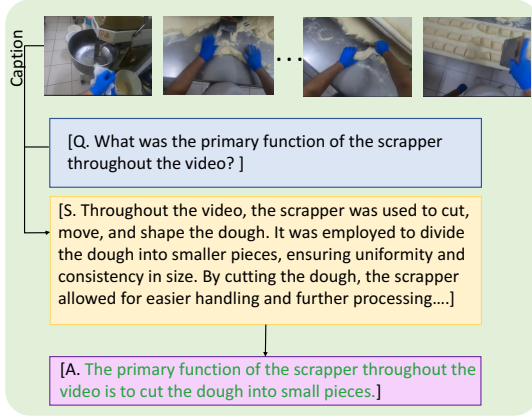
tion and cleans redundant captions to provide clearer textual background for answering the question. The same pattern is observed in Subfigure 7b.

Figure 8 shows two questions that our method fails to answer. In the summarization stage, the LLM answers the question directly instead of using the question to guide the summarization. For example, in Subfigure 8a, all the frames show the camera wearer engaging in actions related to washing dishes, but LLM infers that the person is cleaning the kitchen in the summarization stage. This wrong inference further misdirects the following question answering stage, which leads to an incorrect answer. In Subfigure 8b, LLM concludes that the cup of water is used to dilute the paint because the camera wearer dips the brush into water before dipping it into the paint palette.

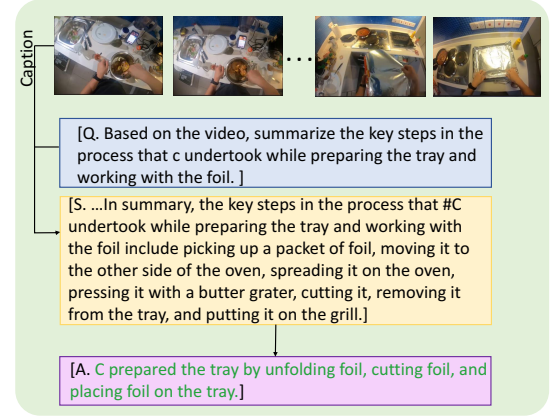
In Figure 9, we also show a question which the standard prompt fails to answer, but the multi-round summarization-based prompt answers correctly. In the video in the example question, we observe the camera wearer involving in activities related to laundry, such as picking up clothes from the laundry basket and throwing them into the washing machine. However, the short-term video captions shown in Subfigure 9a demonstrate the redundancy of actions. The repetitive actions complexes extracting and comprehending the information presented in the caption. For example, excessive captions on picking up clothes can make LLM think that the camera wearer is packing something. Our multi-round summarization-based prompt mitigate this problem by first ask LLM to provide a summary of the captions. The summary shown in Subfigure 9b states clearly that the camera wearer is doing laundry. With the cleaner and more comprehensive summary, the LLM answer the question correctly.

C.4. Question Categories

We provide detailed descriptions of each question category in Table 17. Note that each question can be classified into multiple categories.

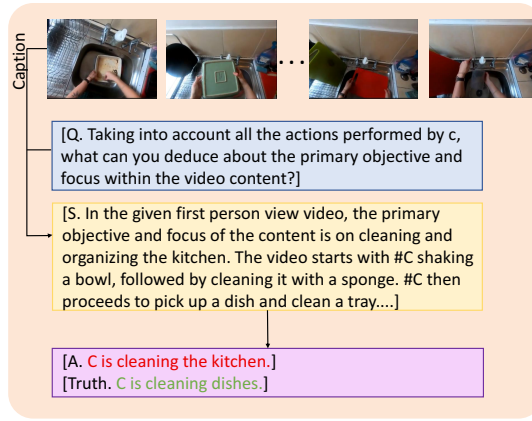


(a)

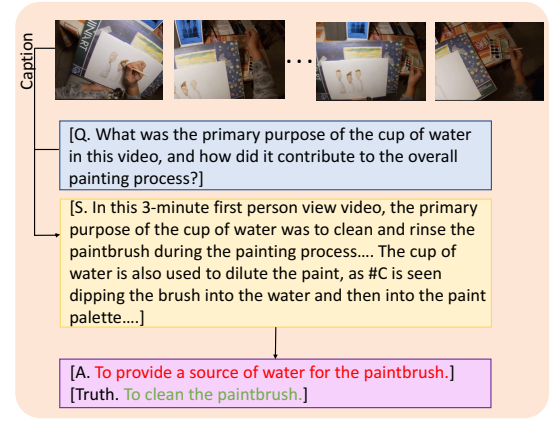


(b)

Figure 7. Success cases of our multi-round summarization-based prompt.

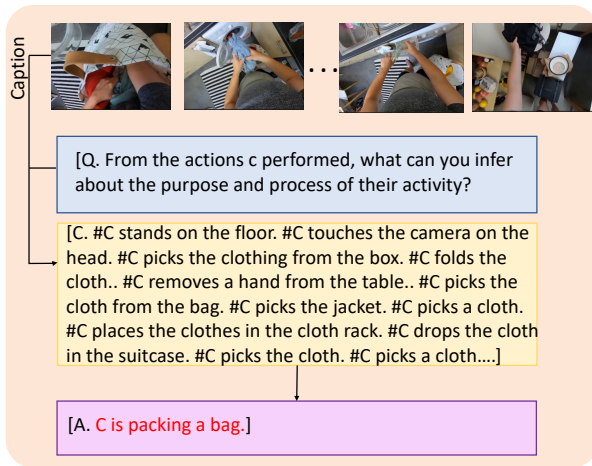


(a)

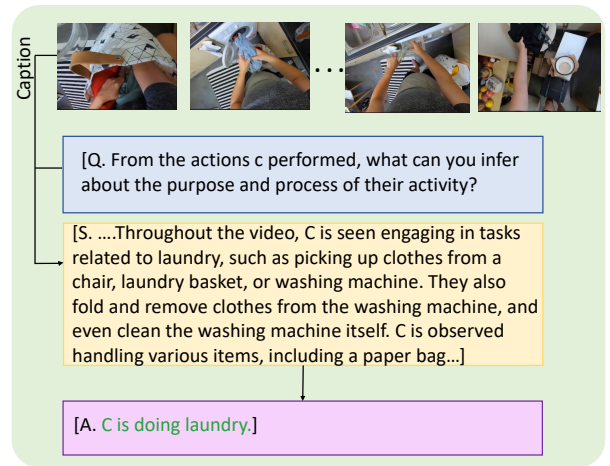


(b)

Figure 8. Failure cases of our framework with multi-round summarization-based prompt.



(a) Standard prompt (wrong answer).



(b) Multi-round summarization-based prompt (correct answer).

Figure 9. Contrast between our standard prompt and our multi-round summarization-based prompt. (a) demonstrates the process of answering the question with a standard prompt, and (b) shows answering the question with our multi-round summarization-based prompt.

Question Category	Description	Examples
Purpose/Goal Identification	primary goals, intentions, summary, or overarching themes of the video	<ol style="list-style-type: none"> 1. Taking into account all the actions performed by c, what can you deduce about the primary objective and focus within the video content? 2. What is the overarching theme of the video, considering the activities performed by both characters?
Tools and Materials Usage	how the character engages with specific tools, materials, and equipment	<ol style="list-style-type: none"> 1. What was the primary purpose of the cup of water in this video, and how did it contribute to the overall painting process? 2. Explain the significance of the peeler and the knife in the video and their respective roles in the preparation process.
Key Action / Moment Detection	identify crucial steps/actions, the influence/rationale of key action/moment/change on the whole task	<ol style="list-style-type: none"> 1. Out of all the actions that took place, identify the most significant one related to food preparation and explain its importance in the context of the video. 2. Identify the critical steps taken by c to organize and prepare the engine oil for use on the lawn mower, and highlight the importance of these actions in the overall video narrative.
Action Sequence Analysis	compare and contrast different action sequences, relationship between different actions, how characters adjust their approach, efficacy and precision, expertise of the character	<ol style="list-style-type: none"> 1. What is the primary sequence of actions performed by c throughout the video, and how do these actions relate to the overall task being performed? 2. Considering the sequence of events, what can be inferred about the importance of precision and accuracy in the character's actions, and how is this demonstrated within the video?
Character Interaction	how characters interact and collaborate, how their roles differ	<ol style="list-style-type: none"> 1. What was the main purpose of the actions performed by both c and the man throughout the video, and how did their roles differ? 2. Describe the general activity in the room and how the different characters and their actions contribute to this environment.

Table 17. **Question categories of EgoSchema.** We manually categorize each question in the EgoSchema Subset into 5 categories. Note that each question may belong to one or more categories.