

VideoAgent: A Memory-augmented Multimodal Agent for Video Understanding

Yue Fan^{*1}, Xiaojian Ma^{*1}, Rujie Wu^{1,2}, Yuntao Du¹, Jiaqi Li¹, Zhi Gao^{1,3}, and Qing Li¹

¹ National Key Laboratory of General Artificial Intelligence, Beijing Institute for General Artificial Intelligence (BIGAI)

² School of Computer Science, Peking University

³ School of Intelligence Science and Technology, Peking University

{maxiaojian,liqing}@bigai.ai

videoagent.github.io

Abstract. We explore how reconciling several foundation models (large language models and vision-language models) with a novel unified memory mechanism could tackle the challenging video understanding problem, especially capturing the long-term temporal relations in lengthy videos. In particular, the proposed multimodal agent *VideoAgent*: 1) constructs a structured memory to store both the generic temporal event descriptions and object-centric tracking states of the video; 2) given an input task query, it employs tools including video segment localization and object memory querying along with other visual foundation models to interactively solve the task, utilizing the zero-shot tool-use ability of LLMs. *VideoAgent* demonstrates impressive performances on several long-horizon video understanding benchmarks, on averaged increasing 6.6% on NExT-QA and 26.0% on EgoSchema over baselines, closes the gap between open-sourced models and private counterparts including Gemini 1.5 Pro.

Keywords: Video understanding · LLMs · tool-use · multimodal agents

1 Introduction

Understanding videos and answering free-form queries (question answering, content retrieval, *etc.*) remains a major challenge in computer vision and AI [1, 12, 15, 22, 23, 26, 31, 44]. Notably, much of the recent progress has achieved by the end-to-end pretrained large transformer models, especially those are developed upon the powerful large language models (LLMs) [4, 11, 22, 31], *i.e.* multimodal LLMs. However, there have been increasing concerns about their capabilities to handle long-form videos with rich events and complex spatial-temporal dependencies [7, 9, 16, 24, 33]. Specifically, the computation, especially memory cost could grow significantly and even become prohibitively expensive when processing lengthy videos [26, 32]. Also, the self-attention mechanism could sometimes struggle to capture the long-range relations [25]. These issues have hindered further advancement in applying sophisticated foundation models to video understanding.

^{*} Denotes equal contribution.

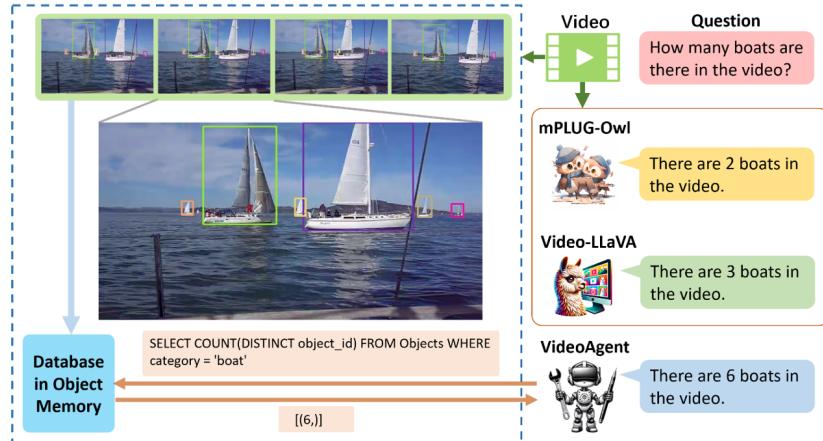


Fig. 1: A comparison between *VideoAgent* and end-to-end video-language models on video QA. Without a unified memory as a structured representation for videos, end-to-end models could struggle with capturing basic spatial-temporal details, especially when asked about objects and lengthy videos. While *VideoAgent* can utilize a curated set of tools to perform sophisticated queries about the *temporal memory* (not shown) and *object memory*, and respond with the correct answer.

More recently, thanks to the tool-use capabilities of LLMs [3, 20], there has been rapid development of a new class of multimodal understanding approaches: *multimodal agents* [6, 12, 23, 34]. The key idea is prompting LLMs into solving the multimodal tasks by invoking several **tool** foundation models (object detection, visual question answering, *etc.*) interactively. These methods have great potential as they are mostly training-free and flexible with tool sets. However, extending them to video understanding, especially on long-form videos is **non-trivial**. Simply adding video foundation models as tools could still suffer from the computation cost and attention limitation issues [11, 22]. Other research has explored more sophisticated prompting strategies with better tools [14, 30, 37], but they usually lead to complicated pipelines and the performances of these methods still fail to match their end-to-end counterparts possibly due to a lack of video-specific agent design.

In this paper, we introduce a simple yet effective LLM-based multimodal tool-use agent *VideoAgent* for video understanding tasks. Our **key insight** is to represent the video as a structured unified memory, therefore facilitating strong spatial-temporal reasoning and tool use of the LLM, and matching/outperforming end-to-end models. Our memory design is **motivated** by the principle of being minimal but sufficient: we've found that the overall event context descriptions and temporally consistent details about objects could cover the most frequent queries about videos. As a result, we design two memory components: 1) *temporal memory*, which stores text descriptions of each short (2 seconds) video segment sliced from the complete video; 2) *object memory*, where we track and store the occurrence of objects and persons in the video. To answer a query, the

LLM will decompose it into several subtasks and invoke the tool models. The unified memory is centered around the following tools:  *caption retrieval*, which will return all the event descriptions between two query time steps;  *segment localization*, which retrieves a short video segment of a given textual query by comparing it against the event descriptions within the temporal memory;  *visual question answering*, which answers a question given a retrieved video segment;  *object memory querying*, which allows sophisticated object state retrieval from the object memory using SQL queries. Finally, the LLM will aggregate the response of the interactive tool use and produce an answer to the input query.

We conduct extensive evaluations of *VideoAgent* on several video understanding tasks, including free-form query localization with Ego4D NLQ [5], generic video question answering with NExT-QA [35] and egocentric question answering with EgoSchema [15], a recent benchmark focusing on complex questions about long-form videos. We compare *VideoAgent* against both the canonical end-to-end video-language foundation models and other multimodal agents. Results demonstrate the advantages of *VideoAgent*: on averaged increasing 6.6% on NExT-QA and 26.0% on EgoSchema over baselines. Our further investigation has examined the role played by the unified memory and tool selection.

To summarize, our contributions are as follows:

- We propose a unified memory mechanism to build structured representations for long-form videos, including a *temporal memory* that stores segment-level descriptions and an *object memory* that tracks the state of objects appearing in the video.
- Based on the unified memory, we design *VideoAgent*, an LLM-powered multimodal agent for video understanding. It decomposes the input task queries and interactively invokes tools to retrieve information from the memory until reaches the final response.
- We perform thorough evaluations of *VideoAgent* on multiple video understanding benchmarks against both end-to-end video-language models and multimodal agent baselines, demonstrating the effectiveness of *VideoAgent*. The additional qualitative and ablation analysis further confirms the crucial design choices we've made.

2 VideoAgent

2.1 Overview

We illustrate the proposed *VideoAgent* in Appendix D.5. It begins with converting the input video into a unified representation: *temporal memory* (Sec. 2.2) and *object memory* (Sec. 2.3). For any incoming task, it interactively invokes tools to collect information from the memory and the raw video segments, and ultimately produces a response (Sec. 2.4). The memory construction and task-solving (inference) procedures are summarized in Algorithm 1 and Algorithm 2, respectively.

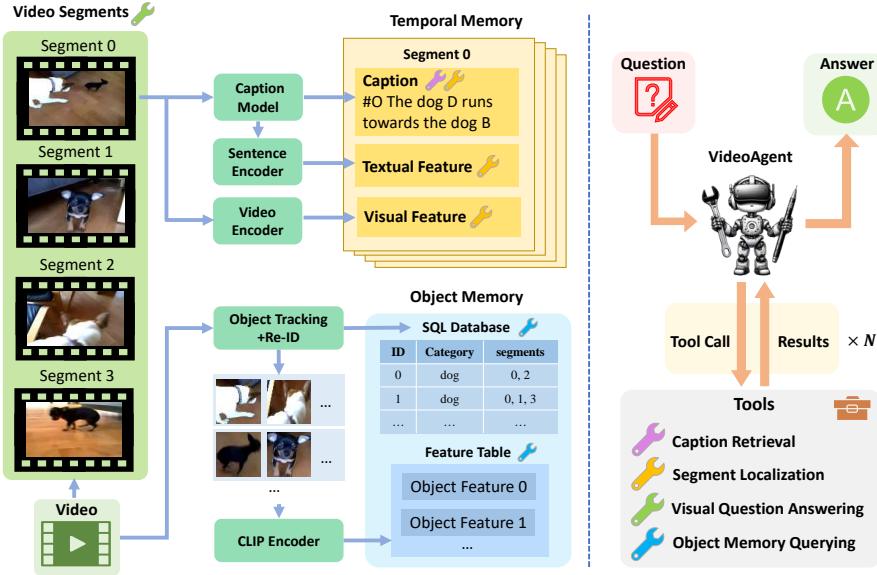


Fig. 2: An overview of *VideoAgent*. Left: We first translate an input video into structured representations: a temporal memory and an object memory; Right: the LLM within *VideoAgent* will be prompted to solve the given task by interactively invoking tools (☞). Our considered tools primarily work with the memory (*e.g.* ☞ interacts with the caption part of the temporal memory while ↪ looks up the object memory).

2.2 Temporal Memory \mathcal{M}_T

The temporal memory is designed to store overall event context descriptions and features of videos. Given n video segments $[v_1, \dots, v_n]$ sliced from a video V , we extract video segment caption s_{caption} , video segment feature e_{video} and the caption text embedding e_{caption} :

Video segment caption. We use a pretrained video captioning model called LaViLa [45] to produce captions for each video segment. Specifically, it takes 4 frames from a 2-second segment to produce a short caption sentence. Typical LaViLa captions can be "#C C cuts a wood with a wood cutter" and "#O The man Y pushes a stroller on the road with his left hand", where "#C" and "#O" is used to denote whether the caption sentence is about the camera wearer or someone other than the camera wearer, therefore making LaViLa captions effective in both egocentric and generic videos.

Video segment feature and caption feature. To obtain the video segment feature, we adopt the video encoder of ViCLIP [28] to encode video segments. We uniformly sampled 10 frames from each video segment as the input to ViCLIP, and save the generated feature of the segment. For the caption feature, we choose



Fig. 3: A visualization of object tracking and re-ID. 6 frames from a video are displayed in order. The cup (light green box) and the milk bottle (pink box) are successfully re-identified in different postures.

`text-embedding-3-large`⁴ offered by OpenAI to compute the embedding of the video segment caption we obtained from LaViLa.

2.3 Object Memory \mathcal{M}_O

In addition to the general video event context stored in the temporal memory, it is also crucial to explicitly capture the temporally consistent details: *e.g.* the presence of people, objects, and the surroundings, *etc.* The intuition is that most queries about videos are object(person)-related; therefore, the occurrence of objects (and people) are tracked and stored in the *object memory*. Specifically, object memory constitutes a feature table that connects object visual features with unique object identifiers, and a SQL database that stores the object(person) occurrence information across the video. Details on the construction can be found below:

Tracking and re-identification. At the heart of our object memory construction pipeline is tracking all the objects across the video, and re-identifying (re-ID) previously appeared objects to eliminate object duplication. We pipeline an object detection model RT-DETR [13] with a multi-object tracker ByteTrack [43] for the object discovery and tracking part. This combination produces tracking IDs, categories, and bounding boxes of the tracked objects in each video frame. For the re-ID part, the key idea is to first compute the features of all the objects that have been discovered and tracked, then group them based on their feature similarities. This gives us a unique object ID for each object. More specifically, object features are generated on object images cropped from 10 randomly sampled frames where the object appears; we also follow a recent study [27] to use an ensemble of

⁴ <https://platform.openai.com/docs/guides/embeddings>

CLIP [19] and DINOv2 [18] feature similarity:

$$\text{CLIP}(i, j) = \frac{1}{1 + \exp[-20 * (\text{cosine}(e_i^{\text{CLIP}}, e_j^{\text{CLIP}}) - 0.925)]},$$

$$\text{DINOv2}(i, j) = \frac{1}{1 + \exp[-4.1 * (\text{cosine}(e_i^{\text{DINOv2}}, e_j^{\text{DINOv2}}) - 0.5)]},$$

$$\text{sim}(i, j) = 0.15 * \text{CLIP}(i, j) + 0.85 * \text{DINOv2}(i, j),$$

where $\text{cosine}(\cdot, \cdot)$ denotes cosine similarity, $e_i^{\text{CLIP}}, e_j^{\text{CLIP}}$ and $e_i^{\text{DINOv2}}, e_j^{\text{DINOv2}}$ are the CLIP and DINOv2 features of the object i and j , respectively. The hyperparameters above (coefficients and biases) are tuned with a simple grid search. We provide more details on the re-ID algorithm in *Appendix*. An example of how our tracking and re-ID pipeline manages to handle the temporally discontinuous object presence in a kitchen can be found in Fig. 3.

Feature table. Assuming we’ve identified m objects from a video V and their presences in the n segments $V = [v_1, \dots, v_n]$ have been confirmed as well. We simply store both its CLIP feature used by the re-ID algorithm and the unique object ID s_{id} obtained after tracking and re-ID in a table. This allows us to use free-form language queries (*e.g.* “red cup”) to search for objects in the video.

SQL database. Further, we build a relational database with three fields: object ID s_{id} , object category s_{category} , and indices of video segments $\{I_1, \dots, I_t\}$ where the object has appeared. Later, this database can be queried using SQL code and support sophisticated querying logic.

2.4 Tools and Inference

Compared to counterparts that offer a large collection of tools and usually result in ambiguity in tool calling and complex tool-use pipeline, in *VideoAgent*, our design principle is to provide a minimal but sufficient tool set with a focus on querying the memory. We find this simplifies the inference procedures as well as leads to better performances. We consider the following tools (☞☛☛☛):

☞ **Caption retrieval.** The goal is to extract the captions from specified video segments. Concretely, given the temporal memory \mathcal{M}_T , a start and an end time step t_{start} , and t_{end} as arguments, the tool `caption_retrieval()` simply retrieves these captions from the temporal memory directly. Due to the context limit, the longest time window allowed is 15 segments, *i.e.* $t_{\text{end}} < t_{\text{start}} + 15$.

☛ **Segment localization.** The goal is to localize a video segment given a text query s_{query} . The tool `segment_localization()` will compare the text feature of s_{query} against the video features in the temporal memory \mathcal{M}_T . Specifically, we consider an ensemble of the query-video similarity (made possible by ViCLIP [28], a pretrained video-text CLIP model) and the query-caption similarity (both text features are computed by `text-embedding-3-large` offered by OpenAI). Ensemble details can be found in *Appendix*. Top-5 video segments will be returned by this tool.

Algorithm 1: Memory construction of *VideoAgent*.

Input: video V , video captioning model $\text{video_cap}(\cdot)$, video embedding model $\text{video_emb}(\cdot)$, text embedding model $\text{text_emb}(\cdot)$, video object tracker with re-identification $\text{object_track_reid}(\cdot)$

Output: temporal memory \mathcal{M}_T , object memory \mathcal{M}_O

- 1 Initialize $\mathcal{M}_T = \emptyset$, $\mathcal{M}_O = \emptyset$;
- 2 Slicing video into n short segments $V = [v_1, v_2, \dots, v_n]$ (each segment spans approximately 2 seconds);
- 3 **for** v_i in $[v_1, v_2, \dots, v_n]$ **do**
- 4 $s_{\text{caption}} \leftarrow \text{video_cap}(v_i)$;
- 5 $e_{\text{video}} \leftarrow \text{video_emb}(v_i)$;
- 6 $e_{\text{text}} \leftarrow \text{text_emb}(s_{\text{caption}})$;
- 7 $\mathcal{M}_T = \mathcal{M}_T + (s_{\text{caption}}, e_{\text{video}}, e_{\text{text}})$
- 8 results $\leftarrow \text{object_track_reid}(V)$;
- 9 **for** S in results **do**
- 10 $s_{\text{id}}, s_{\text{category}}, \{I_1, \dots, I_k\} \leftarrow S$ // object ID, object category, all k video segments where object o_i appears (I indicates video segment index);
- 11 $e_{\text{object}}^{(i)} \leftarrow \frac{1}{k} \sum_{j=1}^k e_{o_j}^{\text{CLIP}}$ // See Sec. 2.3;
- 12 $\mathcal{M}_O = \mathcal{M}_O + (s_{\text{id}}, s_{\text{category}}, e_{\text{object}}^{(i)}, \{I_1, \dots, I_t\})$;
- 13 **return** $\mathcal{M}_T, \mathcal{M}_O$;

👉 **Visual question answering.** The goal is to answer a given question s_{question} about a short video segment at time t_{target} , allowing to gather extra information that is not covered by the captions in temporal memory or states in object memory. Concretely, we run Video-LLaVA [11] when the tool `visual_question_answering()` is called.

🔗 **Object memory querying.** The goal is to perform sophisticated information retrieval about objects that appeared in the video from the object memory \mathcal{M}_O . Specifically, when calling the tool `object_memory_querying()` with a text query s_{query} (e.g. “How many red cups did I take out from the fridge?”), relevant object descriptions will first be extracted from the query (e.g. “red cup”); next, we compare the text feature of the descriptions (obtained from CLIP [19]) against the object features from the feature table in \mathcal{M}_O to obtain the object IDs likely correspond to the descriptions; finally, the LLM will write SQL code based on both s_{query} and the retrieved object IDs to query the database in \mathcal{M}_O and obtain the needed information (segments that the objects appeared, etc.). After being further processed by the LLM, a response to s_{query} will be returned.

The inference procedure of *VideoAgent* is rather straightforward. Starting with a history buffer h initialized with the input query q , *VideoAgent* decides which tool to use, calls the tool with the produced arguments, appends the results to the buffer, and repeats until it decides to stop or a maximum number of

Algorithm 2: Inference of *VideoAgent*.

Input: task instruction q , temporal memory \mathcal{M}_T , object memory \mathcal{M}_O ,
LLM $\text{LLM}(\cdot)$, a set of tools (see Sec. 2.4)

Output: response a

```

1 Initialize history  $h = [q]$ ;
2 Initialize inference step count  $c = 0$ ;
3 while  $c < \text{MAX\_STEP}$  do
4   action, input =  $\text{LLM}(h)$ ;
5   if  $\text{action} == \text{"caption\_retrieval"}$  then
6      $t_{\text{start}}, t_{\text{end}} \leftarrow \text{input}$ ;
7     results  $\leftarrow \textcolor{purple}{\mathcal{F}} \text{caption\_retrieval}(t_{\text{start}}, t_{\text{end}}, \mathcal{M}_T)$ ;
8   else if  $\text{action} == \text{"segment\_localization"}$  then
9      $s_{\text{query}} \leftarrow \text{input}$ ;
10    results  $\leftarrow \textcolor{orange}{\mathcal{F}} \text{segment\_localization}(s_{\text{query}}, \mathcal{M}_T)$ ;
11  else if  $\text{action} == \text{"visual\_question\_answering"}$  then
12     $s_{\text{question}}, t_{\text{target}} \leftarrow \text{input}$ ;
13    results  $\leftarrow \textcolor{green}{\mathcal{F}} \text{visual\_question\_answering}(s_{\text{question}}, t_{\text{target}})$ ;
14  else if  $\text{action} == \text{"object\_memory\_querying"}$  then
15     $s_{\text{query}} \leftarrow \text{input}$ ;
16    results  $\leftarrow \textcolor{blue}{\mathcal{F}} \text{object\_memory\_querying}(s_{\text{query}}, \mathcal{M}_O)$ ;
17  else if  $\text{action} == \text{"stop"}$  then
18     $\text{break}$ ;
19   $h = h + [(action, input, results)]$ ;
20   $c = c + 1$ ;
21 return  $a = \text{LLM}(h)$ ;

```

steps is reached. Finally, a response will be made based on the content in the history buffer. We provide an example of this procedure in Fig. 4. *VideoAgent* is implemented using LangChain⁵ with GPT-4 as the main LLM.

3 Capabilities and Analysis

We evaluate *VideoAgent* on various long-form video understanding benchmarks including EgoSchema (Sec. 3.1), Ego4D Natural Language Queries (Sec. 3.2), and NExT-QA (Sec. 3.3), the performances are compared against state-of-the-art end-to-end video-language models and multimodal agents.

3.1 EgoSchema

Overview. EgoSchema [15] is a dataset focusing on complex questions on long-form videos. Given a 3-minute video and a question, the model needs to select the

⁵ <https://www.langchain.com/>

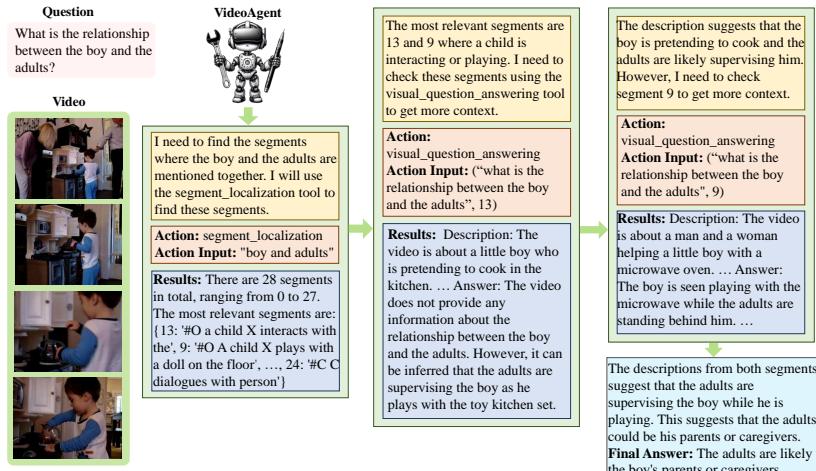


Fig. 4: An examples of the *VideoAgent* inference. Given a question, *VideoAgent* executes multiple tool-use steps until it reaches the answer. The yellow, red, and blue blocks in each step denote the chain of thought, action to be taken, and results of tool use.

Table 1: Results on the EgoSchema dataset. Top row: results on the full EgoSchema test set; Bottom row: results on the EgoSchema 500 subset due to the evaluation cost.

EgoSchema (full set)					
Method	FrozenBiLM	VIOLET	mPLUG-Owl	InternVideo	Gemini 1.5 Pro
Acc	26.9	19.9	30.2	32.0	63.2
EgoSchema (subset, 500 questions)					
Method	SeViLA	Video-LLaVA	mPLUG-Owl	ViperGPT	<i>VideoAgent</i>
Acc	25.8	36.8	33.8	15.8	62.8

true answer out of 5 possible options, and the evaluation metric is the accuracy *acc* on all the questions. The questions in EgoSchema typically involve video-level reasoning such as “describe the general activity in the room and how the different characters and their actions contribute to this environment”, posing great challenges to existing video understanding models. The full EgoSchema test set contains around 5000 questions, and the official subset of EgoSchema contains 500 questions. By considering cost, *VideoAgent* is tested on the 500-question subset, compared with state-of-the-art methods SeViLA [39], Video-LLaVA [11], mPLUG-Owl [38] and ViperGPT [23]. On the full set of EgoSchema, the performance of FrozenBiLM [36], VIOLET [2], mPLUG-Owl [38] and InternVideo [29] are reported by the benchmark [15], along with the performance of Gimini 1.5 Pro provided in their technical report⁶.

⁶ https://storage.googleapis.com/deepmind-media/gemini/gemini_v1_5_report.pdf

Main results. Tab. 1 shows the performances of various models on EgoSchema. *VideoAgent* significantly outperforms other state-of-the-art video understanding models such as SeViLA and Video-LLaVA to nearly 30 percent, achieving an accuracy of 62.8 on the 500 questions, closing to the performance of Gemini 1.5 Pro on the full set. The strong performance of *VideoAgent* on Egoschema proves that our proposed memory-based multimodal tool-use agent can solve complex video tasks on long-form videos better than end-to-end VLM and agent counterparts. Our insights into this can be found in the below two points.

Unified memory facilitates stronger reasoning. The questions in EgoSchema are rather complex in terms of the underlying reasoning about the lengthy videos. Therefore, strong spatial-temporal reasoning is essential. What canonical approaches like video-language models (Video-LLaVA, etc.) or counterpart multimodal agents (ViperGPT) have in common is the lack of a unified memory as a structured representation for the videos. Without such representation, the reasoning has to be either implicit (as in video-language models) or quite limited by the available tools (as in ViperGPT), results in worse performances than ours.

Holistic video understanding with flexible tool-use. Given a typical question such as "how did c's behavior evolve throughout the video, and what stages of engagement with the tasks can you identify?", it is hard to derive a descriptive text from the question and use it for video grounding, which is a common way for video-language models to select limited keyframes (4-32 frames for most models) for the visual input. However, apart from the `segment_localization`, *VideoAgent* can also use `caption_retrieval` to grab the main context of the video and decide which segments are critical, therefore tackling this obstacle.

3.2 Ego4D Natural Language Queries

Overview. The task of Ego4D Natual Language Queries [5] is to locate a temporal window (9 seconds on average) in the video (9 minutes on average) that can best answer a query. The model is allowed to retrieve k candidate windows, and the Intersection Over Union (IoU) between each candidate window and the ground-truth window is calculated. The metric $Rk@t$ is used to evaluate the proportion of the predictions where at least one out of k candidates achieves an IoU greater than t . Among the supervised methods, 2D-TAN [42] and VSLNet [41] are two baselines provided by the benchmark; GroundNLQ [8] ranked first in the Ego4D NLQ challenge 2023. *VideoAgent* is evaluated zero-shot with different variants of the `segment_localization` tool using 1) ViCLIP visual features only; 2) textual features based on LaViLa captions or Ego4D ground-truth narrations; 3) a combination of both textual features and visual features (LaViLa+ViCLIP, Ego4D+ViCLIP). The retrieved segments are padded with 3 seconds on both ends for computing the IoU with the ground-truth window. For LaViLa+ViCLIP and Ego4D+ViCLIP, the ensemble weights of video-text and text-text similarities are 18 : 11 and 7 : 8 respectively found by grid search on the training set.

Main results. Tab. 2 presents the results on the validation set of Ego4D NLQ. It can be inferred that a combination of both textual features and visual features

Table 2: Comparison between supervised baselines and *VideoAgent* with different tool implementation variants on Ego4D NLQ validation set.

Method	EGO4D NLQ Val.			
	<i>R1@0.3</i>	<i>R1@0.5</i>	<i>R5@0.3</i>	<i>R5@0.5</i>
	Supervised			
2D-TAN	5.04	2.02	12.89	5.88
VSLNet	5.45	3.12	10.74	6.63
GroundNLQ	27.20	18.91	54.42	39.98
Zero-Shot (<i>VideoAgent</i> with  segment_localization variants)				
ViCLIP	8.40	3.97	17.36	8.50
LaViLa	10.07	4.19	22.53	10.58
Ego4D	16.41	6.96	31.96	15.01
LaViLa+ViCLIP	11.13	4.76	25.31	12.08
Ego4D+ViCLIP	17.39	7.47	33.05	15.73

Table 3: Comparison between two zero-shot approach: *VideoAgent* and LifeLongMemory [30] on Ego4D NLQ. *LifelongMemory only reports results with a customized metric (*R@0.3*), but its numbers on standard metrics *R1@0.3* and *R5@0.3*, although not reported, must be less or equal than those on *R@0.3*.

Method	<i>R1@0.3</i>	<i>R5@0.3</i>	<i>R@0.3</i>
LifeLongMemory(Ego4D)	*	*	15.99
LifeLongMemory(LaViLa)	*	*	9.74
<i>VideoAgent</i> (Ego4D)	16.41	31.96	-
<i>VideoAgent</i> (LaViLa)	10.07	22.53	-

results in better video grounding. Although having a performance gap with the supervised GroundNLQ, *VideoAgent* outperforms 2D-TAN and VSLNet and achieves good performance considering its simple architecture and zero-shot characteristics.

Caption features vs. visual features. From the comparison among ViCLIP, LaViLa and Ego4D in Tab. 2, it can be inferred that it is more effective to use the caption–query similarities for video grounding than using video–query similarities. Higher quality captions (LaViLa→Ego4D) will also lead to better caption-based video grounding performance.

Similarity-based vs. LLM-based localization. Tab. 3 presents a comparison between *VideoAgent* and LifeLongMemory [30]. Given a query, LifeLongMemory uses a LLM (GPT-4) to digest and refine the captions of the video segments, and outputs a list of candidate windows to the query based on the captions selected by the LLM, *i.e.* segment localization is completely done by the LLM. LifeLongMemory adopts a customized *R@0.3* metric (instead of the standard *R1@0.3* and *R5@0.3*) to calculate the proportion of the predictions where at least one out of all the LLM-generated candidates achieves an *IoU* greater than 0.3. It can be inferred from Tab. 3 that given the same caption type (Ego4D or LaViLa),

Table 4: Results on NExT-QA. We compare baselines on both the original full set as reference and the subset (600 questions) due to the evaluation cost.

Method	NExT-QA			
	Temporal	Causal	Descriptive	Average
Val. Set				
InternVideo	43.4	48.0	65.1	49.1
SeViLA(zero-shot)	61.3	61.5	75.6	63.6
TCR(pre-training)	-	-	-	66.1
Val. Subset (600)				
ViperGPT	17.2	19.0	26.7	21.0
mPLUG-Owl	36.0	41.0	52.5	43.2
Video-LLaVA	42.0	53.5	65.0	53.5
SeViLA(zero-shot)	56.0	66.5	70.0	64.2
<i>VideoAgent</i>	60.0	76.0	76.5	70.8

the performance of *VideoAgent* on $R1@0.3$ where only 1 candidate is allowed for a query, has already surpassed the performance of LifeLongMemory on $R@0.3$. By providing 5 candidates for a query, the performance of *VideoAgent* will exceed LifeLongMemory by more than two-fold. This indicates that similarity-based segment localization is more effective than the LLM-based method.

3.3 NExT-QA

Overview. NExT-QA [35] is a video question answering benchmark containing temporal, causal and descriptive questions. For each question, the model should choose the correct option out of 5 options. The accuracy acc is computed for each type of the questions. The validation set contains around 5000 questions in total. For the reason of cost, we randomly sampled 200 questions for each type and obtained a subset of 600 questions in total to test the performance of *VideoAgent*. Methods directly compared with *VideoAgent* on this subset include ViperGPT [23], mPLUG-Owl [38], Video-LLaVA [11] and SeViLA [39]. The results of three representative methods InternVideo [29], SeViLA [39] and TCR [9] on the full validation set are also provided.

Main results. Tab. 4 shows the main results on NExT-QA. In all, *VideoAgent* achieves the strongest performances among all comparative methods. Particularly, on the challenging causal questions that require strong temporal understanding and reasoning, *VideoAgent* outperforms SeViLA, one of the state-of-the-art models on NExT-QA, for nearly 10 percent. Besides, the comparison between *VideoAgent* and Video-LLaVA, which is used by the  video_question_answering tool, indicates that our *VideoAgent* allows such video-language model to work better as part of the multimodal tool-use agent than being used alone.

Settings for ablation studies. We extract 50 questions for each question type from the 600-question subset, resulting in a subset of 150 questions in total, to evaluate the contributions of different components in *VideoAgent* as ablation

Table 5: The effectiveness of different components of *VideoAgent* on NExT-QA subset. ✓ and ✗ indicates whether or not the tool is included. "w/ re-ID" uses an object memory constructed with re-ID, while "w/o re-ID" uses an object memory that might include duplicated objects.

Type	VQA	Grounding	Captions	Database	Tem.	Cau.	Des.	Avg.
1	GPT-4V	✓	✓	w/ re-ID	64.0	78.0	82.0	74.7
2	Video-LLaVA	✓	✓	w/ re-ID	60.0	74.0	80.0	71.3
3	Video-LLaVA	✓	✓	✗	46.0	64.0	78.0	62.7
4	✗	✗	✓	w/ re-ID	48.0	52.0	68.0	56.0
5	✗	✗	✓	w/o re-ID	46.0	46.0	54.0	48.7
6	✗	✗	✓	✗	34.0	46.0	42.0	40.7

studies. Tab. 5 shows the performances of 6 ablations of *VideoAgent*, with each equipped with a unique set of tools among `visual_question_answering`, `segment_localization`, `caption_retrieval` and `object_memory_querying`, denoted as ‘VQA’, ‘Grounding’, ‘Captions’ and ‘Database’ in Tab. 5.

The necessity of caption retrieval. The `caption_retrieval` tool lays the foundation for *VideoAgent* since it provides the basic information about the main context of the video. With `caption_retrieval` only, *VideoAgent* of type 6 achieves an average result of 40.7 already, which is comparable to the performance 43.2 of the end-to-end video-language model mPLUG-Owl on the 600-question subset.

Object memory boosts all question types. The comparison between type 2 and 3 indicates that a reliable object memory can substantially help with temporal and causal questions since it offers crucial temporally consistent object information across video segments, facilitates object-related temporal localization, and enhances the agent’s understanding of the video. The performance gap between type 4 and type 5 suggests that with the object re-ID algorithm, the performance on descriptive questions (mostly about quantity) will be significantly improved, validating the effectiveness of object re-ID.

VQA and segment localization offer the most bonus. By comparing type 3 and 6, it can be seen that simultaneously adding `visual_question_answering` and `segment_localization` boost the caption-only *VideoAgent* by 22 percent on the average performance, compared to 15.3 percent boost by adding the object memory (inferred from type 4 and 6). Moreover, by switching from Video-LLaVA to GPT-4V in `visual_question_answering` (type 1 and 2), the performance will be raised by 3.4 percent, indicating that accurate visual details identified by the powerful VQA model will aid in better question answering performance.

4 Related Work

4.1 Multimodal LLMs for video understanding

Since LLMs have demonstrated an excellent ability to process and understand natural language [4, 17], several recent works have explored extending them

to multimodal setting, especially for images and videos [1, 10, 11, 21, 26, 44]. LaViLa [45] manages to create a massive and diverse set of text as automatic video narrators for video-text contrastive representation pretraining. Video-LLaMA [40] enables video comprehension by capturing the temporal changes in visual scenes and integrating audio-visual signals for better cross-modal training. As we discussed in Sec. 1, many of these multimodal foundation models could struggle with long-form video understanding. To remedy this, LSTP [31] utilize spatial and temporal sampler modules to extract optical flow based temporal features and aligned spatial relations from the video to achieve long-form video understanding; Gemini [26] scales the multimodal models to longer videos with tens of thousands of TPUs and massive private video-text datasets. Albeit the prompt progress made by these end-to-end models, prohibitive computation costs and the inherent limitation of the transformer on long-form videos remain significant in applying these end-to-end learned multimodal foundation models to video understanding.

4.2 Multimodal tool-use agents for video understanding

Another line of research focuses on augmenting LLMs with a set of **tools** to solve multimodal tasks without costly training. In particular, LLMs within these **multimodal agents** are prompted to produce a step-by-step plan to address the original task, and interactively invoke several multimodal foundation models (“tools”), *e.g.* captioning, VQA, *etc.* VisProg [6] pilots this direction by equipping the GPT-3 planner with a large collection of visual tools, solving complex real-world visual reasoning problems. Applying these agents to video understanding requires careful design as many of the tool models do not guarantee generalization to videos. LifeLongMemory [30] employs natural language video narrations to create a text-based episodic memory and prompt LLMs to reason and retrieve required information for the downstream task. DoraemonGPT [37] introduces a sophisticated prompting strategy with Monte Carlo Tree Search (MCTS) to invoke both tools and a structured memory to solve video understanding tasks. These multimodal agents have great potential but so far they mostly struggle with attaining on-par performances to their end-to-end foundation model counterparts on common benchmarks, likely due to the complicated pipelines and lack of video-specific design.

5 Conclusions

We’ve presented *VideoAgent*, a multimodal tool-use agent that reconciles several foundation models (large language models and vision-language models) with a novel unified memory mechanism for video understanding. Compared to end-to-end video-language model and tool-use agent counterparts, *VideoAgent* adopts a minimalist tool-use pipeline and does not require expensive training, while offering comparable or better empirical results on challenging long-form video understanding benchmarks including EgoSchema, Ego4D NLQ, and NExT-QA.

Possible future direction includes more exploration of real-world applications in robotics, manufacturing, and augmented reality.

References

1. Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)* (2022)
2. Fu, T.J., Li, L., Gan, Z., Lin, K., Wang, W.Y., Wang, L., Liu, Z.: 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*. pp. 22898–22909 (2023)
3. Gao, Z., Du, Y., Zhang, X., Ma, X., Han, W., Zhu, S.C., Li, Q.: Clova: A closed-loop visual assistant with tool usage and update. *arXiv preprint arXiv:2312.10908* (2023)
4. Gong, R., Huang, Q., Ma, X., Vo, H., Durante, Z., Noda, Y., Zheng, Z., Zhu, S.C., Terzopoulos, D., Fei-Fei, L., et al.: Mindagent: Emergent gaming interaction. *arXiv preprint arXiv:2309.09971* (2023)
5. Grauman, K., Westbury, A., Byrne, E., Chavis, Z., Furnari, A., Girdhar, R., Hamburger, J., Jiang, H., Liu, M., Liu, X., 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*. pp. 18995–19012 (2022)
6. Gupta, T., Kembhavi, A.: Visual programming: Compositional visual reasoning without training. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 14953–14962 (2023)
7. Han, T., Xie, W., Zisserman, A.: Temporal alignment networks for long-term video. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2906–2916 (2022)
8. Hou, Z., Ji, L., Gao, D., Zhong, W., Yan, K., Li, C., Chan, W.K., Ngo, C.W., Duan, N., Shou, M.Z.: Groundnlq@ ego4d natural language queries challenge 2023. *arXiv preprint arXiv:2306.15255* (2023)
9. Korbar, B., Xian, Y., Tonioni, A., Zisserman, A., Tombari, F.: Text-conditioned resampler for long form video understanding. *arXiv preprint arXiv:2312.11897* (2023)
10. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597* (2023)
11. Lin, B., Zhu, B., Ye, Y., Ning, M., Jin, P., Yuan, L.: Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122* (2023)
12. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. *arXiv preprint arXiv:2304.08485* (2023)
13. Lv, W., Xu, S., Zhao, Y., Wang, G., Wei, J., Cui, C., Du, Y., Dang, Q., Liu, Y.: Detrs beat yolos on real-time object detection. *arXiv preprint arXiv:2304.08069* (2023)
14. Maaz, M., Rasheed, H., Khan, S., Khan, F.S.: Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424* (2023)

15. Mangalam, K., Akshulakov, R., Malik, J.: Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems* **36** (2024)
16. Miech, A., Zhukov, D., Alayrac, J.B., Tapaswi, M., Laptev, I., Sivic, J.: Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 2630–2640 (2019)
17. OpenAI: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
18. Oquab, M., Dariseti, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., et al.: Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193 (2023)
19. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: *International Conference on Machine Learning (ICML)* (2021)
20. Schick, T., Dwivedi-Yu, J., Dessi, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* **36** (2024)
21. Shafiuallah, N.M.M., Paxton, C., Pinto, L., Chintala, S., Szlam, A.: Clip-fields: Weakly supervised semantic fields for robotic memory. arXiv preprint arXiv:2210.05663 (2022)
22. Song, E., Chai, W., Wang, G., Zhang, Y., Zhou, H., Wu, F., Guo, X., Ye, T., Lu, Y., Hwang, J.N., et al.: Moviechat: From dense token to sparse memory for long video understanding. arXiv preprint arXiv:2307.16449 (2023)
23. Surís, D., Menon, S., Vondrick, C.: Vipergpt: Visual inference via python execution for reasoning. arXiv preprint arXiv:2303.08128 (2023)
24. Tapaswi, M., Zhu, Y., Stiefelhagen, R., Torralba, A., Urtasun, R., Fidler, S.: Movieqa: Understanding stories in movies through question-answering. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4631–4640 (2016)
25. Tay, Y., Dehghani, M., Abnar, S., Shen, Y., Bahri, D., Pham, P., Rao, J., Yang, L., Ruder, S., Metzler, D.: Long range arena: A benchmark for efficient transformers. arXiv preprint arXiv:2011.04006 (2020)
26. Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., et al.: Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 (2023)
27. Tong, S., Liu, Z., Zhai, Y., Ma, Y., LeCun, Y., Xie, S.: Eyes wide shut? exploring the visual shortcomings of multimodal llms. arXiv preprint arXiv:2401.06209 (2024)
28. Wang, Y., He, Y., Li, Y., Li, K., Yu, J., Ma, X., Li, X., Chen, G., Chen, X., Wang, Y., et al.: Internvid: A large-scale video-text dataset for multimodal understanding and generation. arXiv preprint arXiv:2307.06942 (2023)
29. Wang, Y., Li, K., Li, Y., He, Y., Huang, B., Zhao, Z., Zhang, H., Xu, J., Liu, Y., Wang, Z., et al.: Internvideo: General video foundation models via generative and discriminative learning. arXiv preprint arXiv:2212.03191 (2022)
30. Wang, Y., Yang, Y., Ren, M.: Lifelongmemory: Leveraging llms for answering queries in egocentric videos. arXiv preprint arXiv:2312.05269 (2023)
31. Wang, Y., Wang, Y., Wu, P., Liang, J., Zhao, D., Zheng, Z.: Lstp: Language-guided spatial-temporal prompt learning for long-form video-text understanding. arXiv preprint arXiv:2402.16050 (2024)

32. Wiles, O., Carreira, J., Barr, I., Zisserman, A., Malinowski, M.: Compressed vision for efficient video understanding. In: Proceedings of the Asian Conference on Computer Vision. pp. 4581–4597 (2022)
33. Wu, C.Y., Krahenbuhl, P.: Towards long-form video understanding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1884–1894 (2021)
34. Wu, C., Yin, S., Qi, W., Wang, X., Tang, Z., Duan, N.: Visual chatgpt: Talking, drawing and editing with visual foundation models. arXiv preprint arXiv:2303.04671 (2023)
35. Xiao, J., Shang, X., Yao, A., Chua, T.S.: Next-qa: Next phase of question-answering to explaining temporal actions. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9777–9786 (2021)
36. Yang, A., Miech, A., Sivic, J., Laptev, I., Schmid, C.: Zero-shot video question answering via frozen bidirectional language models. Advances in Neural Information Processing Systems **35**, 124–141 (2022)
37. Yang, Z., Chen, G., Li, X., Wang, W., Yang, Y.: Doraemongpt: Toward understanding dynamic scenes with large language models. arXiv preprint arXiv:2401.08392 (2024)
38. Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al.: mplug-owl: Modularization empowers large language models with multimodality. arXiv preprint arXiv:2304.14178 (2023)
39. Yu, S., Cho, J., Yadav, P., Bansal, M.: Self-chained image-language model for video localization and question answering. Advances in Neural Information Processing Systems **36** (2024)
40. Zhang, H., Li, X., Bing, L.: Video-llama: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858 (2023)
41. Zhang, H., Sun, A., Jing, W., Zhou, J.T.: Span-based localizing network for natural language video localization. arXiv preprint arXiv:2004.13931 (2020)
42. Zhang, S., Peng, H., Fu, J., Luo, J.: Learning 2d temporal adjacent networks for moment localization with natural language. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 12870–12877 (2020)
43. Zhang, Y., Sun, P., Jiang, Y., Yu, D., Weng, F., Yuan, Z., Luo, P., Liu, W., Wang, X.: Bytetrack: Multi-object tracking by associating every detection box. In: European Conference on Computer Vision. pp. 1–21. Springer (2022)
44. Zhao, H., Cai, Z., Si, S., Ma, X., An, K., Chen, L., Liu, Z., Wang, S., Han, W., Chang, B.: Mmicl: Empowering vision-language model with multi-modal in-context learning. arXiv preprint arXiv:2309.07915 (2023)
45. Zhao, Y., Misra, I., Krähenbühl, P., Girdhar, R.: Learning video representations from large language models. In: CVPR (2023)

In this Appendix, we will first detail the implementation of object re-ID method. Then, the tasks included in *VideoAgent* and the corresponding models will be listed, followed by the experiment settings of *VideoAgent* and all the comparative methods. Finally, cases of the inference of *VideoAgent* will be illustrated.

A Object Re-Identification

Based on the tracking results, object re-identification (re-ID) aims at merging the occurrence of the same object in different period (diverse tracking IDs). The following algorithm shows the procedure of object re-ID. It receives a set of tracking IDs, and output a set of Re-ID groups, where each Re-ID group contains several tracking IDs that belong to the same object, representing a unique object ID in the database.

Algorithm 3: Object Re-Identification by Grouping.

```

Input: video  $V$ , tracking IDs  $\{t_1, t_2, \dots, t_n\}$ 
Output: a list of RE-ID groups  $G = \{U_1, U_2, \dots, U_m\}$ 
1 Initialize tracking IDs  $T = \{t_1, t_2, \dots, t_n\}$  to be examined;
2 Initialize the set of re-ID groups  $G = \{\}$ ;
3 for frame  $f$  in  $V$  do
4   for  $t_i$  appears in  $f$  and  $t_i \in T$  do
5     for Re-ID group  $U$  in  $G$  do
6       if  $\forall t_j \in U$ , share-no-frame( $t_i, t_j$ ) and  $\forall t_j \in U$ ,  $\text{sim}(t_i, t_j) > 0.5$ 
          and  $\exists t_j \in U$ ,  $\text{sim}(t_i, t_j) > 0.62$  then
7         remove  $t_i$  from  $T$ ;
8         add  $t_i$  to  $U$ ;
9         break;
10      if  $t_i \in T$  then
11        remove  $t_i$  from  $T$ ;
12        create a new group  $U = \{t_i\}$ ;
13        add  $U$  to  $G$ ;
14 output  $G$ ;

```

For each video frame, the algorithm checks every tracking ID in the frame that has not been examined and try to assign it to any existing Re-ID group. A tracking ID t_i should satisfy three conditions in order to be merged to a Re-ID group U : 1) It should not co-exist with any tracking IDs in U , since the same object only has one bounding box in each frame; 2) It should has $\text{sim}(t_i, t_j) > 0.5$ for all t_j in U ; where sim refer to the CLIP and DINOv2 feature similarity in the paper; 3) At least one tracking ID t_j in group U satisfies $\text{sim}(o_i, o_j) > 0.62$. If the

tracking ID t_i cannot be merged to any existing Re-ID group, then the algorithm will spare a new re-ID group initialized with t_i . The results of object re-ID are used to construct the SQL database, with each re-ID group corresponding to a unique object ID in the database.

B Tasks and Models

Tab. 6 shows the different tasks in *VideoAgent* and their corresponding models. For each task, the granularity level of the details is also shown. For instance, in the task of segment captioning, the details of captions usually include the actions of the characters and the primary objects in the video that the characters are interacting with.

Table 6: The methods and the granularity-level of the extracted information in different tasks.

Task	Method	Detail Granularity
Memory		
Segment Captioning	LaViLa	action, primary object
Object Tracking	RT-DETR+ByteTrack	object category
Object Re-ID	CLIP+DINOv2	object feature
Tools		
Video Grounding	ViCLIP+Text-Embedding-3-Large	action, primary object
Visual Question Answering	Video-LLaVA	action, object

C Experiment Settings

C.1 Settings of *VideoAgent*

Prompt of *VideoAgent* The tool-use capabilities of the LLM (GPT-4) is facilitated using LangChain⁷. The LLM is prompted by the following text for the video question answering task.

You are tasked with answering a multiple-choice question related to a video. The question has 5 choices, labeled as 0, 1, 2, 3, 4. The video is segmented into 2-second segments, each with an integer ID starting from zero and incrementing in chronological order. Each segment has a caption depicting the event. There is an object memory that records the appearing objects in each segment. The object memory is maintained by another agent. You have access to the following tools:

{tools}

ATTENTION:

1. the segment captions with prefix '#C' refer to the camera wearer, while captions with prefix '#O' refer to someone other than the camera wearer.
2. You can use both 'visual_question_answering' and 'object_memory_querying' to answer questions related to objects or people.
3. The 'visual_question_answering' may have hallucination. You should pay more attention to the description rather than the answer in 'visual_question_answering'.
4. The input to the tools should not contain the name of any other tool as well as the token '.
5. Its easier to answer the multiple-choice question by validating the choices.
6. If the information is too vague to provide an accurate answer, make your best guess.

Use the following format:

Question: the input question you must answer

Thought: you should always think about what to do

Action: the action to take, should be one of [{tool_names}]

Action Input: the input to the action

Observation: the result of the action... (this Thought/Action/Action Input/Observation can repeat N times)

Thought: I now know the final answer

⁷ <https://www.langchain.com/>

Final Answer: the correct choice label (0, 1, 2, 3, 4) to the original input question

Begin!

Question: {**input**} Thought: {**agent_scratchpad**}

In the above prompt format, **tools** refer to a set of tool names and their functional description, including:

caption_retrieval: Given an input tuple (start_segment, end_segment), get all the captions between the two segment IDs, 15 captions at most. end_segment < start_segment + 15.

segment_localization: Given a single string description, this tool returns the total number of segments and the top-5 candidate segments with the highest caption-description similarities.

visual_question_answering: Given an input tuple (question, segment_id), this tool will focus on the video segment starting from segment_id-1 to segment_id+1. It returns the description of the video segment and the answer of the question based on the segment.

object_memory_querying: Given an object-related question such as ‘what objects are in the video?’, ‘how many people are there in the video?’, this tool will give the answer based on the object memory. This tool is not totally accurate.

input refers to the multiple-choice question input, including a question and 5 options. **agent_scratchpad** is a list maintained by LangChain that stores the intermediate steps of the agent.

Object Memory Querying The **object_memory_querying** tool is achieved by another LLM agent (GPT-4) specialized in SQL writing, equipped with the following tools:

- **database_querying(program)**: return the results from the object memory database by executing the SQL *program*.
- **open_vocabulary_retrieval(description)**: return the possible object IDs that satisfy the object *description*.

Given an object-related query raised by the central agent, the memory agent will get the relevant object IDs by open vocabulary retrieval, translate the query into SQL program, fetch the results from the database in object memory by running the SQL program, and return the natural language answer to the central agent.

Experiment Settings of *VideoAgent* For NExT-QA and EgoSchema, we use the above prompt for testing the performance of *VideoAgent*. For Ego4D NLQ, the ensemble proportion of video-text and text-text similarities for LaViLa+ViCLIP is 18:11, and that for Ego4D+ViCLIP is 7:8. The ensemble proportions is found by grid search on the training set of Ego4D NLQ according to the maximal overall performance on $R1@0.3$, $R1@0.5$, $R5@0.3$ and $R5@0.5$.

C.2 Settings of Comparative Methods

In the experiments, we test the performance of the following methods by our own. The experiment settings for different comparative methods are detailed as follows.

- SeViLA: The default settings provided in their code are adopted for evaluation. The video frame number is set to 32, and the key frame number is set to 4.
- Video-LLaVA: The default settings provided in their code are adopted for evaluation. The input frame number is set to 8.
- mPLUG-Owl: We follow the evaluation procedure in EgoSchema dataset paper [15], which prompts mPLUG-Owl by ‘Given question <question text>, is answer <answer text> correct?’ along with the video frames. The option with the highest softmax score of the token ‘Yes’ in the output text will be viewed as the answer of mPLUG-Owl. The input frame number is set to 5 according to the best mPLUG-Owl settings provided in EgoSchema dataset paper.
- ViperGPT: GPT-3.5 is adopted as the code generator. 4 frames are uniformly sampled from the video and the generated code is run on the 4 frames to gather information for answering the question.

D Case Study

In this section, the successful cases of *VideoAgent* on both NExT-QA and EgoSchema are illustrated. In each step towards the final answer, the LLM first reasons about the action to take, and then outputs the action (tool) and its input. The tool will be executed and return the textual results to the LLM. This procedure will iterate until the LLM gets the final answer.

D.1 Case 1

In this case, the LLM uses the tool `segment_localization`, `caption_retrieval` and `visual_question_answering` to answer the question. Due to the man in red only occupies a tiny area of the screen, Video-LLaVA does not find the man and produce an answer with hallucination to the question "what does the man do next". However, the LLM will synthesize all the information and produce a correct answer. The video can be found here⁸.



Question: "what does the man in red do after looking at a distance from the plane at the start"
 0: "shakes his head"
 1: "happy"
 2: "point forward"
 3: "count down"
 4: "moves away"

Entering new AgentExecutor chain...

I need to find the segment where the man in red is looking at a distance from the plane. Then I can use the `visual_question_answering` tool to find out what he does next.

Action: `segment_localization`

Action Input: "man in red looking at a distance from the plane"

Results: There are 44 segments in total, ranging from 0 to 43. 22: '#C C stares the drone on the ground', 15: '#C C looks around the area', 41: '#C C looks around the area', 39: '#O A man x adjusts a drone on the', 40: '#C C stares the drone on the'

⁸ https://youtu.be/5tCWCMCWJKw?si=-IKmlA20_2SqaI_W

The segment with the highest similarity to the description is segment 22, but it doesn't mention a plane. The next highest is segment 15, but it also doesn't mention a plane. The segments 39 and 40 mention a drone, which could be the "plane" referred to in the question. I'll check the captions around these segments to get more context.

Action: caption_retrieval

Action Input: (37, 42)

Results: {37: '#C C looks around the parking', 38: '#C C looks around the', 39: '#O A man x adjusts a drone on the', 40: '#C C stares the drone on the', 41: '#C C looks around the area', 42: '#C C looks around the'}

The captions from segments 37 to 42 show a man adjusting a drone and then staring at it. This could be the "man in red looking at a distance from the plane" referred to in the question. I'll use the visual_question_answering tool to find out what he does next.

Action: visual_question_answering

Action Input: ("what does the man do next?", 40)

Results: Description: The video is about a large airplane that is parked on the runway. It shows the plane from different angles, including a close-up of the front of the plane. Answer: The man in the video walks towards the plane and opens the door.

The man in the video walks towards the plane and opens the door. This action is not listed in the choices. The closest option is "moves away", which could be interpreted as the man moving away from his original position towards the plane.

Final Answer: 4

Finished chain.

D.2 Case 2

In this case, `segment_localization` and `visual_question_answering` are used to answer the questions. Although some details provided by Video-LLaVA are incorrect, the LLM is able to choose the correct answer based on the captions of the video and the results from Video-LLaVA. The video can be found here⁹.



Question: "how does the boy keep the bat stable on the ground"

- 0: "lean against his body"
- 1: "sit on cushion"
- 2: "put on stone"
- 3: "swing chair is protected"
- 4: "wear shoes"

Entering new AgentExecutor chain...

The question is asking about a specific action performed by a boy in the video. I need to find the segment where the boy is interacting with the bat and then use the `visual_question_answering` tool to understand how he keeps the bat stable. First, I will use the `segment_localization` tool to find the segments where the boy and the bat are mentioned together.

Action: `segment_localization`

Action Input: "boy bat"

Results: There are 12 segments in total, ranging from 0 to 11. {7: '#O A man Y holds the baseball bat with both', 6: '#O The man B hits the ball with the rake in his right hand', 9: '#O man Y holds a head warmer', 10: '#O A man Z wears a helmet on his head with both hands', 2: '#O A

⁹ <https://youtu.be/vz1PCFqdtQQ?si=X9vATb1ClBVM8oMM>

man X walks around the field'}

The segment _localization tool has returned several segments, but only segment 7 seems to be relevant to the question as it mentions a man (possibly the boy) holding a baseball bat. I will use the visual_question_answering tool on this segment to get more details about how the boy keeps the bat stable.

Action: visual_question_answering

Action Input: ("how does the boy keep the bat stable?", 7)

Results: Description: The video is about a man who is practicing his archery skills in a park. He is seen holding a bow and arrow and shooting at a target. Answer: The boy keeps the bat stable by holding it firmly in his hand and keeping it steady. He also uses his body movements to keep the bat stable.

The visual_question_answering tool provided an answer that the boy keeps the bat stable by holding it firmly in his hand and using his body movements. This corresponds to the first option in the multiple-choice question.

Final Answer: 0

Finished chain.

D.3 Case 3

In this case, the `visual_question_answering` tool contributes most to the successful question answering, since it identifies the wedding ring on the woman's left hand, although it mistakes the woman's action in the video (delivering a talk to others) as 'playing rock paper scissors with a man'. The video can be found here¹⁰.



Question: "how do you know that the woman with black spectacles at the start is married"

- 0: "wearing a wedding gown"
- 1: "holding flowers"
- 2: "wearing a wedding ring"
- 3: "grandparent and grandchild"
- 4: "walking with the groom"

Entering new AgentExecutor chain...

I need to find the segment where the woman with black spectacles appears at the start of the video. Then I can use the `visual_question_answering` tool to determine if she is married based on the given choices.

Action: `segment_localization`

Action Input: "woman with black spectacles at the start"

Results: There are 44 segments in total, ranging from 0 to 43. {30: '#O the woman A touches the glasses on her face with her left hand.', 11: '#O The woman X touches her glasses with her right hand', 14: '#O person X adjusts her glasses ', 16: '#O woman X looks at the television', 39: '#C C stares at a woman'}

¹⁰ <https://youtu.be/EFMdUrOELqw?si=6iXTKJCW0x96a7Ap>

The woman with black spectacles appears in the first segment (0). Now I will use the visual_question_answering tool to determine if she is married based on the given choices.

Action: visual_question_answering

Action Input: ("how do you know that the woman is married", 0)

Results: Description: The video shows a man and a woman playing a game of rock paper scissors. The woman wins, and the man gets angry. Answer: The video does not provide any information about the woman's marital status. However, it is possible to observe that she is wearing a wedding ring on her left hand.

The woman is wearing a wedding ring which indicates that she is married. This matches with choice 2.

Final Answer: 2

Finished chain.

D.4 Case 4

In this case, the object memory is queried by the LLM (achieved by a memory agent mentioned in Appendix C.1) to answer the quantity of the elephants, and output the correct answer. The video can be found here¹¹.



Question: "how many elephants are there"

- 0: "one"
- 1: "four"
- 2: "three"
- 3: "six"
- 4: "two"

Entering new AgentExecutor chain..

I should use the 'object_memory_querying' tool to find out how many elephants are in the video.

Action: object_memory_querying

Action Input: 'how many elephants are there in the video?'

Results: There are 2 elephants in the video.

I now know the final answer.

Final Answer: 4

Finished chain.

¹¹ <https://youtu.be/0B87EvVi-as?si=DST7FJfGtPmW6ZcB>

D.5 Case 5

This is a typical question in EgoSchema, which requires strong textual understanding and reasoning abilities of the model to correctly answer it. By retrieving the captions, the LLM can gather the information of the ingredients in various segments, and uses `visual_question_answering`¹² tool on the corresponding segment to investigate the function of each ingredient. The video can be found here¹².



Question: 'Analyze the choice of ingredients and their importance in the video. how do they contribute to the final dish?'

- 0: 'The tissue paper is an important ingredient in the dish. it is used to clean up any spills.'
- 1: "The can cover, a crucial component in the dish, serves to maintain the minced beef's freshness effectively by protecting it."
- 2: 'The wax paper, being a crucial ingredient in the dish, is utilized effectively to securely wrap the can to prevent any potential leakage.'
- 3: 'The minced beef, spring onions, and mayonnaise are all important ingredients in the dish. the minced beef provides protein, the spring onions provide vitamins and minerals, and the mayonnaise provides flavor.'
- 4: "The transparent glass bowl serves as an important crucial ingredient in the dish's preparation. it is efficiently used to thoroughly mix the various ingredients together well."

Entering new AgentExecutor chain...

¹² <https://youtu.be/s0LU2K7Ukwo?si=RNEIfZjGCmWeM-1t>



The question is asking about the importance of ingredients in the final dish. I need to find out what ingredients are used in the video and how they contribute to the final dish. I will start by using the caption_retrieval tool to get an overview of the events in the video.

Action: caption_retrieval

Action Input: (0, 14)

Results: {0: '#C C scoops sauce from the container with the spoon in her right hand', 1: '#C C pours the tomato paste from the spoon into the bowl of macaroni and cheese with her right hand', 2: '#C C scoops sauce from the sauce container with the spoon in her right hand', 3: '#C C pours the spice into the plate of scrambled eggs with the spoon in her right hand', 4: '#C C applies the ketchup from the spoon in his right hand on the salad in the bowl', 5: '#C C scoops sauce from the sauce container with the spoon in her right hand', 6: '#C C pours the tomato sauce into the bowl on the dining table with the spoon in her right hand', 7: '#C C pours milk from the bottle in her left hand into the cup of milk in her right hand', 8: '#C C puts some sauce from the sauce container on the spoon with her right hand.', 9: '#C C holds the spoon in the sauce container with her left hand.', 10: '#C C drops the container on the table with her left hand', 11: '#C C closes the trash', 12: '#C C walks towards the dining table', 13: '#C C drops the tissue paper on the table with her left hand', 14: '#C C picks a tissue paper from the floor with her left'}

The captions indicate that the dish involves sauce, tomato paste, macaroni and cheese, spice, scrambled eggs, ketchup, salad, tomato sauce, milk, and tissue paper. The tissue paper seems to be used for cleaning rather than as an ingredient. I need to find out more about how these ingredients are used. I will use the visual_question_answering tool to get more details about the use of these ingredients.

Action: visual_question_answering

Action Input: ('What is the role of the sauce in the dish?', 0)

Results: Description: The video is about a person preparing a meal in a kitchen. They are seen cutting up food and mixing it in a bowl. Answer: The sauce is poured over the noodles and mixed in with the other ingredients, adding flavor and texture to the dish.

The sauce is used to add flavor and texture to the dish. Now, I will use the visual_question_answering tool again to find out more about the role of the other ingredients.

Action: visual_question_answering

Action Input: ('What is the role of the spice and scrambled eggs in the dish?', 3)

Results: Description: The video is about a person preparing and eating a meal in a spacious and well-lit room. Answer: The spice and scrambled

eggs are added to the dish to enhance the flavor and texture. The spice adds a burst of flavor and aroma, while the scrambled eggs provide a creamy and tender texture that complements the other ingredients.

The spice and scrambled eggs are used to enhance the flavor and texture of the dish. Now, I will use the visual_question_answering tool again to find out more about the role of the remaining ingredients.

Action: visual_question_answering

Action Input: ('What is the role of the tomato sauce and milk in the dish?', 6)

Results: Description: The video is about a person preparing and eating a meal in a spacious and well-lit kitchen. The person is seen cutting up food on a table, mixing it in a bowl, and then eating the food. Answer: The tomato sauce and milk are poured over the ice cream, adding a sweet and tangy flavor to the dish.

The tomato sauce and milk are used to add a sweet and tangy flavor to the dish. However, the answer seems to be hallucinating about ice cream, which was not mentioned in the captions. I will use the visual_question_answering tool again to find out more about the role of the remaining ingredients.

Action: visual_question_answering

Action Input: ('What is the role of the tomato paste and macaroni and cheese in the dish?', 1)

Results: Description: The video is about a person preparing and eating a meal in a spacious and well-lit room. Answer: The tomato paste and macaroni and cheese are added to the dish to enhance its flavor and texture. The tomato paste adds a rich and savory taste, while the macaroni and cheese provides a creamy and comforting element.

The tomato paste and macaroni and cheese are used to enhance the flavor and texture of the dish. The tomato paste adds a rich and savory taste, while the macaroni and cheese provides a creamy and comforting element. Now, I will use the visual_question_answering tool again to find out more about the role of the tissue paper.

Action: visual_question_answering

Action Input: ('What is the role of the tissue paper in the dish?', 13)

Results: Description: The video is about a person cleaning a bathroom. Answer: The tissue paper is used to clean the dish.

The tissue paper is used for cleaning, not as an ingredient in the dish. Now, I can compare the information I have gathered with the choices provided in the question.



Final Answer: 3

Finished chain.