

VTimeLLM: Empower LLM to Grasp Video Moments

Bin Huang, Xin Wang*, Hong Chen, Zihan Song, Wenwu Zhu*
Tsinghua University

{huangb23, h-chen20, songzh23}@mails.tsinghua.edu.cn

{xin_wang, wwzhu}@tsinghua.edu.cn

Abstract

Large language models (LLMs) have shown remarkable text understanding capabilities, which have been extended as Video LLMs to handle video data for comprehending visual details. However, existing Video LLMs can only provide a coarse description of the entire video, failing to capture the precise start and end time boundary of specific events. In this paper, we solve this issue via proposing VTimeLLM, a novel Video LLM designed for fine-grained video moment understanding and reasoning with respect to time boundary. Specifically, our VTimeLLM adopts a boundary-aware three-stage training strategy, which respectively utilizes image-text pairs for feature alignment, multiple-event videos to increase temporal-boundary awareness, and high-quality video-instruction tuning to further improve temporal understanding ability as well as align with human intents. Extensive experiments demonstrate that in fine-grained time-related comprehension tasks for videos such as Temporal Video Grounding and Dense Video Captioning, VTimeLLM significantly outperforms existing Video LLMs. Besides, benefits from the fine-grained temporal understanding of the videos further enable VTimeLLM to beat existing Video LLMs in video dialogue benchmark, showing its superior cross-modal understanding and reasoning abilities.¹

1. Introduction

Large language models (LLMs) have garnered significant attention due to their exceptional capabilities in text understanding and generation [23, 31]. However, harnessing the potential of LLMs for understanding and reasoning over multimodal data, especially videos, still remains a substantial challenge. This is because analyzing videos requires deep understanding of both visual details and temporal dynamics for models. Several preliminary at-

*Corresponding Authors.

¹Our project page is at <https://github.com/huangb23/VTimeLLM>

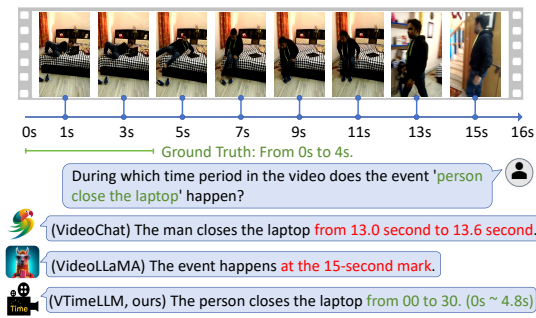


Figure 1. Existing Video LLMs, such as VideoChat and VideoLaMA, exhibit a deficiency in boundary awareness, leading to challenges in accurately capturing the precise timestamps of specific events.

tempts [15, 20, 34] utilize LLMs for video understanding. Nevertheless, these works predominantly focus on generation of generic video captions and can merely offer surface-level summaries of the content, thus failing to capture the relationships between specific moment boundaries and the bounded events, as shown in Figure 1.

To tackle the problem, in this paper, we investigate improving the boundary-aware ability of Video LLM, which faces the following two challenges.

1. There is a scarcity of large-scale video datasets with accurate boundary annotations to train the Video LLM for temporal alignment.
2. It is non-trivial to design effective temporal-related video tasks for training LLM to understand the content of multiple moments within videos.

To address these challenges, we propose VTimeLLM, a novel Video LLM that can perceive fine-grained segments in videos with better temporal reasoning ability. VTimeLLM consists of i) a visual encoder and a visual adapter to process the input video, and ii) a tailored LLM to understand both text and video content, which is trained via a novel boundary-aware three-stage training strategy. Specifically, visual features are aligned with LLM’s semantic space through image-text training in the first stage.

In the second stage, we then design the single-turn and multi-turn related question answering (QA) tasks to endue VTimeLLM with the awareness of time boundaries and the ability to understand the corresponding events bounded within the boundaries. We employ a large-scale video-text dataset containing multiple segments together with their roughly annotated labels for training VTimeLLM with the QA tasks. Finally, in the third stage, we further create a high-quality dialogue dataset for instruction tuning, which simultaneously aligns VTimeLLM with human intention and enables VTimeLLM to conduct temporal understanding for video segments more precisely. Extensive experiments show that VTimeLLM significantly outperforms existing Video LLMs in time-related video understanding tasks, such as Temporal Video Grounding and Dense Video Captioning. In addition, benefiting from the fine-grained temporal understanding of videos, VTimeLLM is able to beat existing Video LLMs in video dialogue benchmark, demonstrating its superiority in cross-modal understanding and reasoning for videos. Our contributions in this paper are listed as follows,

- We propose VTimeLLM, the first boundary-aware Video LLM, to the best of our knowledge.
- We propose the boundary-aware three-stage training strategy, which consecutively leverages i) large-scale image-text data for feature alignment, ii) large-scale multi-event video-text data together with the temporal-related single-turn and multi-turn QA to enhance the awareness of time boundary, and iii) instruction tuning on the high-quality dialog dataset for better temporal reasoning ability.
- We conduct extensive experiments to demonstrate that the proposed VTimeLLM significantly outperforms existing Video LLMs in various fine-grained temporal-related video tasks, showing its superior ability for video understanding and reasoning.

2. Related Works

2.1. Multimodal Large Language Model

Image LLMs To enable Large Language Models (LLMs) to comprehend visual information, significant efforts have been made to align visual and linguistic modalities. BLIP-2 [14] introduced the concept of Q-Former, utilizing learnable query vectors to extract visual features from frozen image encoders. MiniGPT-4 [37] demonstrated that further fine-tuning with detailed image descriptions significantly enhances its usability. LLaVA [16] explored diverse multi-modal instruction-following data, aiming to construct a general-purpose visual assistant. Recent endeavors, such as Kosmos-2 [21] and VisionLLM [26], delved into more detailed aspects of image comprehension, including referring and grounding, significantly enhancing the capability to describe intricate image details.

Video LLMs Driven by the success of Image LLM, researchers have naturally extended their focus from single-frame images to multi-frame videos, leading to the emergence of Video-compatible LLMs like VideoChat [15], Video-LLaMA [34], and Video-ChatGPT [20]. These models employ a two-stage training strategy. In the first stage, large-scale datasets align video features with the feature space of LLMs. In the second stage, a limited amount of GPT-annotated or human-annotated datasets are used for instruction tuning. While these models exhibit impressive overall video comprehension, their abilities to describe specific video segments and perform temporal reasoning remain limited. The limitation arises mainly due to the nature of datasets used in the first training stage, such as WebVid [2], which usually consist of one-event videos and noisy textual annotations. Moreover, the scarcity of high-quality, temporally annotated data in the second stage poses a challenge for models to conduct temporal reasoning. To bridge this gap, our approach, VTimeLLM, introduces a boundary perception stage between these two stages. This stage enables the model to precisely locate events within videos and describe multiple distinct events accurately, empowering our model to grasp fine-grained details of video moments.

2.2. Fine-Grained Video Understanding

Fine-grained video understanding, the ability to precisely locate and comprehend specific events within a video, is a crucial challenge for video analysis. When integrated with natural language, there are two primary tasks: Temporal Video Grounding [1, 8] and Dense Video Captioning [12, 25].

Temporal Video Grounding Temporal Video Grounding aims to identify corresponding video segments for given textual inputs. Traditional approaches can be categorized into two types: proposal-based [4, 29, 32] and proposal-free methods [9, 30, 33]. Proposal-based techniques generate candidate proposals before ranking them based on relevance. In contrast, proposal-free methods directly predict the start and end boundaries of the target moment.

Dense Video Captioning Dense Video Captioning is a more intricate task, demanding both temporal localization and captioning for all events within an untrimmed video. Earlier methods [6, 11, 12] employed a two-stage process involving temporal localization followed by event captioning. Recent developments [25, 28, 36] in this field have witnessed a shift towards joint training of captioning and localization modules. For instance, Vid2Seq [28], enhances a language model by incorporating specific time tokens, enabling the model to generate event boundaries and textual descriptions within the unified output sequence.

Both these two tasks share a fundamental requirement: the alignment of video segments with semantic context. Leveraging the power of LLM with the help of our training strategy, our VTimeLLM model unifies these tasks and has demonstrated remarkable effectiveness. Concurrently, VTimeLLM enables natural language interaction with humans, establishing itself as an excellent assistant for comprehending video content.

3. VTimeLLM: Being Aware of Time Boundaries in Videos

In this section, we introduce VTimeLLM, which is designed to grasp precise video moments for LLMs. We first provide a detailed description of the model architecture, and then our innovative boundary-aware three-stage training framework, as shown in Figure 2.

3.1. Architecture

To enable the LLM to comprehend videos, our VTimeLLM model incorporates two additional modules within LLM, i.e., the visual encoder and the visual adapter, which transform the visual information into text space.

Visual Encoder Our VTimeLLM model utilizes a frozen CLIP ViT-L/14 [22] as the visual encoder, and for simplicity we denote it as ViT. Given a video $V \in \mathbb{R}^{T \times H \times W \times C}$ with T frames, we uniformly sample $N = 100$ frames, represented as $\tilde{V} \in \mathbb{R}^{N \times H \times W \times C}$, where $\tilde{V}_1 = V_1$ and $\tilde{V}_N = V_T$. Each frame \tilde{V}_i is independently processed through the visual encoder:

$$\{v_i^{cls}, v_i^1, v_i^2, \dots, v_i^p\} = \text{ViT}(\tilde{V}_i), i = 1, 2, \dots, N, \quad (1)$$

where p represents the number of patches in the ViT.

Visual Adapter We utilize the global feature v_i^{cls} as the feature for the i -th frame, and apply a linear layer $f(\cdot)$ to project the features of each frame into the same embedding space as that of LLM:

$$z_i = f(v_i^{cls}), i = 1, 2, \dots, N. \quad (2)$$

Subsequently, $Z = \{z_i\} \in \mathbb{R}^{N \times d}$ becomes the input sequence that LLM can comprehend, and d is the hidden dimension of LLM.

Note that in the visual modules, we do not model the temporal relationships for the frames, inspired by the fact that the LLM itself can receive sequential input embeddings and capture their temporal relations.

Input of LLM To enable the simultaneous processing of video and text inputs, we introduce a special token, '<video>', to represent the video content. More intuitively,

if we want to ask the LLM questions about the video, we can use a mixed description, e.g., “*This is a video <video> can you describe this video ?*”. With this mixed description, the embedding layer will transform the textual words except for <video> into embedding as the LLM originally does, and we will obtain a textual embedding list $[w_1, w_2, \dots, w_M]$ for the text description, where $w_i \in \mathbb{R}^d$ is the embedding for each word and M is the word number. Then, the video feature sequence will be inserted into the embedding list, at the position of <video> as follows,

$$input = [w_1, \dots, w_{j-1}, Z, w_j, \dots, w_M], \quad (3)$$

where $j - 1$ and j are the indexes of the words that are close to the special token <video> in the original sentence. (In the previous example, w_{j-1} corresponds to the word “*video*”, w_j corresponds to the word “*can*”). Then LLM can further encode the input embedding list $input$ to understand the video and the user query.

Output for Temporal Boundaries We employ the text format ‘*from s to e*’ to denote a video moment, where s and e represent the starting and ending frame indexes of the moment, ranging from 00 to 99, with each number corresponding to a specific frame.

3.2. Boundary-aware Training

In contrast to the previous typical two-stage training approaches [15, 20, 34], consisting of alignment and instruction tuning, our approach introduces an additional stage to improve the temporal understanding ability of the model. Specifically, the first stage, feature alignment, aims to train the visual adapter, to align video features with LLM’s semantic space. The second stage, boundary perception, focuses on enabling LLM to develop attentional capabilities for specific moments, facilitating the understanding of various events occurring within the video. The third stage, instruction tuning, allows LLM to align with human intent and enabling more precise event localization and description. In the following sections, we will elaborate on the training methods and datasets utilized for each of these three stages.

3.2.1 Stage 1: Feature Alignment

Data Organization In the feature alignment stage, we employ the image-text LCS-558K dataset as curated by LLaVA [16]. This dataset is meticulously filtered to achieve a more balanced distribution of conceptual coverage. Comprising image-text pairs, we deliberately choose not to incorporate datasets containing video-text pairs with the following two considerations. Firstly, contemporary large-scale video-text datasets contain substantial textual noise,

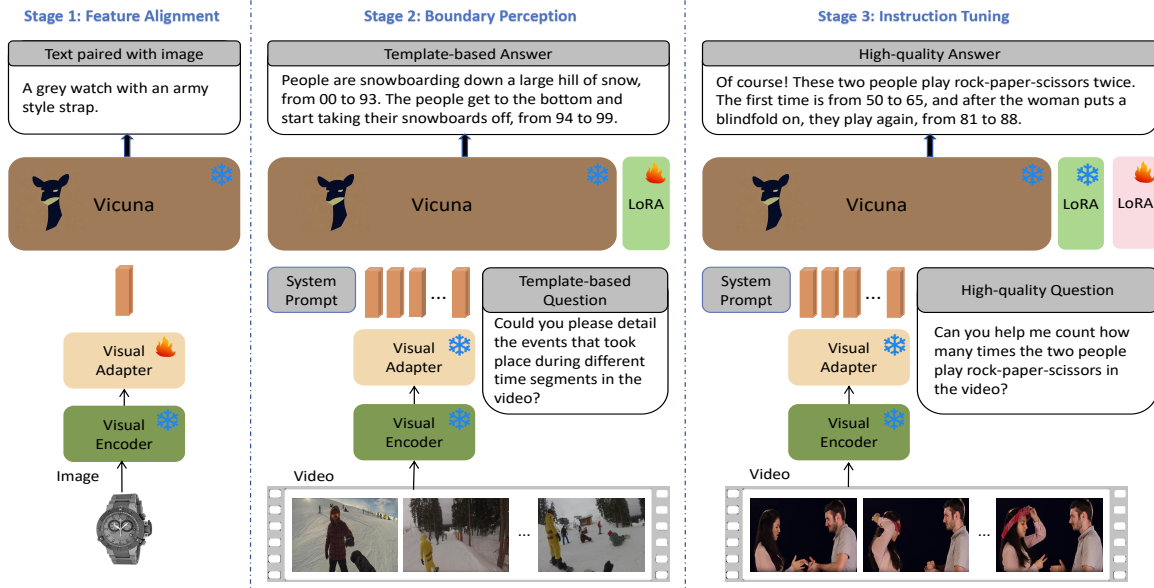


Figure 2. Our boundary-aware three-stage training framework. In the first stage, visual features are aligned with LLM’s semantic space through image-text training. In the second stage, we transform a large-scale multi-event dataset into a QA format based on templates, training VTimeLLM to possess temporal boundary awareness and understand events within the boundaries. In the third stage, we create a high-quality dialogue dataset for instruction tuning, which aligns VTimeLLM with human intent and also enables more precise video temporal understanding.

which severely impedes the alignment between visual features and textual semantics. Secondly, the transformation from visual information to text space usually suffers from information loss, e.g. when captioning an image or a video into “a dog is running on the grass”, we may lose information about the visual details (such as the color) of the dog. Comparatively, the loss of information resulting from summarizing an image into a few words is less than that of videos. Our experiments also demonstrate the superiority of using image datasets for alignment over video dataset (a filtered subset [19] of WebVid2M [2]), and even a combination of both.

Training Strategy For each image-text pair $\langle I, T \rangle$ in the dataset, a special token $\langle \text{image} \rangle$ is directly appended before the text T , the embedding of this token is extracted with the visual encoder and the visual adapter as follows, denoted as Z_I :

$$Z_I = f(\text{ViT}(I)^{cls}), \quad (4)$$

and we can obtain the embedding sequence:

$$\text{input} = [Z_I, w_1, w_2, \dots, w_M]. \quad (5)$$

Then, we can use it to train the visual adapter f , with the original auto-regressive training objective of the LLM.

3.2.2 Stage 2: Boundary Perception

After the training in the first stage, the LLM model becomes proficient in understanding visual information. In the second stage, we enhance the model’s capabilities to comprehend sequential image frames, i.e., video, encompassing the semantic understanding of video segments while ensuring alignment with the corresponding boundaries.

Data Organization Due to the time-consuming nature of manually annotating timestamps and semantics for video segments, there is currently a lack of large-scale multi-event video-text datasets. Traditional methods align video segments with text transcripts generated by Automatic Speech Recognition (ASR). However, this approach faces challenges due to the lack of synchronicity and consistency between actions performed and spoken content, leading to weak correlations and inaccuracies in boundary annotations.

Recently, we identified the InternVid-10M-FLT [27] dataset, which offers a viable solution for our boundary-aware training. This dataset employs an entirely automated process to segment and annotate video clips, eliminating the need for manual intervention. Consequently, a single video may contain multiple event annotations. To ensure suitability for our study, we selected specific videos, each not exceeding 120 seconds in length. These videos encom-

Box 1: Examples of both single-turn and multi-turn QA for a video containing three events.

Single-turn QA.

Q_1 : Could you please detail the events that took place during different time segments in the video?

A_1 : T_1 , from s_1 to e_1 . T_2 , from s_2 to e_2 . T_3 , from s_3 to e_3 .

Multi-turn QA.

Q_1 : Can you describe what occurred from s_2 to e_2 in the video? A_1 : T_2 .

Q_2 : Tell me about the events from s_3 to e_3 . A_2 : T_3 .

Q_3 : During which frames in the video can we observe T_1 happening? A_3 : From s_1 to e_1 .

Box 2: The inputs to VTimeLLM in Stage 2 and Stage 3

A chat between a curious user and an artificial intelligence assistant. The assistant gives helpful, detailed, and polite answers to the user's questions.

USER: This is a video with 100 frames: <video>\n Q_1 ASSISTANT: A_1 </s>

USER: Q_2 ASSISTANT: A_2 </s>.....

pass multiple non-overlapping event annotations, each lasting more than 3 seconds, and the average duration of these events exceeds 8% of the video length. Thus, we curate a dataset comprising 134k videos, where each video contains multiple events and their rough temporal annotations and descriptions.

In each video, a series of events $\{s_i, e_i, T_i\}$ is contained, where s_i and e_i represent the start and end timestamps of a segment, ranging from 00 to 99. T_i corresponds to its textual description. To transform these events into dialogue data $\{Q_1, A_1, Q_2, A_2, \dots\}$ suitable for training LLM, we devise two types of QA dialogues: single-turn and multi-turn, constituting 20% and 80% respectively. In Box 1, we have provided examples of both single-turn and multi-turn QA dialogues for a video containing three events. Specifically, the task of single-turn QA is dense video captioning. Q_1 prompts a question requiring a comprehensive description of all events and their corresponding timestamps, while A_1 outputs the respective textual descriptions and timestamps in a specified format as shown in the upper box of Box 1. On the other hand, multi-turn QA involves segment captioning and temporal video grounding tasks, demanding the description generation given timestamps or timestamps generation given descriptions, as shown in the lower box of Box 1. In multi-turn QA, each event will be randomly queried for one of these two tasks, and the questions are not necessarily presented in the order of the events' occurrence. We design 10 templates for each task to transform events into QA dialogues, which can be found in the appendix.

Training Strategy We format these QA pairs according to the original LLM's format, keeping the initial system prompts intact. Moreover, we insert the statement "This is a video with 100 frames: <video>\n" before the first question. For illustration purposes, the input for VTimeLLM

is presented in Box 2. With the reformatted sequences, we continue to employ the auto-regressive training objective, where the loss is computed exclusively on the tokens within the answer of the QA dialogues (A_1, A_2, \dots). To enhance training efficiency, we utilize LoRA [10] for fine-tuning the LLM. During this stage, we keep the visual adapter f frozen. Consequently, the only trainable parameters are the newly applied LoRA modules.

3.2.3 Stage 3: Instruction Tuning

Following the training in the second stage, our VTimeLLM model demonstrates the ability to comprehend all events within the video and align them with the corresponding timestamps. Despite the diverse templates employed, the model's output still tends to overfit the answers, which behaves more like a multi-task pretrained model while losing chatting ability with the user, e.g., when we input "What color is the coat of the man" to the model, it may response "from 00 to 10". Additionally, the labels of the video-text data in the second stage are originally annotated in an automated way, which are not so accurate and noisy. To tackle the two problems, in the third stage, we incorporate high-quality dialogue data for instruction tuning, enabling the model to follow human instructions for more accurate video temporal comprehension and reasoning.

Data Organization In this stage, we select a subset from ActivityNet Captions [12] and DiDeMo [1] datasets, and transform it into a high-quality QA dialogue dataset with the assistance of Large Language Models. In contrast to InternVid which employs automated segmenting and labeling, these two datasets are entirely manually annotated, resulting in descriptions that are more detailed and temporal boundaries that are more accurate. Specifically, we carefully selected a subset of videos from the training set

Table 1. The results of existing Video LLMs in temporal video grounding and dense video captioning tasks.

Model	Temporal Grounding								Dense Captioning		
	ActivityNet				Charades-STA				ActivityNet		
	R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7	mIoU	SODA_c	CIDEr	METEOR
VideoChat-7B [15]	8.8	3.7	1.5	7.2	9.0	3.3	1.3	6.5	0.9	2.2	0.9
VideoLLaMA-7B [34]	6.9	2.1	0.8	6.5	10.4	3.8	0.9	7.1	1.9	5.8	1.9
VideoChatGPT-7B [20]	26.4	13.6	6.1	18.9	20.0	7.7	1.7	13.7	1.9	5.8	2.1
VTimeLLM-7B	44.0	27.8	14.3	30.4	51.0	27.5	11.4	31.2	5.8	27.6	6.8
VTimeLLM-13B	44.8	29.5	14.2	31.4	55.3	34.3	14.7	34.6	5.9	27.2	6.7

of ActivityNet Captions. These videos contained a minimum of three non-overlapping events, collectively covering over 90% of the video duration, amounting to approximately 4.2k videos. Similarly, a subset of videos is being selected for the DiDeMo dataset, each containing at least two non-overlapping events and covering 40% of the video duration. This process results in a total of about 4k videos for the DiDeMo subset. Subsequently, we also transform these videos, which contain a series of events $\{s_i, e_i, T_i\}$, into QA dialogues. However, results from the second stage of training indicate that template-based conversations lead to model overfitting. Therefore, we utilize LLM for this transformation. Specifically, we provide these events to the LLM, prompting it to assume the role of an AI visual assistant capable of analyzing the video and generate a dialogue about the video between itself and a user. The prompt can be found in the appendix. This approach results in QA dialogues that are grammatically correct, linguistically coherent, and may encompass a variety of tasks. We generate two distinct sets of dialogues for each video, yielding a final dataset comprising around 16k high-quality QA dialogues. Additionally, we observe that introducing a comparable number of other video instruction tuning datasets further enhances the descriptive capabilities of the model, with minimal impact on temporal understanding abilities. Therefore, we add an extra 20k QA pairs from the VideoInstruct100K [20] dataset. Overall, in this stage, a total of approximately 36k QA dialogues are used for training, which is significantly smaller than the dataset used in the second stage.

Training Strategy We merge the LoRA module trained in the second stage with the original model and introduce a new LoRA module, which serves as the only trainable parameters. All other training details remain consistent with those of the second stage.

4. Experiment

4.1. Experiment Setup

Tasks, Dataset, and Evaluation Metrics To assess the capability of VTimeLLM in comprehending various event

segments, we mainly conduct evaluations on two tasks: Temporal Video Grounding and Dense Video Caption.

For the Temporal Video Grounding task, we utilize datasets from ActivityNet Captions [1] and Charades-STA [8]. We calculate the Intersection over Union (IoU) between the time segments generated by the model and the corresponding ground truth time segments. We report mean IoU (mIoU) and recall@1, $\text{IoU} \geq m$ (R@m) metric, where m values are set at $\{0.3, 0.5, 0.7\}$.

In the case of Dense Video Captioning, we employ the ActivityNet Captions [1] dataset. The evaluation process encompasses two categories of metrics. Firstly, we employ SODA_c [7], a metric specifically tailored for dense video caption tasks, taking into account the video’s storyline. Secondly, we compute matched pairs between the generated events and the ground truth across IoU thresholds of $\{0.3, 0.5, 0.7, 0.9\}$, and calculate captioning metrics based on these matched pairs [28]. We report CIDEr [24] and METEOR [3] averages under different IoU thresholds to provide a comprehensive analysis.

Implementation Details In our study, we use Vicuna v1.5 [5] as the Large Language Model and train two versions: 7B and 13B. We use a total batch size of 128 throughout the training process. The AdamW [18] optimizer is applied with a cosine learning rate decay and a warm-up period. In the first training stage, the total epoch number is 1 with a learning rate of 1×10^{-3} , and the subsequent second and third stages we will train for 2 epochs each with a learning rate of 1×10^{-4} . The LoRA parameters are set to $r = 64$ and $\alpha = 128$. Thanks to the efficiency of LoRA, we can complete the training of the 7B model within 30 hours with 1 RTX-4090 GPU.

4.2. Main Results

We evaluate the capabilities of existing Video LLMs in temporal video grounding and dense video captioning tasks, as shown in Table 1. Detailed information about the evaluation process can be found in the appendix. VTimeLLM-7B outperforms these Video LLMs of the same size by a significant margin. Upon further scaling up the model to 13B parameters, we observe minor changes in performance on

Table 2. Ablation study of the three-stage training strategy. In Stage 1, “I” and “V” represent the utilization of image or video datasets, and “I+V” signifies the merging of both datasets. In Stage 2, we compare options for Freezing or Tuning the visual adapter. In Stage 3, “Reuse” indicates the reuse of the LoRA from Stage 2, while “Addition” signifies the addition of a new LoRA module. “X” represents the absence of training in this stage.

Row	Stage1	Stage2	Stage3	Temporal Grounding								Dense Captioning		
				ActivityNet				Charades-STA				ActivityNet		
				R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7	mIoU	SODA_c	CIDEr	METEOR
1	I	Freeze	X	36.1	21.4	10.5	25.5	48.4	23.7	11.5	32.2	4.6	17.0	5.4
2		Tune		35.6	21.4	10.2	24.8	48.0	24.4	12.1	31.6	4.6	17.0	5.5
3	V	Freeze	X	33.5	17.9	8.1	23.9	46.6	17.7	6.7	30.5	4.0	14.1	5.0
4		Tune		34.8	18.2	8.5	24.4	47.9	19.9	8.0	31.4	4.2	14.1	5.2
5	I+V	Freeze	X	31.3	16.5	6.8	22.2	47.5	21.9	9.2	30.6	4.1	13.4	5.0
6		Tune		33.5	17.4	6.9	23.2	47.1	21.4	8.6	30.6	4.0	13.6	4.8
7	X	Tune	X	42.2	22.7	11.5	29.8	40.0	4.9	0.0	27.5	3.7	10.2	5.0
8	I	X	Addition	31.0	18.1	7.7	22.9	37.5	21.8	6.1	22.8	4.2	16.0	5.0
9	I	Freeze	Reuse	39.3	26.6	13.0	28.1	49.7	29.8	13.3	30.9	5.2	23.2	6.1
10	I	Freeze	Addition	44.0	27.8	14.3	30.4	51.0	27.5	11.4	31.2	5.8	27.6	6.8

ActivityNet tasks, while the temporal grounding ability improves on Charades-STA. It is worth noting that our training dataset does not include Charades-STA training data, indicating that increasing the scale of our VTimeLLM model enhances its out-of-distribution generalization ability.

We provide several possible explanations to account for the poor performance of other models: firstly, both VideoChat and VideoLLaMA extract only N=8 frames as input, making it challenging for them to achieve a fine-grained understanding of the video content. Secondly, the commonly used LLM (Vicuna) lacks robust positional awareness in input sequences. For instance, when posed with the question “What is the position of the word ‘video’ in the phrase ‘a video clip’?”, it may erroneously respond, “The word ‘video’ appears at position 67.” Relying solely on a limited set of temporally annotated data for instruction tuning is insufficient to address this issue. Therefore, it is essential to integrate boundary-aware training to achieve precise video comprehension.

4.3. Ablation Study

In this section, we provide detailed ablations about our three-stage training strategy through experiments on the 7B model, as illustrated in Table 2. In the ablation, our most concerned questions and their results are provided in the following.

Q1: How to train a good visual adapter? In contrast to other Video LLMs, we utilize a pure image modality for the first stage and find it to be superior across all metrics than using a pure video modality (Rows 1, 2 vs Rows 3, 4). This effectiveness of using images to alignment could be attributed to the higher quality and reduced information loss in image datasets. Additionally, using pure images outperforms the fusion of two modal datasets (Rows 1, 2 vs

Rows 5, 6). This could be due to the significant disparity in tasks, where describing a single frame event and describing a sequence of 100 frames events pose distinct challenges for model fitting.

Another question arises during the following stage: should the previously pretrained visual adapter be tuned or frozen? Upon comparing Row 1 vs Row 2, Row 3 vs Row 4, and Row 5 vs Row 6, we observed minor difference in performance between the two approaches. To retain the comprehensive information acquired during the pretraining stage, we opt to freeze the parameters of the visual adapter in the latter two stages.

Q2: Should the LoRA from stage 2 be reused in stage 3?

Upon comparing Row 9 to Row 10, it is evident that in stage 3, merging the LoRA module from the second stage with the LLM parameters and additionally incorporating another LoRA module yields superior results. This approach ensures that the temporal understanding capabilities acquired during stage 2 are effectively preserved within the model.

Q3: Is every training stage necessary?

By comparing Rows 1-6 with Row 7, we observe a substantial disparity in the model’s performance on temporal grounding when training without stage 1. The scores are abnormally high on the ActivityNet dataset, while significantly low on the Charades-STA dataset. Upon careful analysis of the outputs under this setting, we find that the model has not effectively learned to localize events. Instead, it tends to predict a temporal segment spanning nearly the entire video (e.g., from 00 to 95). In such cases, if the ratio of ground truth duration to video length is denoted as x , the IoU with the model’s output is approximately x . The ActivityNet dataset contains a significant number of long samples, with 20% of queries having $x > 0.5$, leading to an inflated evaluation metric.

Table 3. The results of video dialogue on video-based generative performance benchmarking.

Evaluation Aspect	VideoLLaMA	LLaMA-Adapter	VideoChat	VideoChatGPT	BT-Adapter	VTimeLLM
Temporal Understanding	1.82	1.98	1.94	1.98	2.34	2.49
Correctness of Information	1.96	2.03	2.23	2.40	2.68	2.78
Detail Orientation	2.18	2.32	2.50	2.52	2.69	3.10
Contextual Understanding	2.16	2.30	2.53	2.62	3.27	3.40
Consistency	1.79	2.15	2.24	2.37	2.46	2.47
Mean	1.98	2.16	2.29	2.38	2.69	2.85

Conversely, in the Charades-STA dataset, x rarely exceeds 0.5, demanding more precise localization [13]. However, the model without stage 1 training fails to achieve it. Moreover, the model’s performance in dense captioning tasks is unsatisfactory, which also highlights the essential nature of the feature alignment stage.

The necessity of stage 2 can be demonstrated by comparing Row 8 with Rows 9, 10. Despite the higher quality of annotations in stage 3, the limited dataset size hinders the model’s ability to achieve a robust temporal understanding through stage 3 training alone. Models trained solely in stage 3 exhibit inferior performance across various tasks compared to those that have undergone preliminary training in stage 2.

After stage 3 training, the model exhibits comprehensive improvement in the tasks outlined in the table (Row 1 vs Row 10). Furthermore, it regains chatting ability, enabling it to respond to a wide range of questions posed by humans.

4.4. Video Dialogue Performance

Besides the ability for fine-grained video understanding tasks, we explore whether VTimeLLM can address a broader range of questions through dialogue. We employ the Video-ChatGPT [20] benchmark and conduct an evaluation of video-based generative performance. This benchmark covers many questions associated with five key aspects. GPT-3.5 assigns a score, not exceeding 5, to the model-predicted answer based on the question and the correct answer. We present the average scores in Table 3 and compare VTimeLLM with existing Video LLMs, including VideoLLaMA [34], LLaMA-Adapter [35], VideoChat [15], VideoChatGPT [20] and BT-Adapter [17].

Thanks to the fine-grained video comprehension capabilities, VTimeLLM achieves state-of-the-art results in all aspects. The most substantial improvement is observed in the aspect of detail orientation, where VTimeLLM achieves a noteworthy enhancement of +0.41 (15.2%). We attribute this progress to two primary factors. Firstly, the image-based training in stage 1 ensures comprehensive preservation of visual details in individual frames, facilitating a detailed understanding of spatial dimension. Secondly, the temporal-aware training employed in the second and third stages enables VTimeLLM to capture multiple events

within videos, enhancing its ability to depict details of temporal dimension.

To better illustrate the video dialogue performance of VTimeLLM, we present a qualitative example, as shown in Figure 3.

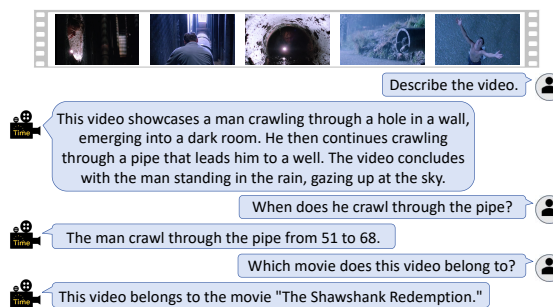


Figure 3. A qualitative example of video dialogue. The video is 160 seconds long.

5. Conclusion

In this work, we introduce VTimeLLM, a Video LLM capable of comprehending multiple events within a video and providing precise temporal boundaries. Specifically, we propose a three-stage temporal-aware training framework. This framework utilizes large-scale image-text data for feature alignment, leverages extensive multi-event video-text data along with temporal-related question-answering to enhance temporal awareness, and employs instruction tuning on a high-quality dialogue dataset to improve temporal reasoning ability. Extensive experiments demonstrate that VTimeLLM outperforms existing Video LLMs significantly across various tasks.

Acknowledgments

This work was supported by the National Key Research and Development Program of China No. 2020AAA0106300, National Natural Science Foundation of China (No. 62250008,62222209), Beijing National Research Center for Information Science and Technology (BNRist) under Grant No.BNR2023RC01003, BNR2023TD03006 and Beijing Key Lab of Networked Multimedia.

References

- [1] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*, pages 5803–5812, 2017. [2](#), [5](#), [6](#)
- [2] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1728–1738, 2021. [2](#), [4](#)
- [3] Satantjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005. [6](#)
- [4] Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. Temporally grounding natural sentence in video. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 162–171, 2018. [2](#)
- [5] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, 2023. [6](#)
- [6] Xuguang Duan, Wenbing Huang, Chuang Gan, Jingdong Wang, Wenwu Zhu, and Junzhou Huang. Weakly supervised dense event captioning in videos. *Advances in Neural Information Processing Systems*, 31, 2018. [2](#)
- [7] Soichiro Fujita, Tsutomu Hirao, Hidetaka Kamigaito, Manabu Okumura, and Masaaki Nagata. Soda: Story oriented dense video captioning evaluation framework. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pages 517–531. Springer, 2020. [6](#)
- [8] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, pages 5267–5275, 2017. [2](#), [6](#)
- [9] Soham Ghosh, Anuva Agarwal, Zarana Parekh, and Alexander Hauptmann. Excl: Extractive clip localization using natural language descriptions. *arXiv preprint arXiv:1904.02755*, 2019. [2](#)
- [10] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. [5](#)
- [11] Vladimir Iashin and Esa Rahtu. A better use of audio-visual cues: Dense video captioning with bi-modal transformer. *arXiv preprint arXiv:2005.08271*, 2020. [2](#)
- [12] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 706–715, 2017. [2](#), [5](#)
- [13] Xiaohan Lan, Yitian Yuan, Xin Wang, Long Chen, Zhi Wang, Lin Ma, and Wenwu Zhu. A closer look at debiased temporal sentence grounding in videos: Dataset, metric, and approach. *ACM Transactions on Multimedia Computing, Communications and Applications*, 19(6):1–23, 2023. [8](#)
- [14] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. [2](#)
- [15] KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. [1](#), [2](#), [3](#), [6](#), [8](#)
- [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. [2](#), [3](#)
- [17] Ruyang Liu, Chen Li, Yixiao Ge, Ying Shan, Thomas H Li, and Ge Li. One for all: Video conversation is feasible without video instruction tuning. *arXiv preprint arXiv:2309.15785*, 2023. [8](#)
- [18] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [6](#)
- [19] Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023. [4](#)
- [20] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023. [1](#), [2](#), [3](#), [6](#), [8](#)
- [21] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. [2](#)
- [22] 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. [3](#)
- [23] 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. [1](#)
- [24] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015. [6](#)
- [25] Jingwen Wang, Wenhao Jiang, Lin Ma, Wei Liu, and Yong Xu. Bidirectional attentive fusion with context gating for dense video captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7190–7198, 2018. [2](#)
- [26] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023. [2](#)
- [27] Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei

- Liu, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2307.06942*, 2023. 4
- [28] Antoine Yang, Arsha Nagrani, Paul Hongsuck Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10714–10726, 2023. 2, 6
- [29] Yitian Yuan, Lin Ma, Jingwen Wang, Wei Liu, and Wenwu Zhu. Semantic conditioned dynamic modulation for temporal sentence grounding in videos. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [30] Yitian Yuan, Tao Mei, and Wenwu Zhu. To find where you talk: Temporal sentence localization in video with attention based location regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9159–9166, 2019. 2
- [31] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022. 1
- [32] Da Zhang, Xiyang Dai, Xin Wang, Yuan-Fang Wang, and Larry S Davis. Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1247–1257, 2019. 2
- [33] Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. Span-based localizing network for natural language video localization. *arXiv preprint arXiv:2004.13931*, 2020. 2
- [34] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 1, 2, 3, 6, 8
- [35] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 8
- [36] Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. End-to-end dense video captioning with masked transformer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8739–8748, 2018. 2
- [37] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 2