# PPLLAVA: VARIED VIDEO SEQUENCE UNDERSTAND ING WITH PROMPT GUIDANCE

Anonymous authors

Paper under double-blind review

#### ABSTRACT

The past year has witnessed the significant advancement of video-based large language models. However, the challenge of developing a unified model for both short and long video understanding remains unresolved. Most existing video LLMs cannot handle hour-long videos, while methods custom for long videos tend to be ineffective for shorter videos and images. In this paper, we identify the key issue as the redundant content in videos. To address this, we propose a novel pooling strategy that simultaneously achieves token compression and instructionaware visual feature aggregation. Our model is termed Prompt-guided Pooling LLaVA, or PPLLaVA for short. Specifically, PPLLaVA consists of three core components: the CLIP-based visual-prompt alignment that extracts visual information relevant to the user's instructions, the prompt-guided pooling that compresses the visual sequence to arbitrary scales using convolution-style pooling, and the clip context extension designed for lengthy prompt common in visual dialogue. Moreover, our codebase also integrates the most advanced video Direct Preference Optimization (DPO) and visual interleave training. Extensive experiments have validated the performance of our model. With superior throughput, PPLLaVA achieves better results on image benchmarks as a video LLM, while achieving state-of-the-art performance across various video benchmarks, excelling in tasks ranging from caption generation to multiple-choice questions, and handling video lengths from seconds to hours. The codes are promised to be made public.

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

#### 1 INTRODUCTION

033 Video Large Language Models (Video LLMs) have made significant advancements over the past 034 year. Given the extensive resources and the scarcity of high-quality video-text data required for video pretraining, performing Image-to-Video transfer on powerful Image-domain Large Language Models has become a more practical approach for most Video LLMs. Building on the most advanced image LLMs (Liu et al., 2023a; 2024a; Dai et al., 2023), existing video LLMs typically ad-037 dress the modal differences between images and videos by video instruction data production (Maaz et al., 2023; Li et al., 2023b; Luo et al., 2023; Zhang et al., 2024b), temporal modeling (Liu et al., 2024c;d; Li et al., 2023c; Huang et al., 2023), or video token aggregation (Jin et al., 2023; Li et al., 040 2023d; Xu et al., 2024). Meanwhile, a wide range of video benchmarks and test tasks offer diverse 041 perspectives and options for evaluating the capabilities of video LLMs, including video question 042 answering (Maaz et al., 2023; Xu et al., 2016; Caba Heilbron et al., 2015; Wu et al., 2017), video 043 dense captioning (Ren et al., 2024), multiple-choice questions (Li et al., 2023c; Fu et al., 2024), and 044 long video assessment (Fu et al., 2024; Song et al., 2024; Zhang et al., 2024a).

For temporal modeling, an intuitive approach is to directly input tokens from each frame into the LLM, a method proven effective in several studies (Liu et al., 2024d;b; Li et al., 2024). However, while this method leverages the LLM's sequence modeling capabilities, it leads to an excessively long visual context. This not only increases computational resource consumption and processing time but also limits the model's ability to handle extended videos. To address this issue, several alternative approaches exist. A commonly adopted method is average pooling across the temporal dimension, frequently seen in early video LLMs (Li et al., 2023b; Maaz et al., 2023; Luo et al., 2023; Liu et al., 2024c). While this approach maintains a constant context length, it significantly diminishes the model's ability to capture temporal dynamics. Models designed specifically for long videos often incorporate unique structures, such as memory mechanisms (Ren et al., 2024; Zhang

069

071



Figure 1: (a) An instance from VideoMME (Fu et al., 2024). The crucial information pertains to only a small portion of the video for different questions. (b) Performance comparison of PPLLaVA with recent strong Video LLM among video benchmarks, image benchmarks, and efficiency. All the models are based on Vicuna-7B.

et al., 2024a; Zhou et al., 2024). Although these designs enable the models to handle hour-long
videos, they offer limited utility for short videos or images. Another approach is the use of conditional token pooling or aggregation (Li et al., 2023d; Xu et al., 2024; Jin et al., 2023). Unlike
global average pooling, this method reduces the context length while preserving some spatiotemporal structure, enabling more effective spatiotemporal modeling.

077 However, pooling inevitably leads to performance loss compared to using the full set. So, how can we reduce the number of tokens while preserving the spatiotemporal modeling capabilities? We believe the solution lies in the inherent characteristics of the video. As proven by many previous works 079 (Han et al., 2022; Liu et al., 2023b; Ma et al., 2022), videos contain significant redundancy, with key information often concentrated in just a few frames, which is particularly true for long videos. For 081 video LLMs, this issue can be more pronounced. As shown in Fig. 1(a), the user's instruction may pertain to only a small portion of the video, with the rest being redundant for correctly answering the 083 question. Therefore, if we can extract crucial video information while compressing tokens, we can 084 maintain or even enhance performance. In this context, Image LLMs have offered valuable inspira-085 tion. The BLIP series (Li et al., 2023a; Dai et al., 2023; Xue et al., 2024) and the LLaVA series (Liu et al., 2023a; 2024a;b; Li et al., 2024) are the two most popular structures in multimodal LLM. BLIP 087 uses a Q-Former for image-to-text mapping, while the LLaVA series employs simple linear projec-880 tion or MLP. Recently, LLaVA-based models have demonstrated that simple mapping can achieve better results with less training (Liu et al., 2024a). However, despite requiring more computation 089 resources and training stages, the Q-Former offers two key advantages: first, it significantly reduces 090 visual tokens by converting them into fewer query tokens; second, through the interaction between 091 text and visual tokens within the Q-Former, it enables more targeted extraction of video features 092 relevant to the user's instructions (Dai et al., 2023). Hence, can we develop a pooling method that retains LLaVA's simple structure and powerful weights while reducing the number of tokens and 094 enabling prompt-aware feature extraction? 095

To this end, we propose Prompt-guided Pooling LLaVA (PPLLaVA), a novel method that combines 096 visual pooling with instruction-aware visual feature extraction. Specifically, PPLLaVA first identifies prompt-relevant visual representations through fine-grained vision-prompt alignment. Then, 098 using the prompt-vision relevance as a 3D convolutional kernel, PPLLaVA can compress the visual tokens to any desired three-dimensional size based on the specified output size or stride. Fi-100 nally, recognizing that CLIP pretraining provides a limited context length and that training video 101 LLMs—particularly for multi-turn dialogues—requires long text contexts, PPLLaVA also employs 102 asymmetric positional embedding extensions to expand the text encoding capacity. As a result, 103 PPLLaVA effectively extracts relevant visual features from both long texts and short phrases while 104 compressing video tokens. PPLLaVA achieves over an 80% compression rate, supports ultra-long 105 video inputs, and simultaneously improves performance on short videos. In fact, PPLLaVA functions similarly to a Q-Former within LLaVA, but it offers several advantages over directly train-106 ing a Q-Former: (1) PPLLaVA introduces far fewer additional parameters and computational over-107 head, amounting to less than one-tenth of a Q-Former. (2) While a Q-Former requires a three-stage pretraining process—contrastive learning, alignment training, and instruction tuning—PPLLaVA
 can be utilized solely during instruction tuning, allowing for seamless transfer from image-domain
 LLMs. (3) PPLLaVA supports flexible output sizes for different modalities, whereas the number of
 queries in a Q-Former is fixed once set. As a result, different Q-Formers typically need to be trained
 separately for images and videos (Zhang et al., 2023; Li et al., 2023c).

113 Extensive experiments on the latest multimodal LLM benchmarks have validated the superiority of 114 PPLLaVA: with superior throughput, PPLLaVA has achieved top results across a wide range of test 115 sets, including MSRVTT (Xu et al., 2016), MSVD (Wu et al., 2017), ActivityNet (Caba Heilbron 116 et al., 2015), VCG Bench (Maaz et al., 2023), MVBench (Li et al., 2023c), and Video-MME (Fu 117 et al., 2024). These benchmarks encompass tasks such as video question answering, detailed video 118 captioning, and video multiple-choice questions, with video lengths ranging from seconds to hours. Furthermore, our codebase has integrated cutting-edge video LLM techniques, including video Di-119 rect Preference Optimization and video-image-multiple image interleave training. As shown in Fig. 120 1(b), compared to recent top Video LLMs, PPLLaVA demonstrates clear advantages across both 121 video and image benchmarks, while responding 7x faster than LLaVA-Next-Video-7B. 122

123 124

125

2 RELATED WORKS

**Image-domain LLMs.** Image-domain pretrained models have long served as the foundation for 126 video understanding (Carreira & Zisserman, 2017; Luo et al., 2022; Liu et al., 2023c). This is partly 127 due to the inherent similarities between image and video modalities and partly because image pre-128 training datasets offer a level of quality, quantity, and diversity that video datasets often lack. In the 129 field of multimodal LLMs, the BLIP and LLaVA series have consistently served as the foundation 130 for various video LLMs. The BLIP series is particularly notable for its Q-Former (Li et al., 2023a), 131 which acts as an intermediary between the vision encoder and the LLM. The Q-Former not only 132 enhances visual encoding but also compresses the number of visual tokens. Building on this foun-133 dation, InstructBLIP further developed the Q-Former's capability to extract instruction-aware visual 134 features, making it a preferred choice for some video LLMs (Zhang et al., 2023; Li et al., 2023d; 135 Liu et al., 2024d; Ren et al., 2024). LLaVA, a pioneer in visual instruction tuning (Liu et al., 2023a), 136 accomplished the mapping from the visual encoder to the LLM using simple linear layers or MLPs. 137 The LLaVA series has been continually updated (Liu et al., 2024a;b; Li et al., 2024), with later versions showing that this straightforward mapping approach can achieve superior results with less 138 data. This simplicity and effectiveness inspired us to use LLaVA as the foundation for our model. 139 Alongside this, we introduced the pooling module that retains LLaVA's efficient structure while also 140 enabling the compression of visual tokens and the extraction of prompt-specific visual features. 141

142 Video LLMs. In the past year, Video LLMs have experienced rapid development since their inception. For video LLMs, updating video instruction data and benchmarks is essential. Video-ChatGPT 143 (Maaz et al., 2023) was the first to introduce a high-quality video instruction training dataset and test 144 set, establishing a benchmark for GPT-assisted evaluation. MVBench (Li et al., 2023c) provides a 145 multiple-choice benchmark that assesses video performance across 20 different tasks. Video-MME 146 (Fu et al., 2024) extends video duration significantly, reaching up to several hours, and serves as 147 a comprehensive multiple-choice video QA benchmark. On the other hand, early Video LLMs 148 (Li et al., 2023b; Zhang et al., 2023; Luo et al., 2023; Maaz et al., 2023; Liu et al., 2024c) typi-149 cally used average pooling to process video sequences with Image LLMs while employing modality 150 perceivers to model temporal sequences. However, this approach significantly limited the model's 151 ability to fully understand video sequences. Alternatively, some models (Liu et al., 2024d;b; Xu 152 et al., 2024) rely on the LLM itself to model video sequences, achieving good video understanding results. Nonetheless, this method is limited to handling a small number of frames and does not 153 support the comprehension of long videos. 154

Understanding long videos is also a hot topic in video LLMs. MovieChat (Song et al., 2024) and
Flash-VStream (Zhang et al., 2024a) use memory structures to process streaming videos, while ChatUniVi (Jin et al., 2023) adopts a clustering approach for token aggregation. LLaMA-VID (Li et al.,
2023d) compresses each video frame into two tokens, capturing both local and global information.
Most similar to our work, PLLaVA (Xu et al., 2024) employs a non-parametric AdaptiveAvgPool3d
function to compress visual tokens. In contrast, our method supports not only token compression but
also the extraction of visual features pertinent to user prompts. Furthermore, our convolution-style
pooling method enables flexible output sizes. Notably, compared to the aforementioned methods,



Figure 2: The overview of PPLLaVA for compressing the video based on user prompts and generating responses on the input video and instructions.

our approach has achieved state-of-the-art results on both long and short video benchmarks, whereas
 the other methods may exhibit slightly lower performance on videos of certain lengths.

181 The diversification of data modalities and formats has also become a prominent direction in research. 182 Beyond the classic image-video instruction tuning, CAT (Ye et al., 2024) introduced mixed training 183 with video and audio, while VideoMME (Fu et al., 2024) emphasized the importance of subtitles. 184 VILA (Lin et al., 2024) and LLaVA-Interleave underscored the value of interleaved training. Besides instruction tuning, Reinforcement Learning from Human Feedback (RLHF) has also been proven to 185 be particularly effective for video LMM. Specifically, VLM-RLAIF (Ahn et al., 2024) and LLaVA-Hound (Zhang et al., 2024b) demonstrated the effectiveness of Proximal Policy Optimization (PPO) 187 and Direct Preference Optimization (DPO), respectively. We have also integrated these cutting-edge 188 techniques into our codebase and demonstrated that they can operate in parallel with PPLLaVA. 189

#### 190 191 3 Methodology

177

178

192

193

#### 3.1 MOTIVATION AND ANALYSIS

194 In the previous section, we discussed that the videos are redundant in both length and content. 195 Vista-LLaMA (Ma et al., 2024) demonstrated that the extensive number of tokens in long videos 196 makes it difficult for LLMs to capture video content. In this section, we further examine whether redundant video content impacts the performance of video LLMs and whether extracting key video 197 content can enhance performance. Inspired by EgoSchema (Mangalam et al., 2024), we adopt the certificate length to measure the redundancy. The certificate length of a video-QA pair is determined 199 by the shortest video sub-clip that can answer the question. Instead of using manual annotation, we 200 employed an automated method to determine the certificate. Specifically, frames are sampled at 2 201 fps, and then the similarity between each frame and the question-answer text is calculated using 202 CLIP-L-336 (Radford et al., 2021). If the similarity exceeds 0.5, the frame is considered relevant to 203 the text. Finally, the proportion of relevant frames to the entire video is calculated as the certificate. 204

Based on the Video-MME 205 dataset, we selected the 100 206 video-QA pairs with the 207 shortest certificate lengths 208 termed Video-MME-redund. 209 We then evaluated the per-210 formance of various models 211 on both the full Video-MME 212 dataset and these selected 213 samples. Additionally, for

Table 1: The study on the impact of video redundancy, we used the Vicuna-
7B version for all models. "Average" and "Manual" refer to the default
average frame sampling and manual frame selection, respectively.

Model	Frames	Tokens	Video-MME-full average	Video-MM average	ME-redund manual
InstructBLIP	32	1024	39.2	36.1 (-3.1)	39.5 (+0.3)
LLaVA-Next	32	4608	41.1	36.9 (-4.2)	42.0 (+0.9)
LLaVA-Next-Video	8	1152	42.9	39.0 (-3.9)	43.5 (+0.6)
LLaVA-Next-Video	32	4608	45.0	41.5 (-3.5)	46.1 (+1.1)
PPLLaVA (ours)	32	1024	49.8	47.6 (-2.2)	50.5 (+0.7)

these 100 samples, we manually selected the frames most relevant to the questions, alongside
 the default frame sampling method. This approach was used to test whether extracting key information enhances video understanding. As shown in Table 1, all models experienced a decline

216 in performance on high-redundancy videos. As an earlier model, InstructBLIP performed as 217 expected, not matching the overall performance of the more advanced LLaVA-Next. However, 218 on high-redundancy videos, InstructBLIP, which has instruction-aware video feature extraction 219 capabilities, declined slower than LLaVA-Next. Furthermore, when manually selected frames 220 were used, all models showed significant performance improvements, highlighting the importance of extracting key video information for enhancing video understanding. Additionally, we clearly 221 observed the importance of including more frames for long videos, such as those in the Video-MME 222 dataset. These findings motivated us to explore token compression to accommodate more video 223 frames while effectively extracting key information. 224

As shown in Fig. 2, PPLLaVA, like most video LLMs, includes a vision encoder, a mapping layer, and a LLM. It also features an additional text encoder paired with the visual encoder. Given a *T*-frame video, we first pass it through the CLIP-ViT visual encoder, obtaining the visual feature  $V \in \mathbb{R}^{T \times W \times H \times D}$ . This feature is then fed into the Prompt-guided Pooling module, where it is compressed by over 90%, resulting in  $V' \in \mathbb{R}^{T' \times W' \times H' \times D}$ . V' is fed into the MLP mapping layer as the final visual input. Importantly, V' not only contains significantly fewer tokens but also condenses information more relevant to the user's instructions. This ensures improved performance while efficiently processing the video input. Next, we will detail how V' is obtained.

Fine-grained Vision-Prompt Alignment. To extract video features relevant to the prompt, we first utilize the original CLIP dual encoders to identify which video features are related to the text. Specifically, we input the user's question into the CLIP text encoder to obtain the text feature  $c \in \mathbb{R}^D$ . Following the CLIP training pipeline, we only use the CLS token of the text. The attention score of the  $(t^{th}, w^{th}, h^{th})$  video token relative to the text feature is then calculated as:

240

225 226

241

242

256 257

263 264  $s_{(t,w,h)} = \frac{\exp(\tau c \cdot f_{clipv}(v_{(t,w,h)}))}{\sum_{t=1}^{T} \sum_{w=1}^{W} \sum_{h=1}^{H} \exp(\tau c \cdot f_{clipv}(v_{(t,w,h)}))},$ (1)

where  $v_{(t,w,h)}$  represents the token at the (t, w, h) position in  $V, \tau$  is the CLIP temperature scale, and  $f_{clipv}$  is the CLIP visual projection, which is typically not used in multimodal LLMs. Note that  $v_{(t,w,h)}$  typically refers to the patch token from the penultimate layer of CLIP, rather than the CLS token from the final layer used during CLIP training. However, since the spatial representations in CLIP's final layers are similar, applying  $f_{clipv}$  still allows the patch tokens to be mapped into the interaction space with the text.

**Prompt-Guided Pooling.** In the previous section, we obtained token-level weights corresponding to the user's prompt, which we use as guidance for pooling the video. Unlike traditional tasks that require only a D-dimensional feature for contrastive learning (Ma et al., 2022; Wang et al., 2022), our approach aims to preserve a certain 3-dimensional structure to enable the LLM to perform temporal modeling. To achieve this, we perform pooling with  $S = \{s_{(t,w,h)}\}$  in a manner similar to 3D convolution. Specifically, we define the spatiotemporal 3D convolution kernel and stride as  $(k_t, k_w, k_h)$  and  $(d_t, d_w, d_h)$ , respectively. The output dimension of V' can then be expressed as:

$$T' = \left(\frac{T - k_t}{d_t}\right) + 1, \quad W' = \left(\frac{W - k_w}{d_w}\right) + 1, \quad H' = \left(\frac{H - k_h}{d_h}\right) + 1. \tag{2}$$

Unlike conventional convolution kernels, our kernel parameters are derived from S. Moreover, the parameters of the kernel are dynamic; as the kernel slides over different positions in V, its parameters are taken from the corresponding positions in S. Finally, the feature at position (t, w, h) in the output V' is calculated as:

$$v'_{(t,w,h)} = \sum_{i=0}^{k_t-1} \sum_{j=0}^{k_w-1} \sum_{k=0}^{k_h-1} v_{(t*d_t+i,w*d_w+j,h*d_h+k)} s_{(t*d_t+i,w*d_w+j,h*d_h+k)}.$$
(3)

By flexibly adjusting the stride and kernel size, we can control the output dimensions. This approach allows us to better accommodate videos of varying lengths and facilitates joint training with images, compared to fixed-output methods.

269 **CLIP Context Extension.** In our method, CLIP-text is the only additional parameter used. Despite having significantly fewer parameters than Qformer, it achieves better performance. However,

270 CLIP-text has a major limitation: its context length is too short (default is 77). While this length 271 is sufficient for objects or simple descriptions, it is inadequate for long prompts or multi-turn di-272 alogues in multimodal LLMs. To address this performance bottleneck, we propose extending the 273 context length of CLIP-text using asymmetric positional embedding extensions. In most cases, ex-274 tending the positional embedding involves randomly initializing new embeddings at the end. A more theoretically sound approach is to perform linear interpolation on the original positional embedding 275 at a rate of r. Assuming the original and target positional embeddings are P and P', respectively, 276 the  $i^{th}$  position of P' can be represented as: 277

- 278
- 279

$$P'_{i} = P_{\lfloor j \rfloor} + (j - \lfloor j \rfloor) \cdot (P_{\lfloor j \rfloor + 1} - P_{\lfloor j \rfloor}), \quad j = i \cdot r, \tag{4}$$

where |j| means taking the floor of j. However, we found linear interpolation yielded inferior re-280 sults to randomly initializing embeddings at the end. We believe this is because CLIP's positional 281 embeddings are well-trained, and globally averaged interpolation disrupts the well-pre-trained in-282 formation. Given that short sentences dominate CLIP's training data, the earlier parts of positional 283 embeddings are more thoroughly trained. Hence, we adopted asymmetric interpolation, applying 284 different interpolation rates at different positions. In the early part of the new positional embedding, 285 we use a large r value to shorten the interpolation distance, while in the later part, we use a smaller r286 value to extend the interpolation distance. This asymmetric approach allows us to effectively extend 287 the context length of CLIP-text while preserving as much of the pre-trained information as possible. 288

289 3.3 TRAINING

Interleave Instruction Tuning. PPLLaVA enables plug-and-play transfer of image-domain LLMs 291 to the video domain. As a result, initialized from well-pretrained image LLM, we can bypass expen-292 sive contrastive or alignment pretraining and proceed directly to instruction tuning. In this stage, we 293 fully fine-tune the LLM, the projection MLP, and the CLIP text encoder. Our instruction datasets 294 include multi-turn and single-turn conversations presented in a conversational format, along with 295 various forms of visual input such as images, videos, and multiple images. For different types of 296 data, we employed an interleaving training approach. Rather than using batches composed of a 297 single data type, we mixed various data types within the same batch. Numerous studies (Li et al., 298 2024; Laurençon et al., 2024; Xue et al., 2024) have demonstrated that this method is the most nat-299 ural approach for handling multimodal data. Additionally, this training method enables the model 300 to simultaneously process both long videos with many frames and single-frame images, greatly 301 enhancing its adaptability to visual sequences of varying lengths.

302 Direct Preference Optimization. (DPO) Video, especially long video-based dialogue, is more 303 prone to hallucinations compared to images. As a result, Reinforcement Learning from Human 304 Feedback (RLHF) (Zhang et al., 2024b; Ahn et al., 2024) has proven particularly effective for video. 305 Therefore, we also implemented this method based on our model. Following LLaVA-Hound (Zhang 306 et al., 2024b). We used detailed video captions as proxies for video content and performed DPO 307 with feedback from the language model serving as a reward. In this stage, all parameters except the 308 LLM were frozen, and only video data was used. This additional phase significantly reduced the occurrence of hallucinations during video-based dialogue. 309

311 4 EXPERIMENTS

310

312

In this section, we have performed comprehensive experimental evaluations of PPLLaVA, covering crucial settings, comparisons, and ablations, while more ablation studies, visualizations, and limitations analysis can be found the appendix.

316 317 4.1 EXPERIMENT SETUP

**Implementation Details.** PPLLaVA is built upon the advanced image-domain LLaVA-Next models (Liu et al., 2024b). To ensure a fair comparison with most models, we chose the Vicuna-7B version. For image and multiple-image inputs, the pooling kernel and strides are set to (1, 3, 3). For video inputs, we uniformly sample 32 frames and set the pooling kernel and strides to (2, 3, 3), compressing the video tokens by over 15 times. During training, both questions and answers are fed into the CLIP text encoder to better capture prompt-vision relevance. For CLIP context extension, when i < 20, r is set to 1, and when  $i \ge 20, r$  is set to 0.25. We train for one epoch using a learning rate

324	Table 2: The results of open-ended QA with GPT-based evaluation, including MSVD-QA
325	MSRVTT-QA, ActivityNet-QA(ANet), and VCG Bench. All the models are based on the Vicuna-
326	7B. † means using DPO or PPO.
327	

327		MS	VD	MSR	VTT	A	Net			VCG	Bench		
328	Method	Acc.	Sco.	Acc.	Sco.	Acc.	Sco.	CI	DO	CU	TU	СО	Avg.
329	VideoChat (Li et al., 2023b)	56.3	2.8	45.0	2.5	26.5	2.2	2.23	2.50	2.53	1.94	2.24	2.29
220	Video-ChatGPT (Maaz et al., 2023)	64.9	3.3	49.3	2.8	35.2	2.7	2.50	2.57	2.69	2.16	2.20	2.42
330	BT-Adapter (Liu et al., 2024c)	67.5	3.7	57.0	3.2	45.7	3.2	2.68	2.69	3.27	2.34	2.46	2.69
331	Video-LLaVA (Lin et al., 2023)	70.7	3.9	59.2	3.5	45.3	3.3	-	-	-	-	-	-
222	MovieChat (Song et al., 2024)	75.2	3.8	52.7	2.6	45.7	3.4	2.76	2.93	3.01	2.24	2.42	2.67
332	Chat-UniVi (Jin et al., 2023)	65.0	3.6	54.6	3.1	45.8	3.2	2.89	2.91	3.46	2.89	2.81	2.99
333	VideoChat2 (Li et al., 2023c)	70.0	3.9	54.1	3.3	49.1	3.3	3.02	2.88	3.51	2.66	2.81	2.98
334	Vista-LLaMA (Ma et al., 2024)	65.3	3.6	60.5	3.3	48.3	3.3	2.44	2.64	3.18	2.26	2.31	2.57
004	LLaMA-VID (Li et al., 2023d)	69.7	3.7	57.7	3.2	47.4	3.3	2.96	3.00	3.53	2.46	2.51	2.89
335	ST-LLM (Liu et al., 2024d)	74.6	3.9	63.2	3.4	50.9	3.3	3.23	3.05	3.74	2.93	2.81	3.15
336	PLLaVA (Xu et al., 2024)	76.6	4.1	62.0	3.5	56.3	3.5	3.21	2.86	3.62	2.33	2.93	2.99
000	CAT (Ye et al., 2024)	-	-	62.1	3.5	50.2	3.5	3.08	2.95	3.49	2.81	2.89	3.07
337	VLM-RLAIF † (Ahn et al., 2024)	76.4	4.0	63.0	3.4	57.3	3.5	3.85	3.45	3.84	3.63	2.8	3.49
338	LLaVA-Next-Video (Liu et al., 2024b)	-	-	-	-	53.5	3.2	3.39	3.29	3.92	2.60	3.12	3.26
000	LLaVA-Next-Video †	-	-	-	-	60.2	3.5	3.64	3.45	4.17	2.95	4.08	3.66
339	PPLLaVA	75.8	3.9	61.9	3.3	56.1	3.4	3.32	3.20	3.88	3.00	3.20	3.32
340	PPLLaVA †	77.1	4.0	64.3	3.5	60.7	3.6	3.85	3.56	4.21	3.21	3.81	3.73

Table 3: Performance on Video-MME with short, medium, and long durations, under the settings of "without subtitles" and "with subtitles". \* means using multi-images during training.

Models	LLM	Short	(%)	Mediu	m (%)	Long	(%)	Overa	ll (%)
Wodels	Params	w/o subs	w/ subs						
Qwen-VL-Chat (Bai et al., 2023)	7B	46.9	47.3	38.7	40.4	37.8	37.9	41.4	41.9
Qwen-VL-Max (Bai et al., 2023)	-	55.8	57.6	49.2	48.9	48.9	47.0	51.3	51.2
InternVL-V1.5 (Chen et al., 2024)	20B	60.2	61.7	46.4	49.1	45.6	46.6	50.7	52.4
Video-LLaVA	7B	45.3	46.1	38.0	40.7	36.2	38.1	39.9	41.6
ST-LLM	7B	45.7	48.4	36.8	41.4	31.3	36.9	37.9	42.3
VideoChat2-Mistral	7B	48.3	52.8	37.0	39.4	33.2	39.2	39.5	43.8
Chat-UniVi-V1.5	7B	45.7	51.2	40.3	44.6	35.8	41.8	40.6	45.9
LLaVA-NeXT-Video	7B	45.9	49.8	40.3	44.3	36.6	41.0	40.9	45.0
LLaVA-NeXT-Video	34B	61.7	65.1	50.1	52.2	44.3	47.2	52.0	54.9
PPLLaVA	7B	56.1	59.7	43.9	48.6	38.4	44.0	46.1	50.0
PPLLaVA*	7B	58.7	62.8	45.6	50.4	42.2	47.4	48.8	53.6

> of 2e-5 and a batch size of 256. We provid both GPU and NPU versions, and the full training takes 24 hours on 16 A100 GPUs or 32 910B NPUs.

**Data Details.** The instruction tuning data includes diverse modalities and sources. We randomly sampled 300k image data from the LLAVA-1.5 training set (Liu et al., 2024a) and used 594k multiple-image data from LLAVA-Interleave (Liu et al., 2024b). The video data includes Kinet-ics (Kay et al., 2017), SthSth-V2 (Goyal et al., 2017), Next-QA (Xiao et al., 2021), CLEVRER (Yi et al., 2019), and LLAVA-Interleave-300k, resulting in a total of 1.36M multimodal training sam-ples. Notably, to ensure fairness in the comparison experiments, we excluded multi-image data and used only 760k image-video data, comparable to the training volume of most video LLMs.

We evaluate our model on six video LLM benchmarks, categorized into two types based on the eval-uation method: GPT-based evaluation and multiple-choice questions. The GPT evaluation mainly involves open-ended QA, including the Video-based Generative Performance Benchmark (VCG Bench) (Maaz et al., 2023), MSVD-QA (Wu et al., 2017), MSRVTT-QA (Xu et al., 2016), and ActivityQA (Caba Heilbron et al., 2015). Consistent with most models, we used the GPT-3.5-turbo-0613 version for testing. The multiple-choice question benchmarks include MVBench (Li et al., 2023c) and Video-MME (Fu et al., 2024). This evaluation method is more objective by eliminat-ing the potential disturbances of GPT. For medium-to-long videos in Video-MME, we sampled 64 frames instead of the 32 frames used in other datasets. Our test corpus encompasses videos of various genres and lengths, offering a comprehensive evaluation of PPLLaVA's performance. 

4.2 **QUANTITATIVE RESULT** 

GPT-Based Evaluation. Table 2 presents the quantitative results for open-ended question-answering, showing that PPLLaVA achieves top performance across all datasets. It also demon-strates a significant performance gap compared to models other than LLaVA-Next-Video, demon-

Method	AS	AP	AA	FA	UA	OE	OI	os	MD	AL	ST	AC	MC	MA	SC	FP	со	EN	ER	CI	Avg.
Video-LLaMA	27.5	25.5	51.0	29.0	39.0	48.0	40.5	38.0	22.5	22.5	43.0	34.0	22.5	32.5	45.5	32.5	40.0	30.0	21.0	37.0	34.1
LLaMA-Adapter	23.0	28.0	51.0	30.0	33.0	53.5	32.5	33.5	25.5	21.5	30.5	29.0	22.5	41.5	39.5	25.0	31.5	22.5	28.0	32.0	31.7
Video-ChatGPT	23.5	26.0	62.0	22.5	26.5	54.0	28.0	40.0	23.0	20.0	31.0	30.5	25.5	39.5	48.5	29.0	33.0	29.5	26.0	35.5	32.7
VideoChat	33.5	26.5	56.0	33.5	40.5	53.0	40.5	30.0	25.5	27.0	48.5	35.0	20.5	42.5	46.0	26.5	41.0	23.5	23.5	36.0	35.5
VideoChat2	66.0	47.5	83.5	49.5	60.0	58.0	71.5	42.5	23.0	23.0	88.5	39.0	42.0	58.5	44.0	49.0	36.5	35.0	40.5	65.5	51.1
ST-LLM	66.0	53.5	84.0	44.0	58.5	80.5	73.5	38.5	42.5	31.0	86.5	36.5	56.5	78.5	43.0	44.5	46.5	34.5	41.5	58.5	54.9
PLLaVA-7B	58.0	49.0	55.5	41.0	61.0	56.0	61.0	36.0	23.5	26.0	82.0	39.5	42.0	52.0	45.0	42.0	53.5	30.5	48.0	31.0	46.6
PLLaVA-13B	66.0	53.0	65.5	45.0	65.0	58.0	64.5	35.5	23.5	30.0	85.0	39.5	45.5	57.0	47.5	49.5	49.0	33.0	53.0	37.0	50.1
PLLaVA-34B	67.5	53.0	82.0	47.0	79.0	68.5	67.5	36.5	37.5	49.5	91.0	40.5	43.0	70.0	51.5	50.0	66.5	39.5	63.5	59.0	58.1
GPT-4V	55.5	63.5	72.0	46.5	73.5	18.5	59.0	29.5	12.0	40.5	83.5	39.0	12.0	22.5	45.0	47.5	52.0	31.0	59.0	11.0	43.5
PPLLaVA-7B	69.0	54.4	69.5	50.5	69.0	87.0	67.0	38.0	35.0	33.0	69.5	37.5	63.5	91.0	47.5	47.5	51.5	27.0	47.5	57.5	57.1
PPLLaVA-7B*	73.5	61.0	83.5	45.5	68.0	87.5	75.5	33.0	37.5	40.0	83.0	37.0	67	96.5	50.5	43.5	59.0	35.5	44.5	63.0	59.2

378 Table 4: Results on MVBench. Models without additional annotation are 7B by default. \* means 379 using multi-images during training.

Table 5: The ablation study on model components. TP means throughput (seconds/video).

Madal	Context			VC	G Be	nch				Video-M	ME (w	/ subs)	
Wodel	Length	CI	DO	CU	TU	CO	Avg	TP	Short	Medium	Long	Overall	TP
LLaVA-Next (Average Pooing)	576	3.05	3.07	3.71	2.62	3.01	3.09	2.9	53.1	41.3	36.0	43.4	3.1
LLaVA-Next (w/o Pooing)	4608	3.23	3.08	3.82	2.75	3.11	3.20	15.0	58.4	45.1	38.8	47.4	15.2
+Prompt-guided Pooling	1024	3.21	3.15	3.80	2.88	3.02	3.21	4.6	59.0	45.6	42.2	48.9	5.3
+CLIP Context Extension	1024	3.32	3.20	3.88	3.00	3.20	3.32	4.6	59.7	48.6	44.0	50.0	5.3

strating the strong text generation capability of our model. Despite using lower-quality data (as 396 LLaVA 1.6 data is not publicly available), PPLLaVA outperforms LLaVA-Next-Video. More im-397 portantly, PPLLaVA uses significantly fewer visual contexts (1024 vs. 4096), resulting in higher 398 throughput. After applying DPO, PPLLaVA also shows consistent improvements and outperforms 399 other models that use DPO or PPO, further proving the adaptability of the PPLLaVA architecture 400 across different training stages. 401

Video-MME. Although Video-MME is a new benchmark, it offers high quality and data diversity. 402 Its inclusion of hour-long videos makes it particularly effective for evaluating models' long video un-403 derstanding capabilities. As shown in Table 3, PPLLaVA achieves the best results on Video-MME, 404 with a notably significant advantage on videos of different lengths compared to other models. The 405 7B model's long video comprehension already surpasses the 34B LLaVA-Next-Video, as PPLLaVA 406 efficiently compresses video tokens, enabling support for a much higher number of frames than 407 LLaVA-Next-Video, thereby enhancing long video understanding capabilities. 408

**MVBench.** MVBench is a multiple-choice benchmark offering a comprehensive set of evalua-409 tion tasks, including Action Sequence (AS), Action Prediction (AP), Action Antonym (AA), Fine-410 grained Action (FA), Unexpected Action (UA), Object Existence (OE), Object Interaction (OI), Ob-411 ject Shuffle (OS), Moving Direction (MD), Action Localization (AL), Scene Transition (ST), Action 412 Count (AC), Moving Count (MC), Moving Attribute (MA), State Change (SC), Fine-grained Pose 413 (FP), Character Order (CO), Egocentric Navigation (EN), Episodic Reasoning (ER), Counterfactual 414 Inference (CI), and the average across all 20 metrics (Avg). PPLLaVA achieves the best average re-415 sults among models, demonstrating a clear advantage and strong adaptability in video understanding 416 across diverse scenarios, especially for moving and action tasks.

417

418 4.3 ABLATIONS AND ANALYSIS 419

420 Model Components. The core of PPLLaVA is its prompt-guided token compression, which en-421 hances both video understanding efficiency and performance. To assess the impact of this feature, 422 we conducted ablation experiments on the overall model components. As shown in Table 5, while 423 the LLaVA-Next Baseline's direct averaging method is the most efficient, its performance is subpar. Directly feeding all tokens into the LLM yields reasonable results but suffers from extremely 424 low throughput. Our Pooling module substantially improves both efficiency and performance. Ex-425 tending the CLIP context further enhances results, particularly in long video understanding. The 426 simultaneous improvement in efficiency and effectiveness underscores the superiority of our model. 427

428 **Pooling Size.** PPLLaVA can flexibly implement pooling at any scale. However, as the pooling kernel and stride increase, while efficiency improves, there will inevitably be performance degrada-429 tion. Therefore, it's crucial to find a pooling size that balances both efficiency and performance. As 430 illustrated in Fig. 3, we first explore the impact of pooling in the spatial dimension. It is evident 431 that when the pooling kernel and stride are small, increasing them significantly improves efficiency,



Figure 3: Spatial pooling effects. We set T = 16 and  $k_t = d_t = 1$ , varying the spatial kernel size and stride.



Figure 4: Temporal pooling effects. We set T = 32 and  $k_w = d_w = k_h = d_h = 3$ , varying the temporal kernel size and stride.

Table 6: The image results. \* means self-implementation.

Model	Resolution	MMMU(val)	MathVista	MMB-ENG	MMB-CN	MM-Vet	SEED-IMG	MME	POPE
LLaVA-1.5-13B	336*336	36.4	27.6	67.8	63.3	36.3	68.2	1531/295	85.93
LLaVA-Next-7B	672*672	35.8	34.6	67.4	60.6	43.9	70.2	1519/332	86.53
VideoLLaVA	336*336	-	-	60.9	-	32.0	-	-	84.40
Chat-Univ-1.5	336*336	-	-	62.7	-	28.3	-	-	85.40
LLaVA-Next-Video *	336*336	34.2	28.9	64.7	56.7	44.0	64.6	1501/ <b>351</b>	83.10
PPLLaVA	336*336	37.9	34.6	68.9	62.0	44.7	70.7	<b>1539</b> /277	88.46

and thanks to the prompt-guided approach, the performance remains almost unaffected. In contrast, 453 as shown in Fig. 4, pooling in the temporal dimension yields smaller efficiency gains compared 454 to spatial scaling, with more noticeable performance degradation as the kernel and stride sizes in-455 crease. When the pooling kernel and stride are large, the efficiency gains tend to plateau, but the 456 decline in effectiveness becomes significantly pronounced, particularly in spatial pooling, where the 457 performance drop is more severe. Considering all factors, for video input, we ultimately selected a 458 pooling kernel and stride of (2, 3, 3) to ensure a substantial improvement in efficiency while main-459 taining stable model performance. 460

Image Performance. Theoretically, further video tuning on top of an image-domain LLM could 461 lead to catastrophic forgetting of pre-trained knowledge and image understanding. The PPLLaVA 462 method can also be seamlessly applied to images. Although images do not have the same need for 463 token compression as videos, and compression may lead to performance loss, the guidance from 464 user prompts can still similarly enhance performance. In Table 6, we present PPLLaVA's results on 465 various popular image LLM benchmarks. Since PPLLaVA was trained on LLaVA-1.5 image data 466 based on LLaVA-Next, we compared the results of these two models. We also compare the image 467 performance with LLaVA-Next-Video and other image-video unified models. As shown, PPLLaVA 468 shows a significant advantage in image performance compared to video models, indicating that 469 PPLLaVA has effectively retained pre-trained knowledge. Compared to image models, despite using a smaller LLM or lower image resolution, PPLLaVA, as a video model, still achieved better 470 results on most benchmarks. Notably, our pooling method reduced the visual tokens to one-ninth 471 of the original count at the same resolution. This demonstrates that PPLLaVA can achieve both 472 performance and efficiency improvements even on image-based tasks, highlighting its potential for 473 lightweight multimodal LLM. 474

474 475 476

443

444

445 446

5 CONCLUSION

477 In this paper, we propose Prompt-guided Pooling LLaVA (PPLLaVA), a novel pooling method that 478 achieves token compression and prompt-aware feature extraction simultaneously. We first observed 479 that current video LLMs struggle to balance performance on both long and short videos. Further 480 analysis revealed that redundant tokens in videos negatively impact video understanding perfor-481 mance. To address this, our model incorporates three key modules: Fine-grained Vision-Prompt 482 Alignment, Prompt-Guided Convolution-Style Pooling, and CLIP Context Extension. These mod-483 ules significantly reduce the visual context while effectively extracting essential visual features. Extensive experiments have demonstrated the effectiveness of PPLLaVA on both images and videos, as 484 it achieves the best results across benchmarks of various tasks and video lengths, ensuring excellent 485 efficiency, with particularly outstanding performance on long videos.

# 486 REFERENCES

498

499

500

501

505

506

507

508

512

513

514

515

516

526

527

528

529

- Daechul Ahn, Yura Choi, Youngjae Yu, Dongyeop Kang, and Jonghyun Choi. Tuning large
   multimodal models for videos using reinforcement learning from ai feedback. *arXiv preprint arXiv:2402.03746*, 2024.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:
   A large-scale video benchmark for human activity understanding. In *Proceedings of the ieee* conference on computer vision and pattern recognition, pp. 961–970, 2015.
  - Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, 2017.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,
   Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to
   commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
  - W Dai, J Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. Instructblip: towards general-purpose vision-language models with instruction tuning. arxiv. *Preprint posted online on June*, 15:2023, 2023.
- Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
  - Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pp. 5842–5850, 2017.
- Tengda Han, Weidi Xie, and Andrew Zisserman. Temporal alignment networks for long-term video.
   In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2906–2916, 2022.
- Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. Vtimellm: Empower llm to grasp video moments. *arXiv preprint arXiv:2311.18445*, 2023.
- Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified vi sual representation empowers large language models with image and video understanding. *arXiv preprint arXiv:2311.08046*, 2023.
  - Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36, 2024.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
  Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv* preprint arXiv:2407.07895, 2024.
- 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*, 2023a.

550

551

554

563

565

572

578

579

580

581

540	KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang,
541	and Yu Qiao. Videochat: Chat-centric video understanding. arXiv preprint arXiv:2305.06355,
542	2023b.
543	

- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, 544 Ping Luo, et al. Mybench: A comprehensive multi-modal video understanding benchmark. arXiv preprint arXiv:2311.17005, 2023c. 546
- 547 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language 548 models. arXiv preprint arXiv:2311.17043, 2023d.
  - Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-Ilava: Learning united visual representation by alignment before projection. arXiv preprint arXiv:2311.10122, 2023.
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pre-552 training for visual language models. In Proceedings of the IEEE/CVF Conference on Computer 553 Vision and Pattern Recognition, pp. 26689–26699, 2024.
- 555 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. arXiv 556 preprint arXiv:2304.08485, 2023a.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 558 tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-559 tion, pp. 26296–26306, 2024a. 560
- 561 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 562 Llava-next: Improved reasoning, ocr, and world knowledge, 2024b.
  - Ruyang Liu, Jingjia Huang, Wei Gao, Thomas H Li, and Ge Li. Mug-stan: Adapting image-language pretrained models for general video understanding. arXiv preprint arXiv:2311.15075, 2023b.
- Ruyang Liu, Jingjia Huang, Ge Li, Jiashi Feng, Xinglong Wu, and Thomas H Li. Revisiting temporal 566 modeling for clip-based image-to-video knowledge transferring. In Proceedings of the IEEE/CVF 567 Conference on Computer Vision and Pattern Recognition, pp. 6555–6564, 2023c. 568
- 569 Ruyang Liu, Chen Li, Yixiao Ge, Thomas H Li, Ying Shan, and Ge Li. Bt-adapter: Video conver-570 sation is feasible without video instruction tuning. In Proceedings of the IEEE/CVF Conference 571 on Computer Vision and Pattern Recognition, pp. 13658–13667, 2024c.
- Ruyang Liu, Chen Li, Haoran Tang, Yixiao Ge, Ying Shan, and Ge Li. St-llm: Large language 573 models are effective temporal learners. arXiv preprint arXiv:2404.00308, 2024d. 574
- 575 Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An 576 empirical study of clip for end to end video clip retrieval and captioning. *Neurocomputing*, 508: 293-304, 2022. 577
  - 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.
- Fan Ma, Xiaojie Jin, Heng Wang, Yuchen Xian, Jiashi Feng, and Yi Yang. Vista-Ilama: Reducing 582 hallucination in video language models via equal distance to visual tokens. In Proceedings of the 583 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13151–13160, 2024. 584
- 585 Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-586 end multi-grained contrastive learning for video-text retrieval. In Proceedings of the 30th ACM International Conference on Multimedia, pp. 638–647, 2022.
- 588 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: 589 Towards detailed video understanding via large vision and language models. arXiv preprint 590 arXiv:2306.05424, 2023. 591
- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic bench-592 mark for very long-form video language understanding. Advances in Neural Information Processing Systems, 36, 2024.

- 594 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 595 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 596 models from natural language supervision. In International conference on machine learning, pp. 597 8748-8763. PMLR, 2021.
- 598 Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multimodal large language model for long video understanding. In Proceedings of the IEEE/CVF Conference 600 on Computer Vision and Pattern Recognition, pp. 14313–14323, 2024. 601
- 602 Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe 603 Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In Proceedings of the IEEE/CVF Conference on Computer Vision 604 and Pattern Recognition, pp. 18221-18232, 2024. 605
- 606 Qiang Wang, Yanhao Zhang, Yun Zheng, Pan Pan, and Xian-Sheng Hua. Disentangled representation learning for text-video retrieval. arXiv preprint arXiv:2203.07111, 2022. 608
- Zuxuan Wu, Ting Yao, Yanwei Fu, and Yu-Gang Jiang. Deep learning for video classification and 609 captioning. In Frontiers of multimedia research, pp. 3–29. ACM, 2017. 610
- 611 Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-612 answering to explaining temporal actions. In Proceedings of the IEEE/CVF conference on com-613 puter vision and pattern recognition, pp. 9777–9786, 2021. 614
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging 615 video and language. In Proceedings of the IEEE conference on computer vision and pattern 616 recognition, pp. 5288-5296, 2016. 617
- 618 Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. Pllava: 619 Parameter-free llava extension from images to videos for video dense captioning. arXiv preprint 620 arXiv:2404.16994, 2024.
- Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj 622 Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal models. arXiv preprint arXiv:2408.08872, 2024. 624
- 625 Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. arXiv 626 preprint arXiv:2403.04640, 2024. 627
  - Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. Clevrer: Collision events for video representation and reasoning. arXiv preprint arXiv:1910.01442, 2019.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language 632 model for video understanding. arXiv preprint arXiv:2306.02858, 2023. 633
  - Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-based real-time understanding for long video streams. arXiv preprint arXiv:2406.08085, 2024a.
- Ruohong Zhang, Liangke Gui, Zhiqing Sun, Yihao Feng, Keyang Xu, Yuanhan Zhang, Di Fu, Chun-638 yuan Li, Alexander Hauptmann, Yonatan Bisk, et al. Direct preference optimization of video large 639 multimodal models from language model reward. arXiv preprint arXiv:2404.01258, 2024b. 640
- 641 Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Na-642 grani, and Cordelia Schmid. Streaming dense video captioning. In Proceedings of the IEEE/CVF 643 Conference on Computer Vision and Pattern Recognition, pp. 18243–18252, 2024.
- 644

621

623

628

629

630

631

634

635

636

637

- 645 646
- 647

Table 7: The ablation study on the Pooling Approach. We report the overall performance of VideoMME (w/ subs).

Table 8: The study on whether PPLLaVA helps long video understanding. We report the Long performance of VideoMME (w/ subs).

Pooling Method	kernel1	kernel2	tokens	Overall	train		te	st	tokono	Long
weighted average	(2,3,3)	-	1024	53.6	frames	kernel	frames	kernel	tokens	Long
separate S-T	-	-	608	44.1	32	(2,3,3)	32	(2,3,3)	1024	45.7
max pooling	(2,3,3)	-	1024	52.0	32	(2,3,3)	64	(4,3,3)	1024	47.4
multiple	(1,6,6)	(8,2,2)	1088	52.8	16	(1,3,3)	16	(1,3,3)	1024	43.5
multiple	(4,3,3)	(2,4,4)	1088	53.2	8	(1,1,1)	8	(1,1,1)	4608	41.2

Table 9: The study on multimodal data with interleave training and DPO training.

Model	Video	Image	Multi-Image	Interleave	DPO	VcgBench	MvBench	VideoMME
LL aVA Navt Vidaa	$\checkmark$	$\checkmark$		$\checkmark$		3.26	-	-
LLa VA-INEXT-VIGEO	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	3.66	-	-
	$\checkmark$					3.20	55.0	48.9
	$\checkmark$	$\checkmark$				3.09	49.8	44.1
PPLLaVA	$\checkmark$	$\checkmark$		$\checkmark$		3.32	57.1	50.0
	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		3.21	59.2	53.6
	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	3.73	55.8	49.3



Figure 5: The visualization of the attention weights used to guide video pooling.

# A MORE ANALYSIS

Prompt-Guided Pooling Approach. Beyond the weighted average pooling detailed in the main text, we experimented with several alternative pooling methods guided by the prompt. First, we applied separate spatiotemporal pooling, conducting pooling operations independently on the tem-poral and spatial dimensions before concatenation. We also explored combinations of different pooling sizes to assess their impact. Lastly, we implemented max pooling using weights derived from the prompt as guidance. As shown in Table 7, spatiotemporal separate pooling demonstrates the worst performance, underscoring the importance of maintaining the 3-dimensional spatiotemporal structure during pooling. Max pooling, though slightly better, still falls short, suggesting that a few prominent features are insufficient to represent the entirety of the video. The combination of various pooling kernels performs similarly to direct weighted averaging when the context length is comparable. Consequently, we opted for weighted averaging, as it provides optimal results while maintaining a simpler structure. 

**Is PPLLaVA Really Helpful for Long Video?** Token compression is a key feature of PPLLaVA, primarily aimed at enhancing the understanding of long videos. To validate PPLLaVA's effectiveness

Figure 6: The visualization of the attention weights used to guide video pooling. The prompts for all videos are "Describe this video.".



Figure 7: Qualitative result of video summary and detailed video description.

in this regard, we conducted an additional ablation experiment. As shown in Table 8, we first
experimented with larger pooling kernels to accommodate more frames. The results indicate that
even with some discrepancies between training and testing, using a greater number of frames still
improves long video comprehension. When fewer frames are used during training, the disadvantage
in long video understanding becomes even more pronounced. In the most extreme case, when no
pooling is applied, even with significantly longer context lengths, the understanding of long videos
is the weakest due to the limited number of frames that can be accommodated.

**Interleave training and DPO.** Multimodal data interleaved training and DPO are two key techniques utilized in this work. We conducted an analysis of their effects and compatibility with





810 PPLLaVA. As shown in Table 9, when different data modalities are not mixed within a batch, adding 811 image or multi-image modalities leads to performance degradation compared to pure video training. 812 This aligns with the conclusions of LLaVA-Next-Video. When data modalities are mixed within 813 the same batch, the additional images enhance performance. However, we found that when further 814 adding multi-image data, the performance on the multiple-choice benchmark improved, but the performance on the caption generation benchmark declined. This indicates that multi-image data can 815 enhance the model's visual knowledge but may reduce its capability in video-based dialogue. In 816 contrast, DPO training has a minimal side effect on multiple-choice benchmarks but significantly 817 improves results in GPT-based evaluation. This highlights DPO's ability to effectively reduce hal-818 lucinations in LLM outputs, leading to higher-quality dialogues. Moreover, when compared to the 819 baseline and LLaVA-Next-Video, the combination of DPO and PPLLaVA yields similar improve-820 ments. This emphasizes the strong compatibility between PPLLaVA and DPO. 821

822

## **B** QUALITATIVE RESULTS

823 824

In Fig. 5, we visualize the attention weights used to guide video pooling based on the user prompts. 825 For the same video, we tried different questions. It can be clearly observed that the model's attention 826 shifts noticeably depending on the question. For example, when the user asks about the girl's feel-827 ings, the attention is significantly focused on her face. Conversely, when asked about the number 828 of 3D objects in the video, the attention shifts more toward the 3D objects. These visualizations 829 demonstrate that while reducing the visual context, PPLLaVA effectively captures the key informa-830 tion in the video. In Fig. 6, we additionally illustrate the attention weights for captioning-related 831 questions, as these questions theoretically provide less informational content. As shown in the fig-832 ure, prompts like "Describe this video," which lack specific references, result in attention weights 833 being evenly distributed across the foreground. This indicates that our model still plays a significant 834 role in handling captioning-related questions. In Fig. 7 and 8, we further present some examples of video dialogue. As shown in Fig. 7, for the famous Sora video, PPLLaVA can accurately and 835 intricately describe details about the protagonist and the environment. For the more complex scene 836 changes in the trailer for Black Myth Wu Kong, PPLLaVA remarkably captures the details of each 837 scene and character. In Fig. 8, PPLLaVA maintains accuracy and consistency across multiple rounds 838 of dialogue and is capable of making reasonable inferences on open-ended questions. 839

840 841

842

## C LIMITATION

Although the 7B PPLLaVA has demonstrated impressive performance, even rivaling that of 34B video LLMs, our biggest regret is that, due to a lack of computational resources, we were unable to train larger-scale LLMs to uncover the limits of this architecture. Additionally, the conflict between the enhanced understanding capabilities brought by multi-image data and the decline in dialogue abilities remains unsolved in this work; a reasonable data allocation ratio might address this issue. We leave these problems for future work.

- 849
- 850
- 851
- 852 853
- 854
- 855
- 856
- 857 858
- 859
- 860
- 861
- 862
- 863