

# From Seconds to Hours: Reviewing MultiModal Large Language Models on Comprehensive Long Video Understanding

Anonymous ACL submission

## Abstract

The integration of Large Language Models (LLMs) with visual encoders has recently shown promising performance in visual understanding tasks, leveraging their inherent capability to comprehend and generate human-like text for visual reasoning. This paper reviews the advancements in MultiModal Large Language Models (MM-LLMs) for long video understanding. We highlight the unique challenges posed by long videos, including fine-grained spatiotemporal details, dynamic events, and long-term dependencies. We summarize the progress in model design and training methodologies for MM-LLMs understanding long videos and compare their performance on various long video understanding benchmarks. Finally, we discuss future directions for MM-LLMs in long video understanding.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable versatility and capability in understanding and generating human-like text by scaling model size and training data (Raffel et al., 2020; Brown, 2020; Chowdhery et al., 2023; Touvron et al., 2023a). To extend these capabilities to visual understanding tasks, various methods have been proposed to integrate LLMs with specific visual modality encoders, thereby endowing LLMs with visual perception abilities (Alayrac et al., 2022; Li et al., 2023a). Single images or multiple frames are encoded as visual tokens and integrated with textual tokens to help MM-LLMs achieve visual understanding. For long-video understanding, MM-LLMs (Dai et al., 2023; Liu et al., 2024c) are designed to process a larger number of visual frames and diverse events, enabling a wide range of real-world applications such as automatically analyzing highlight reels from sports videos, movies, surveillance footage, and egocentric videos in embodied AI. For example, a robot could learn to make a cup

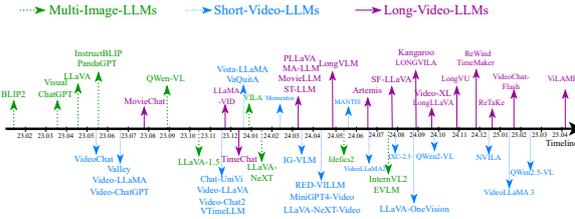


	Image-LLMs	Video-LLMs	Long-Video-LLMs
Task	<ul style="list-style-type: none"> <li>Image understanding: <ul style="list-style-type: none"> <li>Spatial reasoning: e.g. (Changpinyo et al., 2022; Chen et al., 2024a; Mathew et al., 2021; Peng et al., 2024; Sohoni et al., 2020; Wei et al., 2021).</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Short video understanding: <ul style="list-style-type: none"> <li>Spatial reasoning: e.g. (Li et al., 2023b; Ranasinghe et al., 2024).</li> <li>Within-event reasoning: e.g. (Diba et al., 2023; Huang et al., 2018).</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Long video understanding: <ul style="list-style-type: none"> <li>Spatial reasoning: e.g. (Fu et al., 2024a).</li> <li>Within-event reasoning: e.g. (Cheng et al., 2024).</li> <li><b>Between-event reasoning:</b> e.g. (Qian et al., 2024).</li> <li><b>Long-term reasoning:</b> e.g. (Wu et al., 2024).</li> </ul> </li> </ul>
Backbone	<ul style="list-style-type: none"> <li>Visual encoder: CLIP-ViT (Radford et al., 2021), SigLIP-ViT (Zhai301 et al., 2023), etc.</li> <li>LLM: LLaMA (Touvron et al., 2023b), etc.</li> </ul>	<ul style="list-style-type: none"> <li>Visual encoder: CLIP-ViT (Radford et al., 2021), SigLIP-ViT (Zhai301 et al., 2023), etc.</li> <li>LLM: LLaMA (Touvron et al., 2023b), etc.</li> </ul>	<ul style="list-style-type: none"> <li>Visual encoder: CLIP-ViT (Radford et al., 2021), SigLIP-ViT (Zhai301 et al., 2023), etc.</li> <li><b>Long-context LLM:</b> LLaMA3.1 (Dubey et al., 2024), etc.</li> </ul>
Connector	<ul style="list-style-type: none"> <li>Image-level connector: <ul style="list-style-type: none"> <li>Linear-layer-based: e.g. (Liu et al., 2024a; Liu et al., 2024c; Su et al., 2023)</li> <li>Pooling-based: e.g. (Liu et al., 2024b; Maaz et al., 2023; Xu et al., 2024a)</li> <li>Transformer-based: e.g. (Dai et al., 2023; Bai et al., 2023b; Jiang et al., 2024)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Image-level connector: <ul style="list-style-type: none"> <li>Image-Q-Former, Spatial-pooling, etc. e.g. (Liu et al., 2024a; Li et al., 2023b; Maaz et al., 2023; Li et al., 2024f)</li> </ul> </li> <li>Video-level connector <ul style="list-style-type: none"> <li>Video-Q-Former, Temporal-pooling, etc. e.g. (Zhang et al., 2023; Luo et al., 2023)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Image-level connector.</li> <li>Video-level connector.</li> <li><b>Long-video-level connector:</b> <ul style="list-style-type: none"> <li><b>Efficient token-compression:</b> e.g. (Song et al., 2024a; Xu et al., 2024a; Xu et al., 2024b)</li> <li><b>Time-aware design:</b> e.g. (Huang et al., 2024a; Ma et al., 2023b; Qian et al., 2024; Ren et al., 2024)</li> </ul> </li> </ul>
Training	<ul style="list-style-type: none"> <li>Pre-training: Image-text pairs. e.g. (Chen et al., 2015; Sharma et al., 2018; Chen et al., 2023b).</li> <li>Instruction-tuning: Image-language instruction data. e.g. (Chen et al., 2023b; Liu et al., 2024c)</li> </ul>	<ul style="list-style-type: none"> <li>Pre-training: Image-, Short-video-text pairs. e.g. (Chen et al., 2015; Sharma et al., 2018; Chen et al., 2023b; Bain et al., 2021).</li> <li>Instruction-tuning: Image-, short-video-language instruction data. e.g. (Maaz et al., 2023)</li> </ul>	<ul style="list-style-type: none"> <li>Pre-training: Image-, video-, <b>long-video-text</b> pairs. e.g. (Bain et al., 2021; Zhang et al., 2024d).</li> <li>Instruction-tuning: Image-, short-video-, <b>long-video-language instruction data.</b> e.g. (Li et al., 2023c; Huang et al., 2024a; Ren et al., 2024; Qian et al., 2024)</li> </ul>

Figure 2: The comparison of MM-LLMs among Image-, Short-Video-, and Long-Video-LLMs. The **bold content** often highlights special considerations of LV-LLMs for LVU.

We summarize the comparison of MM-LLMs among Image-, Short-Video-, and LV-LLMs in Fig. 2. LV-LLMs build upon advancements in multi-image and short-video MM-LLMs, sharing a similar structure of visual encoders, LLM backbones, and cross-modality connectors. To address the challenges in LVU, LV-LLMs incorporate more efficient long-video-level connectors that bridge cross-modal representations and compress visual tokens to a manageable number (Li et al., 2023c; Zhang et al., 2024d). Additionally, time-aware modules enhance the capture of temporal information (Qian et al., 2024). For pre-training and instruction-tuning, video-text pairs and video-instruction data are essential for MM-LLMs to handle both images and videos with shared spatial perception and reasoning capacity (Li et al., 2023b). Long video training datasets are particularly beneficial for temporal cross-modal semantic alignment and capturing long-term correlations, crucial for LV-LLMs (Song et al., 2024b). Our survey provides a comprehensive summary of recent advances in model design and training methods, tracing the evolution from images to long videos.

Recent surveys on visual understanding tasks typically adopt a single perspective, either from a global view of reviewing MM-LLMs (Yin et al., 2023; Zhang et al., 2024a) or from a local view focusing on image- or video-understanding tasks (Zhang et al., 2024b; Nguyen et al., 2024). While these works provide extensive reviews, they often lack discussing the developmental and inheritance relationships between different tasks and methods. Additionally, existing reviews on video understanding (Tang et al., 2023) focus more on general video understanding rather than the more challenging task of LVU. Long videos are prevalent in education, entertainment, and transportation, necessitating comprehensive automatic understanding with powerful

models (Apostolidis et al., 2021). Our work is among the earliest to summarize and discuss the LVU task from a developmental perspective.

Our survey is structured as follows: firstly, we find that the LVU task is more complex compared with image and short video understanding tasks (Sec.2.1), and summarize the unique challenges of LVU in Sec.2.2. Next, we provide a detailed summary of the developments in MM-LLMs from the perspectives of model architecture (Sec.3) and training methodologies (Sec.4), with an emphasis on the implementation of LV-LLMs for comprehensive LVU. We then compare the performance of video LLMs on LVU benchmarks (Sec.5), offering insights into the existing results of LV-LLMs. Finally, we discuss future research directions in LVU to advance the research field in Sec.6.

## 2 Long Video Understanding

In this section, we elaborate on visual understanding tasks among images, short-videos, and long-videos, and further analyze the challenges for LVU.

### 2.1 Visual Understanding

Visual understanding demands models to interpret visual information, integrating multimodal perception with commonsense reasoning (Johnson et al., 2017; Chen et al., 2024c).

**Image understanding.** As illustrated in Fig. 3 (a), image understanding tasks involve a single image for various visual reasoning tasks, such as image captioning and image-centered question answering (Sharma et al., 2018; Mathew et al., 2021; Changpinyo et al., 2022; Li et al., 2023a; Chen et al., 2024a). These tasks focus solely on spatial information, encompassing both coarse-grained understanding (Ordonez et al., 2011; Sohoni et al., 2020) of global visual context and fine-grained under-

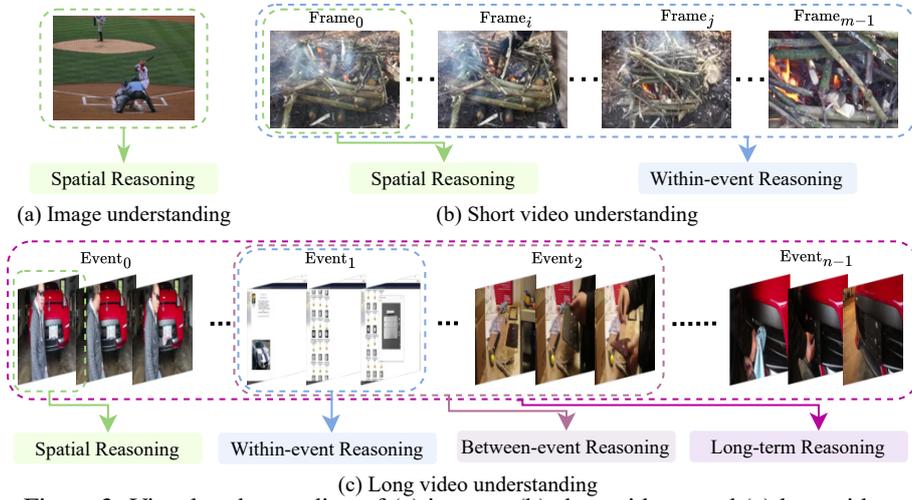


Figure 3: Visual understanding of (a) images, (b) short videos, and (c) long videos.

standing (Wei et al., 2021; Liu et al., 2024b; Peng et al., 2024) of local visual details.

**Short video understanding.** Unlike image understanding tasks, which involve only static visual data, short video understanding also incorporates temporal information from multiple visual frames (Xu et al., 2016; Bain et al., 2021; Li et al., 2023b, 2024e). In addition to spatial reasoning (Ranasinghe et al., 2024), within-event temporal reasoning and spatiotemporal reasoning across frames play crucial roles for short video understanding (Huang et al., 2018; Lin et al., 2019; Diba et al., 2023).

**Long video understanding (LVU).** Long videos typically consist of multiple events, encompassing much richer spatial content and temporal variations compared to short videos (Mangalam et al., 2024; Li et al., 2024f; Song et al., 2024a,b). As summarized in Fig. 3 (c), LVU involves not only spatial and within-event temporal reasoning but also between-event reasoning and long-term reasoning from different video events (Wu et al., 2019; Wu and Krahenbuhl, 2021; Wang et al., 2023a; Zhou et al., 2024; Fang et al., 2024).

## 2.2 Challenges of Long Video Understanding

Compared with images and short videos, long-form videos introduce new challenges to comprehensive visual understanding, as follows:

**Rich fine-grained spatiotemporal details.** Long videos, which cover a wide range of topics, scenes, and activities, contain varying details such as objects, events, and attributes (Fu et al., 2024a; Wu et al., 2024). These details are much richer compared to static images and short videos with multiple similar frames, making LVU more challenging. For instance, fine-grained spatial question answering can be introduced in any frame, while temporal question answering can be introduced between or among frames for long video reasoning tasks (Song

et al., 2024a). MM-LLMs for LVU must capture all relevant fine-grained spatiotemporal details from video frames spanning minutes or even hours, using a limited number of visual tokens.

**Dynamic events with scene transitions and content changes.** Long videos often contain various dynamic events with significant differences in scenes and content (Wu et al., 2024). These events can be semantically related and temporally coordinated according to their order of appearance (Bao et al., 2021), or they can exhibit significant semantic differences due to plot twists (Papalampidi et al., 2019). Between-event reasoning involving multiple events with diverse visual information is crucial for accurate content understanding (Cheng et al., 2024a; Qian et al., 2024). Distinguishing semantic differences and maintaining semantic coherence across varying events are essential for LVU.

**Long-term correlation and dependencies.** Long videos often contain actions and events that span extended periods. Capturing long-term dependencies and understanding how different parts of the video relate to each other over the long period is challenging (Wu et al., 2019). Video LLMs designed for images or short videos typically fail to contextualize the present event in relation to past or future events that are far from the current time (Wu and Krahenbuhl, 2021), as well as in long-term decision-making (Wang et al., 2024b).

## 3 Advances in Model Architecture

In this section, we discuss the advances of MM-LLMs from image-targeted to long-video-targeted models, from the perspective of model architecture.

### 3.1 Visual Encoder and LLM Backbone

MM-LLMs, encompassing both image-targeted and video-targeted models, typically utilize similar

Model	Year	Backbone		Connector			#Frame	#Token
		Visual Encoder	LLMs	Image-level	Video-level	Long-video-level		
InstructBLIP (2023)	23.05	EVA-CLIP-ViT-G/14	FlanT5, Vicuna-7B/13B	Q-Former	-	-	4	32/128
VideoChat (2023b)	23.05	EVA-CLIP-ViT-G/14	StableVicuna-13B	Q-Former	Global multi-head relation aggregator	-	8	/32
Video-LLaMA (2023)	23.06	EVA-CLIP-ViT-G/14	LLaMA, Vicuna	Q-Former	Q-Former	-	8	/32
Video-ChatGPT (2023)	23.06	CLIP-ViT-L/14	Vicuna1.1-7B	Spatial-pooling	Temporal-pooling	-	100	/356
Valley (2023)	23.06	CLIP-ViT-L/14	StableVicuna-7B/13B	-	Transformer and Mean pooling	-	0.5 fps	/256+T
MovieChat (2024a)	23.07	EVA-CLIP-ViT-G/14	LLaMA-7B	Q-Former	Frame mergein, Q-Former	Merging adjacent frames	2048	32/32
Qwen-VL (2023b)	23.08	Openclip-ViT-bigG	Qwen-7B	Cross-attention	-	-	4	/256
Chat-UniVi (2024)	23.11	CLIP-ViT-L/14	Vicuna1.5-7B	Token merging	-	-	64	/112
Video-LLaVA (2023)	23.11	LanguageBind-ViT-L/14	Vicuna1.5-7B	-	-	-	8	256/2048
LLaMA-VID (2023c)	23.11	CLIP-ViT-L/14	Vicuna-7B/13B	-	Context attention and pooling	-	1 fps	2/
VTimeLLM (2024a)	23.11	CLIP-ViT-L/14	Vicuna1.5-7B/13B	Frame feature	Q-Former	-	100	/100
VideoChat2 (2024e)	23.11	EVA-CLIP-ViT-G/14	Vicuna0-7B	-	Q-Former	-	16	/96
Vista-LLaMA (2023a)	23.12	EVA-CLIP-ViT-G/14	LLaVa-Vicuna-7B	Q-Former	Temporal Q-Former	-	16	32/512
TimeChat (2024a)	23.12	EVA-CLIP-ViT-G/14	LLaMA2-7B	Q-Former	Sliding window Q-Former	Time-aware encoding	96	/96
VaQuitA (2023b)	23.12	CLIP-ViT-L/14	LLaVA1.5-LLaMA-7B	-	Video Perceiver, VQ-Former	-	100	/356
Dolphins (2023b)	23.12	CLIP-ViT-L/14	OpenFlamingo	Perceiver Resampler, Gated cross-attention	-	Time embedding	-	-
Momentor (2024)	24.02	CLIP-ViT-L/14	LLaMA-7B	Frame feature, Temporal Perception Module, Grounded Event-Sequence Modeling	-	-	300	/1300
MovieLLM (2024b)	24.03	CLIP-ViT-L/14	Vicuna-7B/13B	-	Context attention and pooling	-	1 fps	2/
MA-LMM (2024)	24.04	EVA-CLIP-ViT-G/14	Vicuna-7B	Q-Former	Memory Bank Compression	Merging adjacent frames	100	/32
PLLaVA (2024a)	23.04	CLIP-ViT-L/14	LLaVA-Next-LLM	-	Adaptive Pooling	-	64	2304
LongVLM (2024)	23.04	CLIP-ViT-L/14	Vicuna1.1-7B	-	Hierarchical token merging	-	100	/305
MiniGPT4-Video (2024a)	24.04	EVA-CLIP-ViT-G/14	LLaMA2-7B, Mistral-7B	Merging adjacent tokens	-	-	90	64/5760
RED-VILLM (2024b)	24.04	Openclip-ViT-bigG	Qwen-7B	Spatial pooling	Temporal pooling	-	100	/1124
ST-LLM (2024c)	24.04	BLIP-2	InstructBLIP-Vicuna1.1-7B	Q-Former	Masked video modeling	Global-Local input	16	/512
LLaVA-NeXT-Video (2024e)	24.04	CLIP-ViT-L/14	Vicuna1.5-7B/13B	Merging adjacent tokens	-	-	32	4608
Mantis-Idetics2 (2024)	24.05	SigLIP-SO400M	Mistral0.1-7B	Perceiver resampler	-	-	8	64/512
VideoLLaMA 2 (2024b)	24.06	CLIP-ViT-L/14	Mistral-7B-Instruct	-	Spatial-Temporal Convolution	-	8	/576
LongVA (2024d)	24.06	CLIP-ViT-L/14	Qwen2-7B-224K	Merging adjacent tokens	Expanding tokens	-	384	55,296
Artemis (2024)	24.06	CLIP-ViT-L/14	Vicuna1.5-7B	-	Average pooling	-	5	/356
VideoGPT+ (2024)	24.06	CLIP-ViT-L/14	Phi3-Mini-3.8B	Adaptive pooling	Adaptive pooling	-	16	/2560
IXC-2.5 (2024c)	24.07	CLIP-ViT-L/14-490	InternLM2-7B	Merging adjacent tokens	Expanding tokens	Frame index	64	400/25600
EVLm (2024b)	24.07	EVA2-CLIP-E-Plus	Qwen-14B-Chat 1.0	Gated cross attention	-	-	-	/16
StowFast-LLaVA (2024b)	24.07	CLIP-ViT-L/14	Vicuna1.5-7B	Merging adjacent tokens	-	Slow and fast pathway	50	3680
LLaVA-Interleave (2024d)	24.07	SigLIP-SO400M	Qwen1.5-0.5B/7B/14B	-	-	-	16	729/11664
Kangaroo (2024d)	24.08	EVA-CLIP-ViT-G/14	LLaMA3-8B	-	3D Depthwise convolution	-	-	-
VITA (2024b)	24.08	InternViT-300M-448px	Mixtral 8x7B	MLP	-	-	16	256/4096
LLaVA-OneVision (2024c)	24.08	SigLIP-SO400M	Qwen2-7B	Merging adjacent tokens	-	-	1 fps	729/
LONGVILA (2024)	24.08	SigLIP-SO400M	Qwen2-1.5B/7B	-	Multi-Modal Sequence Parallelism	-	1024	256/
LongLLaVA (2024e)	24.09	CLIP-ViT-B/32	LLaVA1.6-13B	Merging adjacent tokens	Mamba Layers	Hybrid architecture	256	144/
Qwen2-VL (2024a)	24.09	CLIP-ViT-L/14	Qwen2-1.5B/7B/72B	Merging adjacent tokens	3D convolutions	-	2 fps	66/
Video-XL (2024)	20.09	CLIP-ViT-L/14	Qwen-2-7B	Merging adjacent tokens	Visual Summarization Token and Dynamic Compression	-	128	-
Oryx-1.5 (2024g)	24.10	OryxViT	Qwen-2.5-7B/32B	Variable-Length Self-Attention	Dynamic Compressor	-	64	256/
LongVU (2024)	24.10	SigLIP-SO400M	Qwen2-7B, LLaMA3.2-3B	Selective frame feature reduction	Frame selection and Token Reduction	-	1 fps	-
TimeMarker (2024d)	24.11	LLaVA-Encoder	LLaVA-LLM	Adaptive Token Merge and Temporal Separator Tokens Integration	-	-	128	-
ReWind (2024)	24.11	EVA-CLIP-ViT-G/14	LLaMA2-7B	Pooling, Learnable queries, Cross-attentions and Frame Selection	-	-	64	32/256
NVILA (2024f)	24.12	SigLIP-SO400M	Qwen2-7B/14B	Spatial-to-Channel Reshaping	Temporal Averaging	-	256	/8192
IQViC (2024)	24.12	CLIP-ViT-L/14	Vicuna-v1-7B	Visual Compressor	Temporal Compressor	-	-	/640
ReTaKe (2024c)	24.12	CLIP-ViT-L/14	Qwen2-7B	-	Key frame selection, KV-Cache compression	-	1024	-
InternVL2.5 (2024e)	24.12	InternViT-300M	Qwen2.5	Dynamic High Resolution	-	-	64	356/16384
VideoLLaMA 3 (2025a)	25.01	InternViT-6B	InternLM2.5	Pixel Unshuffle	-	-	180	/10240
LLaVA-Mini (2025c)	25.01	SigLIP-SO400M	Qwen-2.5-7B	Any-resolution Vision Tokenization	Differential Frame Pruner	-	10000	1/10000
VideoChat-Flash (2025b)	25.01	UMT-L@24	LLaMA-3.1-8B-Instruct	Vision Token Compression	-	-	1000	16/
VideoLLaMA 3 (2025a)	25.01	SigLIP-SO400M	Qwen2.5-2B/7B	-	Hierarchical Visual Token Compression	-	180	/16K
LLaVA-ST (2025a)	25.01	SigLIP-SO400M	LlaVA-OneVision-7B	Dynamic Vision Tokens, Video Compressor, Spatial Downsampling	-	-	-	-
Qwen2.5-VL (2025)	25.02	Qwen2.5-VL-ViT	Qwen-2.5	Language-Aligned Positional Embedding	-	-	768	/24576
Vision-R1 (2025)	25.03	Qwen-2.5-VL-7B-Instruct	-	Dynamic FPS Sampling, Absolute Time Encoding	-	-	-	/16K
R1-VL (2025b)	25.03	Qwen2-VL-2B/7B	-	Progressive Thinking Suppression Training, Group Relative Policy Optimization	-	-	-	-
ViLAMP (2025)	25.04	LLaVA-OneVision-7B	-	Step-wise Group Relative Policy Optimization	-	-	-	-
					Differential Keyframe Selection, Differential Feature Merging	-	128	-

Table 1: Comparison of mainstream Video-LLMs across model design choices. The notation A/B in the last column indicates A tokens per frame and a total of B tokens for the entire video.

visual encoders for visual information extraction. LLM backbones are also universal in early MM-LLM methods, while existing LV-LLMs tend to use long-context LLMs in the implementation.

**Visual encoder.** The pretrained visual encoders are responsible for capturing vision knowledge from raw visual data. As summarized in Table 1, image encoders like CLIP-ViT-L/14 (Radford et al., 2021), EVA-CLIP-ViT-G/14 (Sun et al., 2023) and SigLIP-SO400M (Zhai et al., 2023) are widely utilized as visual modality encoders in image- and video-targeted LLMs. Recent work (Li et al., 2024a) shows that the visual representation, including image resolution, the size of visual token, and the pre-training visual resources, play a more important role than the size of the visual encoder.

**LLM backbone.** The LLM is the core module in visual understanding systems, inheriting properties of reasoning and decision-making.

The strength of the LLM typically correlates with superior multimodal capabilities in visual LLMs (Li et al., 2024b,a). For LLMs of equivalent scale, those with superior language capabilities demonstrate enhanced performance, whereas for the same LLMs with varying model sizes, larger models generally achieve better multimodal performance. Additionally, long-context LLMs that extend the context length to hundreds of thousands of tokens support learning with more extensive data (Yang et al., 2024). Recent LV-LLMs effectively transfer the LLM’s long-context understanding ability to the vision modality (Zhang et al., 2024d).

### 3.2 Modality Interface

Connectors between visual encoders and LLMs act as modality interfaces, mapping visual features to the language space. Due to the variability in visual data sources, these connectors are categorized

259 into image-level, video-level, and long-video-level  
260 types. (More image- and short-video-level connec-  
261 tor designs are summarized in Appendix B.3.)

262 **Image-level connectors.** Image-level connectors  
263 aim to map image features to the language space for  
264 processing raw visual tokens and are widely used in  
265 both image- and video-targeted MM-LLMs. These  
266 connectors fall into three categories: (1) single  
267 linear layers (Liu et al., 2024c) or multi-layer per-  
268 ceptors (MLPs) (Liu et al., 2024a) for embedding,  
269 (2) pooling-based methods, and (3) cross-attention  
270 or transformer-based structures like Q-Former (Li  
271 et al., 2023a) and Perceiver Resampler (Jaegle et al.,  
272 2021) for feature compression.

273 **Video-level connectors.** Video-level connectors  
274 are used for extracting sequential visual data and  
275 further compressing visual features. Compared to  
276 the solely image-level connectors in image-targeted  
277 MM-LLMs, video-level connectors are essential  
278 for video-targeted MM-LLMs, including LV-LLMs.  
279 Some methods directly concatenate image tokens  
280 before inputting them to the LLMs, making them  
281 sensitive to the number of frame images (Dai et al.,  
282 2023; Lin et al., 2023). Similar structures used  
283 for token compression in image-level connectors  
284 can be adapted for video-level interfaces, such  
285 as pooling-based and transformer-based structures  
286 (Maaz et al., 2023; Song et al., 2024a; Zhang et al.,  
287 2023; Ma et al., 2023a; Ren et al., 2024a).

288 **Long-video-level connectors.** Long-video-level  
289 connectors focus more on efficient visual data com-  
290 pression and long-term information preserving.  
291 Efficiently compressing visual information requires  
292 not only reducing the input visual tokens to an ac-  
293 ceptable quantity but also preserving the complete  
294 spatiotemporal details contained in long videos.  
295 Videos contain two types of data redundancy: spa-  
296 tial data redundancy within frames and spatiotem-  
297 poral data redundancy across frames (Li et al.,  
298 2022; Chen et al., 2023a; Cheng et al., 2024c). On  
299 the one hand, spatial data redundancy arises when  
300 region-level pixels within frames are the same, lead-  
301 ing to inefficiencies when representing the redun-  
302 dant visual frame through full visual tokens. To  
303 reduce spatial video data redundancy, the LLaVA-  
304 Next-series methods (Zhang et al., 2024e; Li et al.,  
305 2024d; Liu et al., 2024b; Li et al., 2024c) merge  
306 adjacent frame patch tokens, and Chat-UniVi (Jin  
307 et al., 2024) merges similar frame patch tokens.  
308 On the other hand, spatiotemporal data redundancy  
309 includes both cross-frame pixel redundancy and  
310 motion redundancy (Pourreza et al., 2023), where

311 the semantic information is similar among these  
312 redundant video frames. To reduce spatiotemporal  
313 video redundancy, MovieChat (Song et al., 2024a)  
314 and MA-LMM (He et al., 2024) merge frame fea-  
315 tures with higher frame similarity before inputting  
316 them to LLMs. In addition to reducing redundant  
317 information, preserving more video spatiotem-  
318 poral details is crucial for accurate long video rea-  
319 soning (Diba et al., 2023). To balance global and  
320 local visual information and support more frame in-  
321 puts, SlowFast-LLaVA (Xu et al., 2024b) employs  
322 a slow pathway to extract features at a low frame  
323 rate while retaining more visual tokens, and a fast  
324 pathway at a high frame rate with a larger spatial  
325 pooling stride to focus on motion cues.

326 Additionally, time-involved visual data efficiently  
327 manage the temporal and spatial information inher-  
328 ent in long-form videos (Hou et al., 2024).  
329 The time-aware design can enhance the temporal-  
330 capturing capability of video-related LLMs,  
331 which is particularly beneficial for LVU. Both  
332 VTimeLLM (Huang et al., 2024a) and InternLM-  
333 XComposer-2.5 (IXC-2.5) (Zhang et al., 2024c)  
334 use frame indices to enhance temporal relations.  
335 The difference lies in their approach: VTimeLLM  
336 learns temporal information by training with de-  
337 coded text that includes frame indices, while IXC-  
338 2.5 encodes frame indices along with the frame  
339 image context. TimeChat (Ren et al., 2024a) and  
340 Momentor (Qian et al., 2024) inject temporal in-  
341 formation directly into frame features for fine-  
342 grained temporal information capture. Specifically,  
343 TimeChat designs a Time-aware Frame Encoder to  
344 extract visual features with corresponding times-  
345 tamp descriptions at the frame level, while Mo-  
346 mentor utilizes a Temporal Perception Module for  
347 continuous time encoding and decoding, injecting  
348 temporal information into frame features.

349 **Retrieval-based LVU.** A significant proportion  
350 of LVU methods are retrieval-based, addressing  
351 challenges like "noise and redundancy" and "mem-  
352 ory and computation" constraints. R-VLM selects  
353 the most relevant video chunks for question an-  
354 swering (Xu et al., 2023), Goldfish retrieves top-  
355 clips to focus on pertinent segments (Ataallah et al.,  
356 2024b), and DrVideo transforms videos into text  
357 documents to retrieve key frames (Ma et al., 2024).  
358 Video-RAG uses visually-aligned auxiliary texts  
359 for cross-modality alignment (Luo et al., 2024),  
360 while VideoLLaMB employs temporal memory to-  
361kens and a SceneTiling algorithm to preserve se-  
362 mantic continuity (Wang et al., 2024f). VideoAgent



Figure 4: Long video sample for pretraining and instruction-tuning.

(Wang et al., 2024d) leverages a LLM to iteratively compile critical information, using vision-language models to enhance LVU.

## 4 Advances in Model Training

Multimodal LLMs for visual understanding consist of two principal stages: pre-training (PT) for vision-language feature alignment and instruction-tuning (IT) for reasoning response (seen in Appendix B.4).

### 4.1 Pre-training

Vision-language pre-training for MM-LLMs aims to align visual features with the language space using text-paired data. This includes pre-training with image-, short-video-, and long-video-text datasets. Initially introduced for visual LLMs focused on images, **image-text pre-training** is also widely used in video-related understanding tasks. Coarse-grained image-text pair datasets, such as COCO Captions (Chen et al., 2015) and CC-3M (Sharma et al., 2018), are employed for global vision-language alignment. Fine-grained image-text datasets like ShareGPT4V-PT (Chen et al., 2023b) are used for locally spatial semantics alignment. Given the limited changes in semantic content of short videos, short-video-text paired datasets, such as Webvid-2M (Bain et al., 2021), can be used similarly for **short-video-text pre-training**. Similarly, **long-video-text pre-training** is important to capture the temporal semantic alignment of long videos for LVU (Ju et al., 2025). Given the absence of long-term cross-modal correlation in image-text and short-video-text pairs, long-video-text pre-training datasets with pairs of long videos and their corresponding text descriptions are necessary (Argaw et al., 2023). Moreover, as shown in Fig. 4

(a), the scenes and events in long videos vary significantly across frames, necessitating event-level vision-language alignment (Qian et al., 2024) for long-video-text pre-training, which is markedly different from both image-text and short-video-text pre-training (Zhang et al., 2024d).

### 4.2 Instruction-tuning

Instruction-tuning with vision-language sources enables LLMs to follow instructions and generate human-like text. Multimodal vision-language instruction-following data (Dai et al., 2023; Liu et al., 2024c), including both image-text and video-text pairs, are used to align multimodal LLMs with human intent, thereby enhancing their ability to complete real-world tasks.

Similar to the pre-training stage, **image-text instruction-tuning** is employed in various vision-understanding tasks, including image, short-video, and long-video understanding. Basic image-based instruction-following datasets, such as ShareGPT4V-Instruct (Chen et al., 2023b) and LLaVA-Instruct (Liu et al., 2024c), provide high-quality instruction-tuning data for spatial reasoning and chat capabilities. For video-related LLMs, **short-video-text instruction-tuning** is necessary to enable multimodal LLMs to understand temporal sequences, as seen in models like VideoChatGPT (Maaz et al., 2023) and VideoChat (Li et al., 2023b). Short-video-LLMs require both spatial and within-event reasoning instructions to understand the spatial and small-scale temporal content of short videos. However, the limited content and semantic changes in short videos are insufficient for LVU tasks, where frames are more numerous and exhibit significant variation. **Long-video-text instruction-tuning** is specifically introduced to bet-

Model	LLM	Long	VideoVista (131s)	MMBench-Video (165s)	EgoSchema (180s)	LongVideoBench (473s)	MLVU (12mins)	Video-MME (1024s)	LVBench (4101s)
Momentor	LLaMA-7B	✗	–	–	–	–	–	–	–
TimeChat	LLaMA2-7B	✓	–	–	33.0 <sup>‡</sup>	–	30.9 <sup>†</sup>	–	22.3 <sup>♣</sup>
LLaMA-VID	Vicuna-7B	✓	56.87 <sup>‡</sup>	–	38.5 <sup>†</sup>	–	33.2 <sup>†</sup>	–	23.9 <sup>♣</sup>
LLaVA-NeXT-Video	Vicuna1.5-7B	✗	56.66 <sup>‡</sup>	–	43.9 <sup>†</sup>	43.5 <sup>♡</sup>	–	–	–
VideoLLaMA 2 (16)	Mistral-7B-Instruct	✗	60.47 <sup>‡</sup>	–	51.7	–	48.5 <sup>†</sup>	47.9/50.3 <sup>♣</sup>	–
PLLaVA	LLaVA-Next-7B	✓	60.36 <sup>‡</sup>	1.03 <sup>†</sup>	54.4 <sup>†</sup>	39.2 <sup>♡</sup>	–	–	–
LongVA	Qwen2-7B-224K	✓	67.36 <sup>‡</sup>	–	–	–	56.3 <sup>†</sup>	52.6/54.3 <sup>♣</sup>	–
IXC-2.5-7B	InternLM2-7B	✗	68.91 <sup>‡</sup>	1.41	–	–	58.8	55.8/-	–
Kangaroo	LLaMA3-8B	✓	69.50 <sup>‡</sup>	1.44	62.7	54.8	61.0	56.0/57.6 <sup>♣</sup>	39.4
Video-XL	QWen2-7B	✓	70.60	–	–	50.7	–	55.5/61.0	–
TimeMarker	LLaVA-7B	✓	78.40	1.53	–	56.3	–	57.3/62.8	41.3
ReTaKe	Qwen2-7B	✓	–	–	–	–	69.8	63.9/68.9	47.8
VideoLLaMA 3	Qwen-2.5-7B	✓	–	–	63.3	59.8	73.0	66.2/70.3	45.3
VideoChat-Flash	Qwen2-7B	✓	–	–	–	64.2	74.5	64.0/69.4	48.2

Table 2: Comparison of Long-Video-LLMs on LVU benchmarks. Results with <sup>‡</sup> are from the VideoVista benchmark (Li et al., 2024f). Results with <sup>†</sup> are from the Kangaroo (Liu et al., 2024d). Results with <sup>♣</sup> are from Video-MME benchmark (Fu et al., 2024a). Results with <sup>♠</sup> are from LVBench (Wang et al., 2024b). Results with <sup>♡</sup> are from LongVideoBench (Wu et al., 2024).

ter capture and understand long videos. In addition to spatial and within-event reasoning instructions, between-event and long-term reasoning instructions (Ren et al., 2024b; Zeng et al., 2024) are necessary for comprehensive understanding, as shown in Fig. 4 (b). Among the newly introduced long-video instruction-format datasets, Long-VideoQA (Li et al., 2023c), Video-ChatGPT (Maaz et al., 2023) and LongViTU (Wu et al., 2025) are not time-aware. In contrast, VTimeLLM (Huang et al., 2024a), TimeIT (Ren et al., 2024a), and Moment-10M (Qian et al., 2024) are time-aware, incorporating temporal information to enhance reasoning.

## 5 Evaluation, Performance and Analysis

This section presents a performance comparison across popular evaluation datasets with videos of varying lengths, along with our analysis. Additional comparisons are provided in Appendix C.

To address the unique characteristics of long videos, several long video benchmarks have been introduced in recent years, with video lengths varying from hundreds of seconds to thousands of seconds. EgoSchema (Mangalam et al., 2024) is long-form video understanding datasets designed for multiple-choice question answering, after accessing all frames. VideoVista (Li et al., 2024f), MMBench-Video (Fang et al., 2024), and MLVU (Zhou et al., 2024) cover various topics and are designed for fine-grained capability evaluation. LongVideoBench (Wu et al., 2024) introduces referring reasoning questions to address the longstanding issue of single-frame bias in long videos. Video-MME (Fu et al., 2024a) and LVBench (Wang et al., 2024b) contain numerous hour-level videos. Video-MME further categorizes them into short, medium, and long categories, while LVBench aims to challenge models to demonstrate long-term memory and extended comprehension capabilities.

We further compare and analyze the perfor-

mance of LVU, specifically summarizing their performance on long video benchmarks with lengths varying from hundreds of seconds to thousands of seconds. As shown in Table 2, LVU-specific methods typically outperform short video understanding methods on LVU tasks. This indicates that specially designed, powerful video-level connectors are essential for LVU. Additionally, the performance on benchmarks with longer video lengths is generally worse than on those with shorter lengths. For example, the performance of methods across VideoVista and MLVU, Video-MME and LVBench, using the same evaluation metric, shows a decline as video length increases. This suggests that LVU remains a challenging research topic.

## 6 Future Directions

To meet the demands of an AI-driven society with increasingly longer multimodal data, developing more powerful visual LLMs for LVU is crucial.

### 6.1 More Long Video Training Resources

The two-stage training pipeline, consisting of cross-modal alignment pre-training and visual-instruction tuning, is widely employed for training MM-LLMs (Dai et al., 2023; Liu et al., 2024c). However, there are several challenges for LVU:

- **Hour-long video datasets.** The length of newly introduced long-video training data is limited to minutes, restricting effective reasoning for hour-long LVU (Li et al., 2023c).
- **Necessity of long video pre-training.** Fine-grained long-video-language training pairs are lacking compared to image- and short-video-language pairs during pre-training (Song et al., 2024b; Qiu et al., 2024). Exploring the necessity of long-video-language paired datasets is crucial for evaluating the value of capturing long-term correlations in long videos (Zhang et al., 2024d).
- **Large-scale long video instruction-tuning**

**datasets.** Existing long video datasets, mentioned in Sec 4.2 are limited in size. Creating large-scale long-video-instruction datasets is essential for comprehensive long-video understanding.

## 6.2 More Challenging LVU Benchmarks

Recent video understanding benchmarks, such as LongVideoBench (Wu et al., 2024), VideoVista (Li et al., 2024f), and MLVU (Zhou et al., 2024), focus on specific aspects of LVU like long-context interleaved and fine-grained video understanding. However, comprehensive benchmarks that cover frame-level and segment-level reasoning with time and language are necessary but currently unexplored for a thorough evaluation of general LVU methods (Wu et al., 2024). Existing benchmarks, typically at the minute level, fail to adequately test long-term capabilities. LVU methods often suffer from catastrophic forgetting and loss of spatiotemporal details when reasoning with extensive sequential visual information (Wang et al., 2024b), such as hour-level videos. Additionally, most LVU benchmarks focus solely on the visual modality. Incorporating multi-modal data, including audio and language, would significantly benefit LVU tasks.

## 6.3 Powerful and Efficient Frameworks

Visual LLMs for videos need to support more visual frames and preserve more visual details with a fixed number of visual tokens. There are four main considerations when implementing LV-LLMs:

- **Select long-context LLMs as the LLM backbones.** Previous methods have suffered from limited context capacity and required specific fine-tuning to support more tokens (Zhang et al., 2024d). Recent long-context LLMs, such as QWen2 (Yang et al., 2024) and LLaMA-3.1 (Dubey et al., 2024), offer a context window length of 128K and can be utilized in LV-LLM without extensive fine-tuning.
- **Compress visual tokens efficiently with minimal information loss.** Existing methods face issues with insufficient or excessive compression. For example, Chat-UniVi (Jin et al., 2024) uses multi-scale token merging, and LongVA merges adjacent tokens only, while LLaMA-VID (Li et al., 2023c) and MA-LMM (He et al., 2024) compress too much visual information, leading to significant loss of frame details. New frameworks must efficiently compress visual tokens to support more temporal frames and preserve spatiotemporal details. At the image level, adjacent

frames can be merged or represented with fewer visual tokens due to their similarity and redundancy (Kim et al., 2024; Xu et al., 2024b). At the video level, relatively independent video events can be compressed into single visual units with corresponding visual tokens, allowing the inputs to cover long-form visual content effectively. Additionally, retrieval-based methods address challenges like "noise and redundancy" and "memory and computation" constraints by leveraging an LLM to iteratively compile critical information (Xu et al., 2023; Ataallah et al., 2024b).

- **Incorporate time-aware designs.** Enhance video reasoning by incorporating temporal information, as seen in designs like TimeIT (Ren et al., 2024a) and Moment-10M (Qian et al., 2024), to improve temporal information extraction in LVU tasks. Temporal information can be injected at various levels: token level, image level, or event level, significantly enhancing the model's ability to understand and reason about long videos.
- **Utilize infrastructure for memory-intensive training.** To handle the increased data load, it is essential to have infrastructure that supports memory-intensive long-context training. Employ infrastructure capable of supporting long-context training with a large number of GPU devices, as demonstrated by LongViLa (Xue et al., 2024), ensuring efficient training on long-form content.
- **Reasoning for long video understanding.** Recent studies improve MM-LLMs' reasoning via supervised fine-tuning but often fail to truly understand reasoning errors in complex tasks such as long video understanding; inspired by DeepSeek-R1 (Guo et al., 2025), a systematic exploration of rule-based reinforcement learning to enhance reasoning in long video understanding with MM-LLMs is therefore necessary.

## 7 Conclusion

This paper reviews the progress of visual LLMs from images to long videos, highlighting key challenges in long video understanding (LVU) such as capturing fine-grained spatiotemporal details and long-term dependencies amid dynamic scene changes. We discuss advances in model architectures and training strategies for Long-Video-LLMs, survey various video benchmarks, and compare method performances. As the first work focusing on Long-Video-LLMs for LVU, we aim to provide insights and drive future research in this area.

## 610 Limitation

611 We reviewed literature on comprehensive long  
612 video understanding, covering methods, training  
613 datasets, and benchmarks. Due to space constraints,  
614 we omit detailed application scenarios like real-  
615 time processing and multimodal tasks. We will  
616 maintain an open-source repository and add these  
617 contents to complement our survey. The perfor-  
618 mance comparisons are based on final results from  
619 previous papers and official benchmarks, which  
620 vary in training resources, strategies, and model  
621 architectures, making it difficult to analyze spe-  
622 cific models and training differences. We plan to  
623 conduct detailed ablation studies on public bench-  
624 marks for a more direct analysis of model design,  
625 training resources, and methods.

## 626 References

627 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
628 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
629 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
630 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.  
631 *arXiv preprint arXiv:2303.08774*.

632 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,  
633 Antoine Miech, Iain Barr, Yana Hasson, Karel  
634 Lenc, Arthur Mensch, Katherine Millican, Malcolm  
635 Reynolds, et al. 2022. Flamingo: a visual language  
636 model for few-shot learning. *Advances in neural  
637 information processing systems*, 35:23716–23736.

638 Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-  
639 garet Mitchell, Dhruv Batra, C Lawrence Zitnick, and  
640 Devi Parikh. 2015. Vqa: Visual question answering.  
641 In *Proceedings of the IEEE international conference  
642 on computer vision*, pages 2425–2433.

643 Evlampios Apostolidis, Eleni Adamantidou, Alexan-  
644 dros I Metsai, Vasileios Mezaris, and Ioannis Pa-  
645 tras. 2021. Video summarization using deep neu-  
646 ral networks: A survey. *Proceedings of the IEEE*,  
647 109(11):1838–1863.

648 Dawit Mureja Argaw, Joon-Young Lee, Markus Wood-  
649 son, In So Kweon, and Fabian Caba Heilbron. 2023.  
650 Long-range multimodal pretraining for movie un-  
651 derstanding. In *Proceedings of the IEEE/CVF In-  
652 ternational Conference on Computer Vision*, pages  
653 13392–13403.

654 Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman,  
655 Essam Sleiman, Deyao Zhu, Jian Ding, and Mo-  
656 hamed Elhoseiny. 2024a. Minigt4-video: Advanc-  
657 ing multimodal llms for video understanding with  
658 interleaved visual-textual tokens. *arXiv preprint  
659 arXiv:2404.03413*.

660 Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrah-  
661 man, Essam Sleiman, Mingchen Zhuge, Jian Ding,

Deyao Zhu, Jürgen Schmidhuber, and Mohamed El-  
hoseiny. 2024b. Goldfish: Vision-language under-  
standing of arbitrarily long videos. *arXiv preprint  
arXiv:2407.12679*. 662  
663  
664  
665

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hes-  
sel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,  
Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al.  
2023. Openflamingo: An open-source framework for  
training large autoregressive vision-language models.  
*arXiv preprint arXiv:2308.01390*. 666  
667  
668  
669  
670  
671

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
Huang, et al. 2023a. Qwen technical report. *arXiv  
preprint arXiv:2309.16609*. 672  
673  
674  
675

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,  
Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,  
and Jingren Zhou. 2023b. Qwen-vl: A frontier large  
vision-language model with versatile abilities. *arXiv  
preprint arXiv:2308.12966*. 676  
677  
678  
679  
680

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
bin Ge, Sibio Song, Kai Dang, Peng Wang, Shijie  
Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical  
report. *arXiv preprint arXiv:2502.13923*. 681  
682  
683  
684

Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zis-  
serman. 2021. Frozen in time: A joint video and  
image encoder for end-to-end retrieval. In *Proceed-  
ings of the IEEE/CVF international conference on  
computer vision*, pages 1728–1738. 685  
686  
687  
688  
689

Peijun Bao, Qian Zheng, and Yadong Mu. 2021. Dense  
events grounding in video. In *Proceedings of  
the AAAI Conference on Artificial Intelligence*, vol-  
ume 35, pages 920–928. 690  
691  
692  
693

Tom B Brown. 2020. Language models are few-shot  
learners. *arXiv preprint ArXiv:2005.14165*. 694  
695

Fabian Caba Heilbron, Victor Escorcia, Bernard  
Ghanem, and Juan Carlos Nieves. 2015. Activitynet:  
A large-scale video benchmark for human activity  
understanding. In *Proceedings of the ieee conference  
on computer vision and pattern recognition*, pages  
961–970. 696  
697  
698  
699  
700  
701

Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen,  
Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi  
Chen, Pei Chu, et al. 2024. Internlm2 technical re-  
port. *arXiv preprint arXiv:2403.17297*. 702  
703  
704  
705

Soravit Changpinyo, Doron Kukliansky, Idan Szpektor,  
Xi Chen, Nan Ding, and Radu Soricut. 2022. All you  
may need for vqa are image captions. *arXiv preprint  
arXiv:2205.01883*. 706  
707  
708  
709

Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter,  
Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024a.  
Spatialvlm: Endowing vision-language models with  
spatial reasoning capabilities. In *Proceedings of the  
IEEE/CVF Conference on Computer Vision and Pat-  
tern Recognition*, pages 14455–14465. 710  
711  
712  
713  
714  
715

716	Jou-An Chen, Wei Niu, Bin Ren, Yanzhi Wang, and Xipeng Shen. 2023a. Survey: Exploiting data redundancy for optimization of deep learning. <i>ACM Computing Surveys</i> , 55(10):1–38.	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. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.	770 771 772 773 774 775
720	Kaibing Chen, Dong Shen, Hanwen Zhong, Huasong Zhong, Kui Xia, Di Xu, Wei Yuan, Yifei Hu, Bin Wen, Tianke Zhang, et al. 2024b. Evlm: An efficient vision-language model for visual understanding. <i>arXiv preprint arXiv:2407.14177</i> .	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. <i>Journal of Machine Learning Research</i> , 24(240):1–113.	776 777 778 779 780 781
725	Liangyu Chen, Bo Li, Sheng Shen, Jingkang Yang, Chunyuan Li, Kurt Keutzer, Trevor Darrell, and Ziwei Liu. 2024c. Large language models are visual reasoning coordinators. <i>Advances in Neural Information Processing Systems</i> , 36.	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. <i>Journal of Machine Learning Research</i> , 25(70):1–53.	782 783 784 785 786
730	Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023b. Sharegpt4v: Improving large multimodal models with better captions. <i>arXiv preprint arXiv:2311.12793</i> .	Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tjong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	787 788 789 790 791 792
735	Shimin Chen, Xiaohan Lan, Yitian Yuan, Zequn Jie, and Lin Ma. 2024d. Timemarker: A versatile video-llm for long and short video understanding with superior temporal localization ability. <i>arXiv preprint arXiv:2411.18211</i> .	Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 326–335.	793 794 795 796 797
740	Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. <i>arXiv preprint arXiv:1504.00325</i> .	Ali Diba, Vivek Sharma, Mohammad Arzani, Luc Van Gool, et al. 2023. Spatio-temporal convolution-attention video network. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 859–869.	798 799 800 801 802
745	Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024e. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. <i>arXiv preprint arXiv:2412.05271</i> .	Anxhelo Diko, Tinghuai Wang, Wassim Swaileh, Shiyang Sun, and Ioannis Patras. 2024. Rewind: Understanding long videos with instructed learnable memory. <i>arXiv preprint arXiv:2411.15556</i> .	803 804 805 806
751	Chuanqi Cheng, Jian Guan, Wei Wu, and Rui Yan. 2025. Scaling video-language models to 10k frames via hierarchical differential distillation. <i>arXiv preprint arXiv:2504.02438</i> .	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	807 808 809 810 811
755	Dingxin Cheng, Mingda Li, Jingyu Liu, Yongxin Guo, Bin Jiang, Qingbin Liu, Xi Chen, and Bo Zhao. 2024a. Enhancing long video understanding via hierarchical event-based memory. <i>arXiv preprint arXiv:2409.06299</i> .	Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao, Yining Li, Dahua Lin, and Kai Chen. 2024. Mmbench-video: A long-form multi-shot benchmark for holistic video understanding. <i>arXiv preprint arXiv:2406.14515</i> .	812 813 814 815 816
760	Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. 2024b. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. <i>arXiv preprint arXiv:2406.07476</i> .	Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024a. Videomme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. <i>arXiv preprint arXiv:2405.21075</i> .	817 818 819 820 821 822
766	Zheng Cheng, Rendong Wang, and Zhicheng Wang. 2024c. Focuschat: Text-guided long video understanding via spatiotemporal information filtering. <i>arXiv preprint arXiv:2412.12833</i> .	Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, et al. 2024b. Vita: Towards	823 824 825

826	open-source interactive omni multimodal llm. <i>arXiv preprint arXiv:2408.05211</i> .	Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhu Chen. 2024. Mantis: Interleaved multi-image instruction tuning. <i>arXiv preprint arXiv:2405.01483</i> .	881
827			882
828	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .		883
829			884
830		Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. 2024. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13700–13710.	885
831			886
832			887
833	Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. 2024. Ma-Imm: Memory-augmented large multimodal model for long-term video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13504–13514.		888
834			889
835		Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2901–2910.	890
836			891
837	Jianlong Hou, Tao Tao, Jibin Wang, Zhuo Chen, Xuelian Ding, and Kai Wang. 2024. Memotichat: A memory-augmented time-sensitive model for ultra-long video understanding. In <i>2024 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–9. IEEE.		892
840			893
841			894
842			895
843			896
844			897
845	Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. 2024a. Vtimellm: Empower llm to grasp video moments. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14271–14280.	Xuan Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu, and Ying Shan. 2025. Miradata: A large-scale video dataset with long durations and structured captions. <i>Advances in Neural Information Processing Systems</i> , 37:48955–48970.	898
846			899
847			900
848			901
849			902
850	De-An Huang, Vignesh Ramanathan, Dhruv Mahajan, Lorenzo Torresani, Manohar Paluri, Li Fei-Fei, and Juan Carlos Niebles. 2018. What makes a video a video: Analyzing temporal information in video understanding models and datasets. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 7366–7375.	Hamza Karim, Keval Doshi, and Yasin Yilmaz. 2024. Real-time weakly supervised video anomaly detection. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 6848–6856.	903
851			904
852			905
853			906
854			907
855			908
856			909
857	Suyuan Huang, Haoxin Zhang, Yan Gao, Yao Hu, and Zengchang Qin. 2024b. From image to video, what do we need in multimodal llms? <i>arXiv preprint arXiv:2404.11865</i> .	Wonkyun Kim, Changin Choi, Wonseok Lee, and Wonjong Rhee. 2024. An image grid can be worth a video: Zero-shot video question answering using a vlm. <i>arXiv preprint arXiv:2403.18406</i> .	910
858			911
859			912
860			913
861	Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025. Vision-r1: Incentivizing reasoning capability in multimodal large language models. <i>arXiv preprint arXiv:2503.06749</i> .	Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. 2024a. <b>Llava-next: What else influences visual instruction tuning beyond data?</b>	914
862			915
863			916
864			917
865			918
866	Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. 2021. Perceiver: General perception with iterative attention. In <i>International conference on machine learning</i> , pages 4651–4664. PMLR.	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024b. <b>Llava-next: Stronger llms supercharge multimodal capabilities in the wild.</b>	919
867			920
868			921
869			922
870			923
871	Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2017. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2758–2766.	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024c. <b>Llava-onevision: Easy visual task transfer.</b> <i>arXiv preprint arXiv:2408.03326</i> .	924
872			925
873			926
874			927
875			928
876	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024d. <b>Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models.</b> <i>arXiv preprint arXiv:2407.07895</i> .	929
877			930
878			931
879			932
880			933
		Hongyu Li, Jinyu Chen, Ziyu Wei, Shaofei Huang, Tianrui Hui, Jialin Gao, Xiaoming Wei, and Si Liu. 2025a. <b>Llava-st: A multimodal large language model for fine-grained spatial-temporal understanding.</b> <i>arXiv preprint arXiv:2501.08282</i> .	934
			935

936	Jiahao Li, Bin Li, and Yan Lu. 2022. Hybrid spatial-temporal entropy modelling for neural video compression. In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , pages 1503–1511.	990
937		991
938		992
939		993
940	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	994
941		995
942		996
943		997
944		998
945	KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023b. Videochat: Chat-centric video understanding. <i>arXiv preprint arXiv:2305.06355</i> .	999
946		1000
947		1001
948		1002
949	Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. 2024e. Mvbench: A comprehensive multi-modal video understanding benchmark. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 22195–22206.	1003
950		1004
951		1005
952		1006
953		1007
954		1008
955	Sheng Li, Zhiqiang Tao, Kang Li, and Yun Fu. 2019. Visual to text: Survey of image and video captioning. <i>IEEE Transactions on Emerging Topics in Computational Intelligence</i> , 3(4):297–312.	1009
956		1010
957		1011
958		1012
959	Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhan Zhu, Haian Huang, Jianfei Gao, Kunchang Li, Yanan He, Chenting Wang, et al. 2025b. Videochat-flash: Hierarchical compression for long-context video modeling. <i>arXiv preprint arXiv:2501.00574</i> .	1013
960		1014
961		1015
962		1016
963		1017
964	Yanwei Li, Chengyao Wang, and Jiaya Jia. 2023c. Llama-vid: An image is worth 2 tokens in large language models. <i>arXiv preprint arXiv:2311.17043</i> .	1018
965		1019
966		1020
967	Yunxin Li, Xinyu Chen, Baotian Hu, Longyue Wang, Haoyuan Shi, and Min Zhang. 2024f. Videovista: A versatile benchmark for video understanding and reasoning. <i>arXiv preprint arXiv:2406.11303</i> .	1021
968		1022
969		1023
970		1024
971	Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection. <i>arXiv preprint arXiv:2311.10122</i> .	1025
972		1026
973		1027
974		1028
975	Ji Lin, Chuang Gan, and Song Han. 2019. Tsm: Temporal shift module for efficient video understanding. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 7083–7093.	1029
976		1030
977		1031
978		1032
979	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 26296–26306.	1033
980		1034
981		1035
982		1036
983		1037
984	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. <b>Llava-next: Improved reasoning, ocr, and world knowledge</b> .	1038
985		1039
986		1040
987	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024c. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36.	1041
988		1042
989		1043
	Jiajun Liu, Yibing Wang, Hanghang Ma, Xiaoping Wu, Xiaoqi Ma, Xiaoming Wei, Jianbin Jiao, Enhua Wu, and Jie Hu. 2024d. Kangaroo: A powerful video-language model supporting long-context video input. <i>arXiv preprint arXiv:2408.15542</i> .	994
	Ruyang Liu, Chen Li, Haoran Tang, Yixiao Ge, Ying Shan, and Ge Li. 2024e. St-llm: Large language models are effective temporal learners. <i>arXiv preprint arXiv:2404.00308</i> .	995
	Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, et al. 2024f. Nvila: Efficient frontier visual language models. <i>arXiv preprint arXiv:2412.04468</i> .	996
	Zuyan Liu, Yuhao Dong, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. 2024g. Oryx mllm: On-demand spatial-temporal understanding at arbitrary resolution. <i>arXiv preprint arXiv:2409.12961</i> .	997
	Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Da Li, Pengcheng Lu, Tao Wang, Linmei Hu, Minghui Qiu, and Zhongyu Wei. 2023. Valley: Video assistant with large language model enhanced ability. <i>arXiv preprint arXiv:2306.07207</i> .	998
	Yongdong Luo, Xiawu Zheng, Xiao Yang, Guilin Li, Haojia Lin, Jinfa Huang, Jiayi Ji, Fei Chao, Jiebo Luo, and Rongrong Ji. 2024. Video-rag: Visually-aligned retrieval-augmented long video comprehension. <i>arXiv preprint arXiv:2411.13093</i> .	999
	Fan Ma, Xiaojie Jin, Heng Wang, Yuchen Xian, Jiashi Feng, and Yi Yang. 2023a. Vista-llama: Reliable video narrator via equal distance to visual tokens. <i>arXiv preprint arXiv:2312.08870</i> .	1000
	Yingzi Ma, Yulong Cao, Jiachen Sun, Marco Pavone, and Chaowei Xiao. 2023b. Dolphins: Multi-modal language model for driving. <i>arXiv preprint arXiv:2312.00438</i> .	1001
	Ziyu Ma, Chenhui Gou, Hengcan Shi, Bin Sun, Shutao Li, Hamid Rezaatofghi, and Jianfei Cai. 2024. Drvideo: Document retrieval based long video understanding. <i>arXiv preprint arXiv:2406.12846</i> .	1002
	Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Khan. 2024. Videogpt+: Integrating image and video encoders for enhanced video understanding. <i>arXiv preprint arXiv:2406.09418</i> .	1003
	Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. <i>arXiv preprint arXiv:2306.05424</i> .	1004
	Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2024. Egoschema: A diagnostic benchmark for very long-form video language understanding. <i>Advances in Neural Information Processing Systems</i> , 36.	1005

1044	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	1100
1045		1101
1046		
1047		
1048		
1049	Thong Nguyen, Yi Bin, Junbin Xiao, Leigang Qu, Yicong Li, Jay Zhangjie Wu, Cong-Duy Nguyen, See-Kiong Ng, and Anh Tuan Luu. 2024. <a href="#">Video-language understanding: A survey from model architecture, model training, and data perspectives</a> . In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 3636–3657, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.	1102
1050		1103
1051		1104
1052		1105
1053		1106
1054		1107
1055		
1056		
1057		
1058	Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. Im2text: Describing images using 1 million captioned photographs. <i>Advances in neural information processing systems</i> , 24.	1108
1059		1109
1060		1110
1061		1111
1062	Pinelopi Papalampidi, Frank Keller, and Mirella Lapata. 2019. Movie plot analysis via turning point identification. <i>arXiv preprint arXiv:1908.10328</i> .	1112
1063		1113
1064		
1065	Jinlong Peng, Zekun Luo, Liang Liu, and Boshen Zhang. 2024. Frih: Fine-grained region-aware image harmonization. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 4478–4486.	1114
1066		1115
1067		1116
1068		1117
1069	Reza Pourreza, Hoang Le, Amir Said, Guillaume Sautiere, and Auke Wiggers. 2023. Boosting neural video codecs by exploiting hierarchical redundancy. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 5355–5364.	1118
1070		1119
1071		
1072		
1073		
1074		
1075	Jiaxin Qi, Yulei Niu, Jianqiang Huang, and Hanwang Zhang. 2020. Two causal principles for improving visual dialog. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10860–10869.	1120
1076		1121
1077		1122
1078		1123
1079		1124
1080	Long Qian, Juncheng Li, Yu Wu, Yaobo Ye, Hao Fei, Tat-Seng Chua, Yueting Zhuang, and Siliang Tang. 2024. Momentor: Advancing video large language model with fine-grained temporal reasoning. <i>arXiv preprint arXiv:2402.11435</i> .	1125
1081		1126
1082		1127
1083		1128
1084		1129
1085	Jihao Qiu, Yuan Zhang, Xi Tang, Lingxi Xie, Tianren Ma, Pengyu Yan, David Doermann, Qixiang Ye, and Yunjie Tian. 2024. Artemis: Towards referential understanding in complex videos. <i>arXiv preprint arXiv:2406.00258</i> .	1130
1086		1131
1087		1132
1088		1133
1089		1134
1090	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	1135
1091		1136
1092		
1093		
1094		
1095		
1096	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	1137
1097		1138
1098		1139
1099		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155

1156	2024b. MovieLLM: Enhancing long video understanding with ai-generated movies. <i>arXiv preprint arXiv:2403.01422</i> .		
1157		Xiao Wang, Qingyi Si, Jianlong Wu, Shiyu Zhu, Li Cao, and Liqiang Nie. 2024c. Retake: Reducing temporal and knowledge redundancy for long video understanding. <i>arXiv preprint arXiv:2412.20504</i> .	1210 1211 1212 1213
1158			
1159	Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023. Eva-clip: Improved training techniques for clip at scale. <i>arXiv preprint arXiv:2303.15389</i> .		
1160		Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. 2024d. Videoagent: Long-form video understanding with large language model as agent. In <i>European Conference on Computer Vision</i> , pages 58–76. Springer.	1214 1215 1216 1217 1218
1161			
1162			
1163	Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. 2023. Video understanding with large language models: A survey. <i>arXiv preprint arXiv:2312.17432</i> .		
1164		Xidong Wang, Dingjie Song, Shunian Chen, Chen Zhang, and Benyou Wang. 2024e. Longlava: Scaling multi-modal llms to 1000 images efficiently via hybrid architecture. <i>arXiv preprint arXiv:2409.02889</i> .	1219 1220 1221 1222 1223
1165			
1166			
1167			
1168	InternLM Team. 2023. Internlm: A multilingual language model with progressively enhanced capabilities.		
1169		Yizhou Wang, Ruiyi Zhang, Haoliang Wang, Uttaran Bhattacharya, Yun Fu, and Gang Wu. 2023b. Vaquita: Enhancing alignment in llm-assisted video understanding. <i>arXiv preprint arXiv:2312.02310</i> .	1224 1225 1226 1227
1170			
1171	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		
1172		Yuxuan Wang, Cihang Xie, Yang Liu, and Zilong Zheng. 2024f. Videollamb: Long-context video understanding with recurrent memory bridges. <i>arXiv preprint arXiv:2409.01071</i> .	1228 1229 1230 1231
1173			
1174			
1175			
1176			
1177	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023b. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		
1178		Xiu-Shen Wei, Yi-Zhe Song, Oisín Mac Aodha, Jianxin Wu, Yuxin Peng, Jinhui Tang, Jian Yang, and Serge Belongie. 2021. Fine-grained image analysis with deep learning: A survey. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 44(12):8927–8948.	1232 1233 1234 1235 1236 1237
1179			
1180			
1181			
1182			
1183	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023c. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		
1184		Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. 2024. Longvlm: Efficient long video understanding via large language models. <i>arXiv preprint arXiv:2404.03384</i> .	1238 1239 1240 1241
1185			
1186			
1187			
1188			
1189	Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 3156–3164.		
1190		Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krahenbuhl, and Ross Girshick. 2019. Long-term feature banks for detailed video understanding. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 284–293.	1242 1243 1244 1245 1246 1247
1191			
1192			
1193			
1194	Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. 2023a. Selective structured state-spaces for long-form video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 6387–6397.		
1195		Chao-Yuan Wu and Philipp Krahenbuhl. 2021. Towards long-form video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 1884–1894.	1248 1249 1250 1251
1196			
1197			
1198		Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. 2024. Longvideobench: A benchmark for long-context interleaved video-language understanding. <i>arXiv preprint arXiv:2407.15754</i> .	1252 1253 1254 1255
1199			
1200	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .		
1201		Rujie Wu, Xiaojian Ma, Hai Ci, Yue Fan, Yuxuan Wang, Haozhe Zhao, Qing Li, and Yizhou Wang. 2025. Longvitu: Instruction tuning for long-form video understanding. <i>arXiv preprint arXiv:2501.05037</i> .	1256 1257 1258 1259
1202			
1203			
1204			
1205	Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao Dong, Ming Ding, et al. 2024b. Lvbench: An extreme long video understanding benchmark. <i>arXiv preprint arXiv:2406.08035</i> .		
1206		Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 9777–9786.	1260 1261 1262 1263 1264
1207			
1208			
1209			

1265	Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In <i>Proceedings of the 25th ACM international conference on Multimedia</i> , pages 1645–1653.	Rufai Yusuf Zakari, Jim Wilson Owusu, Hailin Wang, Ke Qin, Zaharaddeen Karami Lawal, and Yuezhou Dong. 2022. Vqa and visual reasoning: An overview of recent datasets, methods and challenges. <i>arXiv preprint arXiv:2212.13296</i> .	1318
1266			1319
1267			1320
1268			1321
1269			1322
1270			
1271	Jiaqi Xu, Cuiling Lan, Wenxuan Xie, Xuejin Chen, and Yan Lu. 2023. Retrieval-based video language model for efficient long video question answering. <i>arXiv preprint arXiv:2312.04931</i> .	Xiangyu Zeng, Kunchang Li, Chenting Wang, Xinhao Li, Tianxiang Jiang, Ziang Yan, Songze Li, Yansong Shi, Zhengrong Yue, Yi Wang, et al. 2024. Timesuite: Improving mllms for long video understanding via grounded tuning. <i>arXiv preprint arXiv:2410.19702</i> .	1323
1272			1324
1273			1325
1274			1326
1275			1327
1276	Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msrvtv: A large video description dataset for bridging video and language. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 5288–5296.	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 11975–11986.	1328
1277			1329
1278			1330
1279			1331
1280	Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. 2024a. Pllava: Parameter-free llava extension from images to videos for video dense captioning. <i>arXiv preprint arXiv:2404.16994</i> .	Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang, Hang Zhang, Xin Li, et al. 2025a. Videollama 3: Frontier multimodal foundation models for image and video understanding. <i>arXiv preprint arXiv:2501.13106</i> .	1333
1281			1334
1282			1335
1283			1336
1284			1337
1285	Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin Dehghan. 2024b. Slowfast-llava: A strong training-free baseline for video large language models. <i>arXiv preprint arXiv:2407.15841</i> .	Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024a. <b>MM-LLMs: Recent advances in MultiModal large language models</b> . In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 12401–12430, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.	1339
1286			1340
1287			1341
1288			1342
1289			1343
1290	Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian Tang, Shang Yang, Zhijian Liu, et al. 2024. Longvila: Scaling long-context visual language models for long videos. <i>arXiv preprint arXiv:2408.10188</i> .	Hang Zhang, Xin Li, and Lidong Bing. 2023. Videollama: An instruction-tuned audio-visual language model for video understanding. <i>arXiv preprint arXiv:2306.02858</i> .	1344
1291			1345
1292			1346
1293			1347
1294			1348
1295	Sosuke Yamao, Natsuki Miyahara, Yuki Harazono, and Shun Takeuchi. 2024. Iqvic: In-context, question adaptive vision compressor for long-term video understanding lms. <i>arXiv preprint arXiv:2412.09907</i> .	Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. 2024b. Vision-language models for vision tasks: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> .	1350
1296			1351
1297			1352
1298			1353
1299	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .	Jingyi Zhang, Jiaying Huang, Huanjin Yao, Shunyu Liu, Xikun Zhang, Shijian Lu, and Dacheng Tao. 2025b. R1-vl: Learning to reason with multimodal large language models via step-wise group relative policy optimization. <i>arXiv preprint arXiv:2503.12937</i> .	1354
1300			1355
1301			1356
1302			1357
1303	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. <i>arXiv preprint arXiv:2306.13549</i> .	Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. 2024c. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. <i>arXiv preprint arXiv:2407.03320</i> .	1359
1304			1360
1305			1361
1306			1362
1307	Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. <i>arXiv preprint arXiv:2403.04652</i> .	Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Hao-ran Tan, Chunyuan Li, and Ziwei Liu. 2024d. Long context transfer from language to vision. <i>arXiv preprint arXiv:2406.16852</i> .	1363
1308			1364
1309			1365
1310			1366
1311			1367
1312	Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 9127–9134.	Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang Feng. 2025c. Llava-mini: Efficient image and video large multimodal models with one vision token. <i>arXiv preprint arXiv:2501.03895</i> .	1368
1313			1369
1314			1370
1315			1371
1316			1372
1317			1373

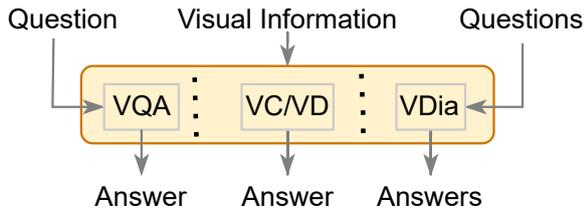


Figure 5: Various visual reasoning tasks.

1374 Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee,  
 1375 Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and  
 1376 Chunyuan Li. 2024e. Llava-next: A strong zero-shot  
 1377 video understanding model.

1378 Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao,  
 1379 Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang,  
 1380 and Zheng Liu. 2024. Mlvu: A comprehensive  
 1381 benchmark for multi-task long video understanding.  
 1382 *arXiv preprint arXiv:2406.04264*.

## 1383 A Multiple Visual Reasoning

1384 Visual reasoning demands models to comprehend  
 1385 and interpret visual information and integrate mul-  
 1386 timodal perception with commonsense understand-  
 1387 ing (Johnson et al., 2017; Chen et al., 2024c). There  
 1388 are three main types of visual reasoning tasks: vi-  
 1389 sual question answering (VQA), visual captioning  
 1390 (VC) or description (VD), and visual dialog (VDia).  
 1391 VQA (Antol et al., 2015; Zakari et al., 2022) in-  
 1392 volves generating a natural language answer based  
 1393 on the input visual data and accompanying ques-  
 1394 tions. VC and VD systems (Vinyals et al., 2015;  
 1395 Sharma et al., 2018; Li et al., 2019) typically gener-  
 1396 ate a concise, natural language sentence that sum-  
 1397 marizes the main content of the visual data and a  
 1398 detailed and comprehensive description of the cor-  
 1399 responding visual data, respectively. VDia (Das  
 1400 et al., 2017; Qi et al., 2020) involves multi-turn  
 1401 conversations, consisting of a series of question-  
 1402 answer pairs centered around the visual content.

## 1403 B Development of MM-LLM Model 1404 Architecture

### 1405 B.1 Multiple MM-LLMs

1406 As illustrated in Fig. 6, MM-LLMs for images,  
 1407 short videos, and long videos share a similar struc-  
 1408 ture comprising a visual encoder, an LLM back-  
 1409 bone, and an intermediary connector. Unlike  
 1410 the image-level connector in image-targeted MM-  
 1411 LLMs, the video-level connector is crucial for in-  
 1412 tegrating cross-frame visual information. In LV-  
 1413 LLMs, designing the connector is more challeng-

1414 ing, requiring efficient compression of amounts  
 1415 of visual information and incorporating temporal  
 1416 knowledge to manage long-term correlations.

### 1417 B.2 Multiple LLM Backbones

1418 Compared to closed-source LLMs like GPT-3/4  
 1419 (Brown, 2020; Achiam et al., 2023) and Gemini-  
 1420 1.5 (Reid et al., 2024), various open-source LLMs  
 1421 are more commonly used in implementing visual  
 1422 LLMs. These include Flan-T5 (Chung et al., 2024),  
 1423 LLaMA (Touvron et al., 2023b,c; Dubey et al.,  
 1424 2024), Vicuna (Chiang et al., 2023), QWen (Bai  
 1425 et al., 2023a), Mistral (Jiang et al., 2023), Open-  
 1426 flamingo (Awadalla et al., 2023), Yi (Young et al.,  
 1427 2024), and InternLM (Team, 2023; Cai et al.,  
 1428 2024).

### 1429 B.3 Various Connector Designs

1430 In addition to the detailed discussed long-video-  
 1431 level connectors, the image-level and video-level  
 1432 connectors are also popular.

1433 **Image-level connectors.** Image-level connectors  
 1434 are used to map image features to the language  
 1435 space for processing raw visual tokens, and they  
 1436 are widely used in both image-targeted and video-  
 1437 targeted MM-LLMs. These connectors can be cate-  
 1438 gorized into three groups: **The first group** directly  
 1439 uses a single linear layer (Liu et al., 2024c) or a  
 1440 multi-layer perceptron (MLP) (Liu et al., 2024a)  
 1441 to map image features into the language embed-  
 1442 ding space. However, this method, which retains  
 1443 all visual tokens, is not suitable for visual under-  
 1444 standing tasks involving multiple images. To ad-  
 1445 dress the limitations of retaining all visual tokens,  
 1446 **the second group** employs various pooling-based  
 1447 methods. These include spatial pooling (Maaz  
 1448 et al., 2023), adaptive pooling (Xu et al., 2024a),  
 1449 semantic-similar token merging (Jin et al., 2024),  
 1450 and adjacent token averaging (Zhang et al., 2024e;  
 1451 Li et al., 2024c). **The third group** utilizes cross-  
 1452 attention or transformer-based structures, such as  
 1453 Q-Former (Li et al., 2023a) and Perceiver Resam-  
 1454 pler (Jaegle et al., 2021), for image feature com-  
 1455 pression. Q-Former is a lightweight transformer  
 1456 structure that employs a set of learnable query vec-  
 1457 tors to extract and compress visual features. Many  
 1458 visual LLMs (Dai et al., 2023; Li et al., 2023b; Ma  
 1459 et al., 2023a; Liu et al., 2024e), following BLIP-2,  
 1460 choose the Q-Former-based connector. Other vi-  
 1461 sual LLMs (Ma et al., 2023b; Jiang et al., 2024)  
 1462 opt for the Perceiver Resampler to reduce compu-  
 1463 tational burden by extracting patch features.

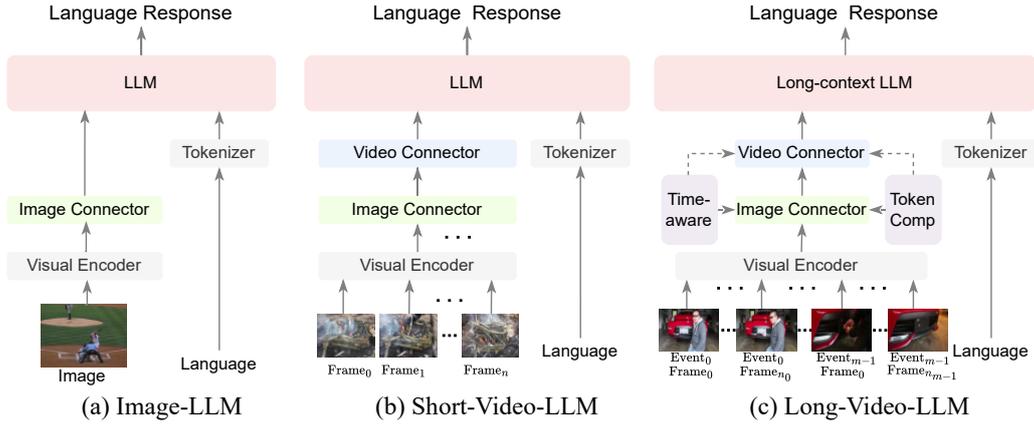


Figure 6: MM-LLMs of (b): Image-LLM, (c) Short-Video-LLM and (c) Long-Video-LLM.

**Video-level connectors.** Video-level connectors are used for extracting sequential visual data and further compressing visual features. Compared to the solely image-level connectors in image-targeted MM-LLMs, video-level connectors are essential for video-targeted MM-LLMs, including LV-LLMs. Some methods directly concatenate image tokens before inputting them to the LLMs, making them sensitive to the number of frame images (Dai et al., 2023; Lin et al., 2023). Similar structures used for token compression in image-level connectors can be adapted for video-level interfaces, such as pooling-based and transformer-based structures. Pooling along the time series dimension is a straightforward way to reduce temporal information redundancy (Maaz et al., 2023; Song et al., 2024a). Transformer-based methods, such as Video Q-Former (Zhang et al., 2023; Ma et al., 2023a; Ren et al., 2024a) and Video Perceiver (Wang et al., 2023b), are effective in extracting video features while reducing data complexity. Additionally, 3D-Convolution-based methods can extract and compress visual data from both the spatial and temporal dimensions (Cheng et al., 2024b; Liu et al., 2024d).

#### B.4 Training Design for LVU

As shown in Table 3, the training devices and resources used in pre-training and supervised fine-tuning are summarized. Adequate computing power and sufficient training data are essential for developing a robust long video understanding model.

### C Video Understanding from Seconds to Minutes

As shown in Table 4, we summarize the general video understanding performance of various vi-

sual LLMs on open-ended video question answering benchmarks, including TGIF-QA (Jang et al., 2017), MSVD-QA, MSRVT-QA (Xu et al., 2017), NEXT-QA (Xiao et al., 2021), and ActivityNet-QA (Yu et al., 2019). Additionally, we consider the VideoChatGPT-introduced video-based generative performance benchmark (Maaz et al., 2023), which evaluates five aspects of video-based text generation: Correctness of Information (CI), Detail Orientation (DO), Context Understanding (CU), Temporal Understanding (TU), and Consistency (CO). The video benchmarks with lengths shorter than 1 minute, such as TGIF-QA, MSVD-QA, MSRVT-QA, and NEXT-QA, are commonly used for short video understanding. In contrast, benchmarks exceeding one minute, such as ActivityNet-QA and the ActivityNet-200-based (Caba Heilbron et al., 2015) generative performance benchmark, are used for long video understanding.

By comparing the performance in Table 4, we can conclude that long video understanding is challenging, with the following findings: (1) Video reasoning with more frames introduces more complex visual information and is more challenging. Methods designed to support long videos, such as LongVA (Zhang et al., 2024d), show better performance compared to being fed with fewer frames on the same video dataset. However, performance decreases when being fed with more frames from the same video dataset for methods without special designs for long videos, like VideoLLaMA2 (Cheng et al., 2024b). (2) Short video understanding methods that perform well on seconds-level video understanding often do not perform well on minutes-level moderately long video understanding, such as RED-VILLM (Huang et al., 2024b) and MiniGPT4-Video (Ataallah et al., 2024a). Long video un-

derstanding methods tend to share consistently good performance on both short and moderately long video benchmarks, such as ST-LLM (Liu et al., 2024e), SlowFast-LLaVA (Xu et al., 2024b), PLLaVA (Xu et al., 2024a), and MovieChat (Song et al., 2024a). This improvement likely stems from better-captured spatiotemporal information in specially designed long video understanding methods.

Analyzing the trade-offs between design choices for modeling short and long videos is crucial. Current long video understanding methods are still suboptimal. Short video understanding methods can outperform long video methods on long video benchmarks if trained with more high-quality data or equipped with larger LLM backbones, such as LLaVA-NeXT-Video (Zhang et al., 2024e) and LLaVA-OneVision (Li et al., 2024c). Conversely, long-video-specific models do not perform well on short video benchmarks. As noted in LLaV-Next (Zhang et al., 2024e), combining different training resources has proven more effective. However, model design must balance these trade-offs: long video models require more frames but fewer visual details compared to short video models. Supporting hour-long videos necessitates powerful visual token compression, which may reduce short video understanding performance. Future versions will include a detailed examination of key aspects such as token compression strategies, handling temporal dynamics, and architectural choices for long video model design.

## D More Application Scenarios on Long Video Understanding

Long video understanding with large models faces several key challenges for more long video applications. Contextual understanding is critical, as long videos require models to maintain temporal coherence and contextual awareness over extended periods (He et al., 2024). Real-time processing (Karim et al., 2024) is essential for applications like surveillance, live event analysis, and embodied AI, necessitating the development of low-latency models capable of processing video streams in real-time. Multi-modal integration is another frontier, as long videos often contain audio, text, and visual information (Zhang et al., 2023; Cheng et al., 2024b). Future models should better integrate these modalities to enhance understanding and provide a more holistic analysis of video content.

Model	Year	Hardware	Training	
			PT	IT
InstructBLIP (2023)	23.05	16 A100-40G	Y-N-N	Y-N-N
VideoChat (2023b)	23.05	1 A10	Y-Y-N	Y-Y-N
Video-LLaMA (2023)	23.06	-	Y-Y-N	Y-Y-N
Video-ChatGPT (2023)	23.06	8 A100-40G	N-N-N	N-Y-N
Valley (2023)	23.06	8 A100 80G	Y-Y-N	Y-Y-N
MovieChat (2024a)	23.07	-	E2E	E2E
Qwen-VL (2023b)	23.08	-	Y-N-N	Y-N-N
Chat-UniVi (2024)	23.11	-	Y-N-N	Y-Y-N
Video-LLaVA (2023)	23.11	4 A100-80G	Y-Y-N	Y-Y-N
LLaMA-VID (2023c)	23.11	8 A100	Y-Y-N	Y-Y-Y
VTimeLLM (2024a)	23.11	1 RTX-4090	Y-Y-N	N-Y-N
VideoChat2 (2024e)	23.11	-	Y-Y-N	Y-Y-N
Vista-LLaMA (2023a)	23.12	8 A100-80GB	E2E	E2E
TimeChat (2024a)	23.12	8 V100-32G	Y-Y-N	N-N-Y
VaQuitA (2023b)	23.12	8 A100-80GB	E2E	E2E
Dolphins (2023b)	23.12	4 A100	N-Y-N	Y-Y-N
Momentor (2024)	24.02	8 A100	Y-Y-N	N-Y-N
MovieLLM (2024b)	24.03	4 A100	Y-Y-N	Y-Y-Y
MA-LMM (2024)	24.04	4 A100	E2E	E2E
PLLaVA (2024a)	23.04	-	Y-N-N	Y-Y-N
LongVLM (2024)	23.04	4 A100 80G	Y-N-N	Y-Y-N
MiniGPT4-Video (2024a)	24.04	-	Y-Y-N	N-Y-N
RED-VILLM (2024b)	24.04	-	Y-N-N	Y-Y-N
ST-LLM (2024e)	24.04	8 A100	E2E	E2E
LLaVA-NeXT-Video (2024e)	24.04	-	Y-Y-N	Y-Y-N
Mantis-Idedics2 (2024)	24.05	16 A100-40G	Y-N-N	N-Y-N
VideoLLaMA 2 (2024b)	24.06	-	Y-Y-N	Y-Y-N
LongVA (2024d)	24.06	8x A100-80G	-	Y-N-N
Artemis (2024)	24.06	8 x A800	Y-Y-N	N-Y-N
VideoGPT+ (2024)	24.06	8 x A100 40G	Y-Y-N	N-Y-N
IXC-2.5 (2024c)	24.07	-	Y-Y-N	Y-Y-N
EVLM (2024b)	24.07	-	Y-Y-N	Y-Y-N
SlowFast-LLaVA (2024b)	24.07	A100-80G	-	-
LLaVA-NeXT-Interleave (2024d)	24.07	-	Y-N-N	Y-Y-N
Kangaroo (2024d)	24.08	-	Y-Y-N	Y-Y-Y
VITA (2024b)	24.08	-	Y-Y-N	Y-Y-N
LLaVA-OneVision (2024c)	24.08	-	Y-N-N	Y-Y-N
LONGVILA (2024)	24.08	256 A100 80G	Y-Y-N	Y-Y-Y
LongLLaVA (2024e)	24.09	24 A800 80G	Y-N-N	Y-Y-N
Qwen2-VL (2024a)	24.09	-	Y-N-N	Y-Y-N
Video-XL (2024)	20.09	8 A800-80G	Y-N-N	Y-Y-N
Oryx-1.5 (2024g)	24.10	64 A800-80G	Y-Y-N	Y-Y-Y
TimeMarker (2024d)	24.11	-	Y-Y-Y	Y-Y-Y
NVILA (2024f)	24.12	128 H100-80G	Y-Y-N	Y-Y-Y
InternVL2.5 (2024e)	24.12	-	Y-Y-N	Y-Y-N
VideoLLaMA 3 (2025a)	25.01	-	Y-Y-N	Y-Y-N
LLaVA-ST (2025a)	25.01	48 A100	N-Y-N	N-Y-N
Qwen2.5-VL (2025)	25.02	-	Y-Y-Y	Y-Y-Y
Vision-R1 (2025)	25.03	64 H800-80G	-	+RL
R1-VL (2025b)	25.03	4 H100-80GB	-	+RL
ViLAMP (2025)	25.04	21 x A100-80G	N-Y-N	N-Y-Y

Table 3: Comparison of mainstream Video-LLMs on training design. "PT" and "IT" denote the two stages of pre-training and instruction-tuning during model training. The letters "Y" (Yes) and "N" (No) indicate whether image, short-video, and long-video language datasets are used in these stages. "E2E" stands for an end-to-end training pipeline.

Model	LLM	Long	TGIF-QA	MSVD-QA	MSRVTT-QA	NeXT-QA	ActivityNet-QA	GPT-based Evaluation( 2mins)							
			(2-5s)	(10-15s)	(10-15s)	(42.9s)	( 2mins)	CI	DO	CU	TU	CO	Average		
InstructBLIP	Vicuna-7B	✗	-	41.8/	22.1/	-	-	-	-	-	-	-	-	-	-
Video-ChatGPT	Vicuna1.1-7B	✗	51.4/3.0	64.9/3.3	49.3/2.8	-	35.2/2.8	2.40	2.52	2.62	1.98	2.37	2.38	-	-
MA-LMM	Vicuna-7B	✓	-	60.6/	48.5/	-	49.8/	-	-	-	-	-	-	-	-
Valley	StableVicuna-7B	✗	-	60.5/3.3	51.1/2.9	-	45.1/3.2	2.43	2.13	2.86	2.04	2.45	2.38	-	-
MovieLLM	Vicuna-7B	✓	-	63.2/3.5	52.1/3.1	-	43.3/3.3	2.64	2.61	2.92	2.03	2.43	2.53	-	-
Vista-LLaMA	Vicuna-7B	✗	-	65.3/3.6	60.5/3.3	60.7/3.4	48.3/3.3	2.44	2.31	2.64	3.18	2.26	2.57	-	-
RED-VILLM	LLaVA-7B	✗	55.9/3.1	68.9/2.8	52.4/2.9	-	39.2/3.0	2.57	2.64	3.13	2.21	2.39	2.59	-	-
Momentor	LLaMA-7B	✗	-	68.9/3.6	55.6/3.0	-	40.8/3.2	-	-	-	-	-	-	-	-
Video-LLaVA	Vicuna1.5-7B	✗	70.0/4.0	70.7/3.9	59.2/3.5	-	45.3/3.3	-	-	-	-	-	-	-	-
Artemis	Vicuna1.5-7B	✓	-	72.1/3.9	56.7/3.2	-	39.3/2.9	2.69	2.55	3.04	2.24	2.70	2.64	-	-
MovieChat	LLaMA-7B	✓	-	75.2/3.8	52.7/2.6	-	45.7/3.4	2.76	2.93	3.01	2.24	2.42	2.67	-	-
VaQuitA	LLaMA-7B	✗	-	74.6/3.7	68.6/3.3	-	48.8/3.3	-	-	-	-	-	-	-	-
RED-VILLM	QWen-VL-7B	✗	62.3/3.3	71.2/3.7	53.9/3.1	-	44.2/3.2	2.69	2.72	3.32	2.32	2.47	2.70	-	-
MiniGPT4-Video	Mistral-7B	✗	72.2/4.1	73.9/4.1	58.3/3.5	-	44.3/3.4	2.97	2.58	3.17	2.38	2.44	2.71	-	-
VTimeLLM	Vicuna-7B	✗	-	-	-	-	-	2.49	2.78	3.10	3.40	2.47	2.85	-	-
MiniGPT4-Video	LLaMA2-7B	✗	67.9/3.7	72.9/3.8	58.8/3.3	-	45.9/3.2	2.93	2.97	3.45	2.47	2.60	2.88	-	-
Chat-UniVi	Vicuna1.5-7B	✗	69.0/3.8	69.3/3.7	55.0/3.1	-	46.1/3.3	2.89	2.91	3.46	2.40	2.81	2.89	-	-
LLaMA-VID	Vicuna-7B	✓	-	69.7/3.7	57.7/3.2	-	47.4/3.3	2.96	3.00	3.53	2.46	2.51	2.89	-	-
LongVLM	Vicuna1.1-7B	✓	-	70.0/3.8	59.8/3.3	-	47.6/3.3	2.76	2.86	3.34	2.39	3.11	2.89	-	-
VideoChat2	Vicuna0-7B	✗	-	70.0/3.9	54.1/3.3	-	49.1/3.3	3.02	2.88	3.51	2.66	2.81	2.98	-	-
SlowFast-LLaVA	Vicuna1.5-7B	✓	78.7/4.2	79.1/4.1	65.8/3.6	64.2/	56.3/3.4	3.09	2.70	3.57	2.52	3.35	3.04	-	-
PLLaVA	LLaVA-Next-7B	✓	77.5/4.1	76.6/4.1	62.0/3.5	-	56.3/3.5	3.21	2.86	3.62	2.33	2.93	3.12	-	-
VideoLLaMA2-16	Mistral-7B-Instruct	✗	-	70.9/3.8	-	-	50.2/3.3	3.16	3.08	3.69	2.56	3.14	3.13	-	-
VideoLLaMA2-8	Mistral-7B-Instruct	✗	-	71.7/3.9	-	-	49.9/3.3	3.09	3.09	3.68	2.63	3.25	3.15	-	-
ST-LLM	Vicuna-7B	✓	-	74.6/3.9	63.2/3.4	-	50.9/3.3	3.23	3.05	3.74	2.93	2.81	3.15	-	-
LongVA-32	Qwen2-7B-224K	✓	-	-	-	67.1/	72.8	3.65	3.08	3.10	3.74	2.28	3.17	-	-
LongVA-64	Qwen2-7B-224K	✓	-	-	-	68.3/	72.8	3.64	3.05	3.09	3.77	2.44	3.20	-	-
LLaVA-NeXT-Video	Vicuna1.5-7B	✗	-	-	-	-	53.5/3.2	3.39	3.29	3.92	2.60	3.12	3.26	-	-
LLaVA-NeXT-Interleave	Qwen1.5-7B	✗	-	-	-	78.2	55.3/3.13	3.51	3.28	3.89	2.77	3.68	3.43	-	-
LLaVA-OneVision	Qwen2-7B	✗	-	-	-	-	56.6/	-	-	-	-	-	3.49	-	-
LongVA-32-DPO	Qwen2-7B-224K	✓	-	-	-	69.3/	72.8	4.07	3.55	3.32	4.09	2.86	3.58	-	-
LLaVA-NeXT-Video-DPO	Vicuna1.5-7B	✗	-	-	-	-	60.2/3.5	3.64	3.45	4.17	2.95	4.08	3.66	-	-
InstructBLIP	Vicuna-13B	✗	-	41.2/	24.8/	-	-	-	-	-	-	-	-	-	-
LLaMA-VID	Vicuna-13B	✓	-	70.0/3.7	58.9/3.3	-	47.5/3.3	3.07	3.05	3.60	2.58	2.63	2.99	-	-
PLLaVA	LLaVA-Next-13B	✓	77.8/4.2	75.7/4.1	63.2/3.6	-	56.3/3.6	3.27	2.99	3.66	2.47	3.09	3.27	-	-
LLaVA-NeXT-Interleave	Qwen1.5-14B	✗	-	-	-	79.1	56.2/3.19	3.65	3.37	3.98	2.74	3.67	3.48	-	-
LLaVA-NeXT-Interleave-DPO	Qwen1.5-14B	✗	-	-	-	77.9	55.0/3.13	3.99	3.61	4.24	3.19	4.12	3.83	-	-
SlowFast-LLaVA	Nous-Hermes-2-Yi-34B	✓	80.6/4.3	79.9/4.1	67.4/3.7	-	59.2/3.5	3.48	2.96	3.84	2.77	3.57	3.32	-	-
LLaVA-NeXT-Video	Nous-Hermes-2-Yi-34B	✗	-	-	-	-	58.8/3.4	3.48	3.37	3.95	2.64	3.28	3.34	-	-
PLLaVA	LLaVA-Next-34B	✓	80.6/4.3	79.9/4.2	68.7/3.8	-	60.9/3.7	3.60	3.20	3.90	2.67	3.25	3.48	-	-
LLaVA-NeXT-Video-DPO	Nous-Hermes-2-Yi-34B	✗	-	-	-	-	64.4/3.6	3.81	3.55	4.24	3.14	4.12	3.77	-	-

Table 4: Comparison of mainstream Video-LLMs on video understanding benchmarks of different lengths. Methods with ✓ in the "Long" column are designed for long videos.