



VideoLLaMA 3: Frontier Multimodal Foundation Models for Image and Video Understanding

Boqiang Zhang[♠], Kehan Li[♠], Zesen Cheng[♠], Zhiqiang Hu[♠], Yuqian Yuan[♠],
Guanzheng Chen[♠], Sicong Leng[♠], Yuming Jiang[♠][◇], Hang Zhang[♠][◇], Xin Li[♠][◇],
Peng Jin, Wenqi Zhang, Fan Wang, Lidong Bing, Deli Zhao

DAMO Academy, Alibaba Group

[♠]Equal Contribution [◇]Project Lead

<https://github.com/DAMO-NLP-SG/VideoLLaMA3>

Abstract

In this paper, we propose VideoLLaMA 3, a more advanced multimodal foundation model for image and video understanding. The core design philosophy of VideoLLaMA3 is vision-centric. The meaning of “vision-centric” is two-fold: the vision-centric training paradigm and vision-centric framework design. The key insight of our vision-centric training paradigm is that **high-quality image-text data is crucial** for both image and video understanding. Instead of preparing massive video-text datasets, we focus on constructing large-scale, high-quality image-text datasets. VideoLLaMA3 has four training stages: 1) **Vision Encoder Adaptation**, which enables the vision encoder to accept images of variable resolutions as input; 2) **Vision-Language Alignment**, which jointly tunes the vision encoder, projector, and LLM with large-scale image-text data covering multiple types (including scene images, documents, and charts) as well as text-only data. 3) **Multi-task Fine-tuning**, which incorporates image-text SFT data for downstream tasks and video-text data to establish a foundation for video understanding. 4) **Video-centric Fine-tuning**, which further improves the model’s capability in video understanding. As for the framework design, to better capture fine-grained details in images, the pretrained vision encoder is adapted to encode images of varying sizes into vision tokens with corresponding numbers, rather than a fixed number of tokens. For video inputs, we reduce the number of vision tokens according to their similarity so that the representation of videos will be more precise and compact. Benefiting from vision-centric designs, VideoLLaMA3 achieves compelling performances in both image and video understanding benchmarks.

1 Introduction

Recent years have witnessed the rapid growth of Large Language Models (LLMs) [1–6], which significantly enhance natural language processing and understanding. The growth of LLMs enables intelligence at the language level. However, to progress further, we need

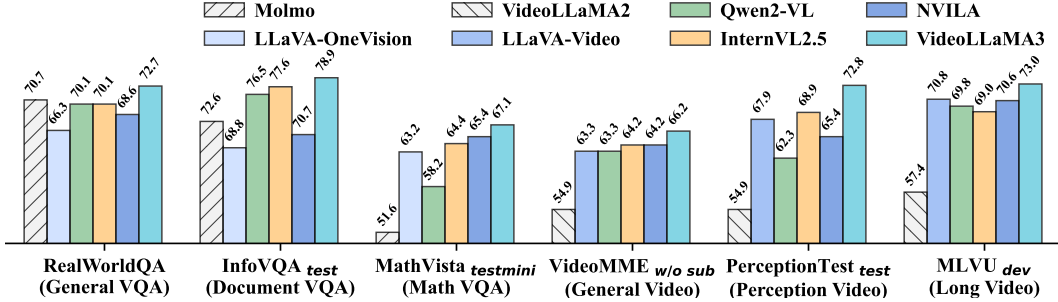


Figure 1: **Performance Comparison** of VideoLLaMA3 with the previous advanced image/video MLLM on various representative benchmarks. Specifically, VideoLLaMA3 not only demonstrates strong video understanding capabilities (VideoMME, PerceptionTest, MLVU) but also maintains excellent document comprehension abilities (DocVQA) and multimodal mathematical reasoning skills (MathVista). Note that LLaVA-OneVision is only used to evaluate image benchmarks, while LLaVA-Video is only used to evaluate video benchmarks.

intelligence that extends beyond language, as the world itself is inherently multimodal. Specifically, the model should be capable of perceiving both static scenes and dynamic environments, which necessitates the ability to understand images and videos. Building upon the success of LLMs, Multimodal LLMs (MLLMs) [7–10] have been proposed for multimodal understanding.

Existing MLLMs [11–33] have made significant progress in multimodal understanding. Image-centric MLLMs [7, 29, 31, 34–37], leveraging high-quality image-text datasets [29, 38–43] that are easier to collect and curate, have demonstrated strong performance in image understanding, such as visual question answering, OCR, and document understanding. Beyond static content such as images, video-centric MLLMs [23, 25, 28, 44] must tackle the added complexity of modeling the temporal dimension of videos, requiring models to handle dynamic content and capture dependencies across frames. This temporal complexity, combined with the need for large-scale video-text datasets that are often of lower quality and harder to annotate, makes video MLLMs more challenging. These challenges underscore the advantages of using image understanding as a foundation for video understanding. By extending the robust visual capabilities of image MLLMs, video models can focus on and better address the unique challenges of temporal and dynamic content modeling.

Inheriting from VideoLLaMA [45] and VideoLLaMA2 [46], VideoLLaMA3, a more advanced multimodal foundation model, is proposed for image and video understanding. We design VideoLLaMA3 in a vision-centric way. Specifically, we propose a vision-centric training paradigm and vision-centric framework designs. For the training paradigm, considering the intrinsic relationship between image and video modalities - where videos are essentially sequences of temporally correlated images, we prioritize the improvement of image understanding, which in turn enhances the performance of video understanding. Moreover, compared to video-text data, image-text data is easier to collect and ensures higher data quality. For vision-centric framework designs, we propose adapting the vision encoder to handle images of any resolution during the image understanding enhancement stage and tuning the encoder to efficiently embed video inputs.

Our vision-centric training paradigm consists of four stages (Figure 2): **1) Vision Encoder Adaptation:** This stage aligns the vision encoder’s feature space with LLMs. Inputs to the vision encoder are adapted from fixed to dynamic resolutions. Scene images with short captions are used to enhance the encoder’s performance, while document and scene text images are used to enable the encoder to capture fine-grained visual details. **2) Vision-Language Alignment:** This stage establishes the foundation for multimodal understanding using detailed image-text data. Scene images are annotated with detailed captions, and document and chart data include extensive explanations. To enhance spatial reasoning, fine-grained image-text data with bounding boxes are utilized. A small amount of text-only data is included to retain the model’s language capabilities. All parameters are unfrozen

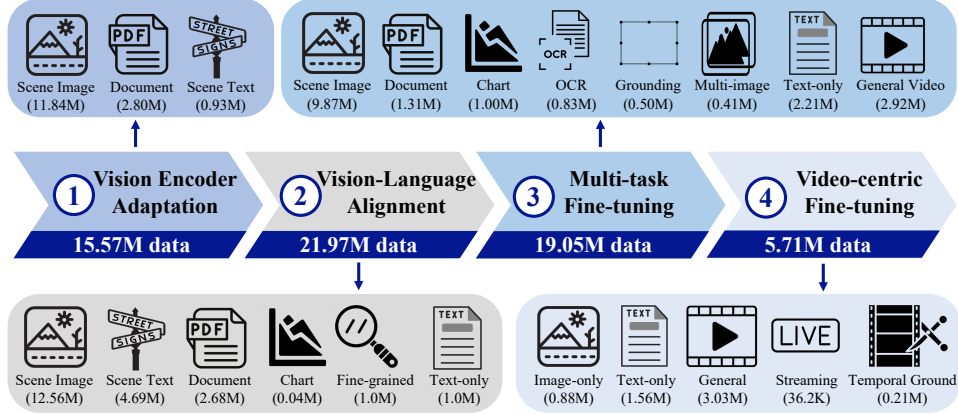


Figure 2: **Training paradigm of VideoLLaMA3.** The training of VideoLLaMA3 has four stages: (1) Vision Encoder Adaptation, (2) Vision-Language Alignment, (3) Multi-task Fine-tuning, and (4) Video-centric Fine-tuning.

during this stage. **3) Multi-Task Fine-tuning:** In this stage, the model is fine-tuned for downstream tasks, such as interactive question answering. Image-text data with questions and answers are employed, along with general video caption data to prepare the model for video perception. The use of general video caption data also surprisingly improves the performance of image understanding. **4) Video-centric Fine-tuning:** This final stage enhances the model’s performance in video understanding and video question answering. Training data includes general videos, streaming videos, videos annotated with temporal grounding information, image-only and text-only data.

On the model side, we enhance the vision encoder with two vision-centric designs: 1) we adapt the vision encoder to take images with dynamic resolutions as inputs, and 2) we lift the vision encoder to receive videos and compress the video tokens into more compacted representations. In previous methods [29, 31, 34, 47, 48], vision tokens are either with fixed numbers or with numbers among several fixed choices, which is an inflexible and unnatural way to represent images. To alleviate this limitation, we adapt the pretrained vision encoder to receive images with variable shapes. This is achieved by replacing the fixed positional embeddings with the Rotary Position Embedding (RoPE). We finetune the vision encoder in the vision encoder adaptation stage so that it can accommodate dynamic inputs. In this way, enabling it to process high-resolution images and images with unusual aspect ratios with minimal information loss. As for video inputs, we consider the redundant information in videos and propose to reduce the number of vision tokens to represent a video. The advantages of vision token compression are two-fold. One is to make the visual embeddings of videos more compact and precise so that the model can focus more on the dynamic parts of videos. The other is to save computation demands during training and inference for video understanding.

Thanks to the vision-centric training paradigm and framework designs, our proposed VideoLLaMA3 achieves state-of-the-art performance on both image and video understanding benchmarks (Figure 1). Notably, in image understanding, the performance in chart understanding and vision-related math problems surpasses state-of-the-art models by a large margin. While in video understanding, our model achieves state-of-the-art performance in most benchmarks including general video understanding, long video understanding, temporal reasoning and grounding.

To summarize, the key contributions of VideoLLaMA3 include:

- We propose VideoLLaMA3, a more advanced multimodal foundation model, for both image and video understanding. The model achieves state-of-the-art performance on most image and video understanding benchmarks. Notably, VideoLLaMA3 has significant improvements compared to previous versions of VideoLLaMA.

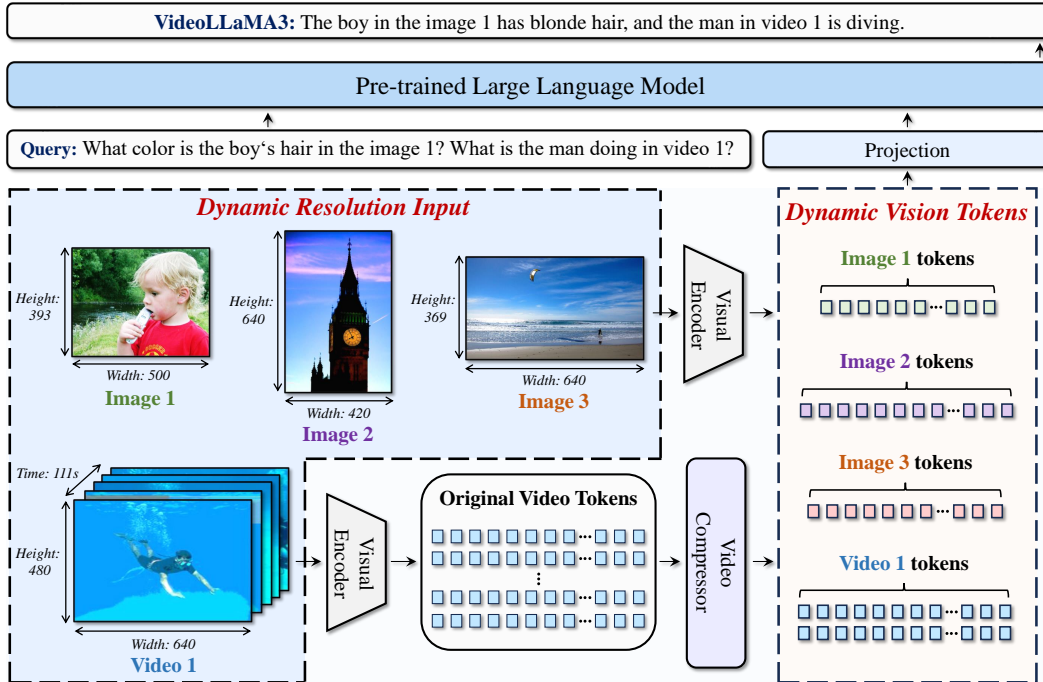


Figure 3: **The overall pipeline of our VideoLLaMA3.** There are two key technical points: ❶ **Any-resolution Vision Tokenization (AVT)**: AVT converts images or videos of any resolution into a set of 1-D token sequences, enabling compatibility with varying amounts of input images and videos of different resolutions, thereby supporting more flexible vision input; ❷ **Differential Frame Pruner (DiffFP)**: Serving as a video compressor, DiffFP eliminates video content with minimal differences between adjacent frames. This approach enhances video processing efficiency, particularly for long-form videos.

- We propose the vision-centric training paradigm. Specifically, we propose to improve video understanding capabilities through large-scale image understanding pretraining.
- We propose two vision-centric framework designs to adapt vision encoders to represent images and videos better.

2 Methodology

As shown in Figure 3, on the model side, VideoLLaMA3 consists of two key technical points: ❶ **Any-resolution Vision Tokenization (AVT)** and ❷ **Differential Frame Pruner (DiffFP)**. When it comes to data, since we propose to improve video understanding capabilities based on image understanding, we also develop a pipeline for constructing high-quality re-captioned image dataset.

2.1 Any-resolution Vision Tokenization

In MLLMs, visual inputs are extracted into vision tokens for multimodal understanding. The common practice [47, 48] is to extract visual inputs with a pre-trained ViT-based vision encoder. The pre-trained vision encoder only receives images with fixed resolutions, which introduces information loss. To alleviate information loss, AnyRes techniques [29, 31, 34] are proposed to split images into patches with fixed resolutions. Although AnyRes techniques increase the number of vision tokens, it is still inflexible and neglects the position relationship within an image when extracting vision tokens. In VideoLLaMA3, we adopt the idea of Any-resolution Vision Tokenization (AVT) [7, 49] to dynamically process images and videos of any resolution. Concretely, we adapt the pre-trained vision encoder (ViT-based architectures)

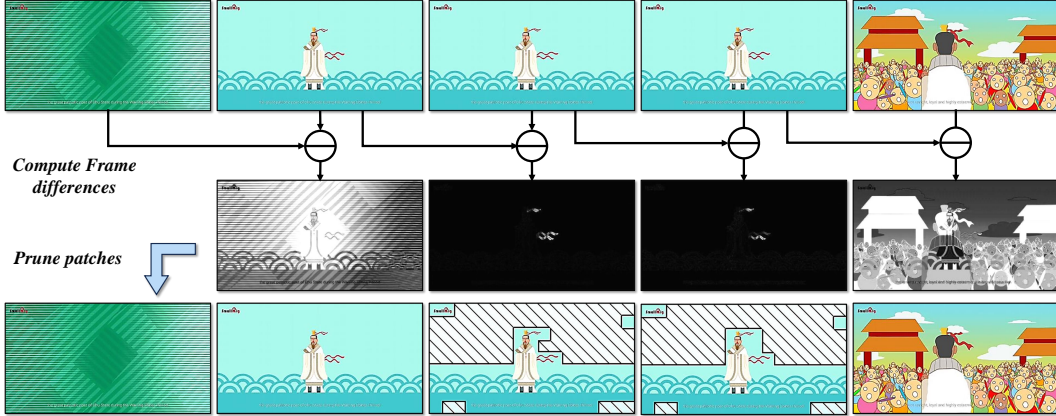


Figure 4: **The calculation flow of our DiffFP.** We prune video tokens based on patch similarities in pixel space, removing patches with smaller distances to the previous frame.

to handle variable resolutions by employing a strategy to replace the absolute position embeddings in ViT with 2D-RoPE [50]. With AVT, images and videos of different resolutions are better represented with more details included in vision tokens. To make the vision encoder compatible with AVT, we fine-tune the vision encoder and the projector in the stage of Vision Encoder Adaptation (i.e., stage #1 in Figure 2) using scene images, document data, and scene images with texts.

2.2 Differential Frame Pruner

For videos, inputs usually have much more tokens than image inputs after tokenization. To reduce the computation demand for videos, we apply a per-frame 2×2 spatial downsampling by bilinear interpolation to limit the context length within a certain range. Besides, considering that videos consist of frames with overlapping content, representing videos by stacking vision tokens from each frame leads to lengthy and redundant tokens. To further reduce the number of tokens of videos, we propose the Differential Frame Pruner (DiffFP) to prune the video tokens. Inspired by RLT [51], we compare the 1-norm distance between temporally consecutive patches within the pixel space. We consider temporally consecutive patches with smaller distances to be redundant, and the later patches can be pruned. Specifically, as shown in Figure 4, we first calculate the 1-norm distance between consecutive frames in the pixel space and then remove patches whose distances fall below a pre-defined threshold. Following RLT [51], we set the default threshold to 0.1.

2.3 Construction of High-Quality Image Re-Caption Dataset

To train our VideoLLaMA3, we constructed a high-quality image re-caption dataset, VL3-Syn7M. All images in this dataset are sourced from COYO-700M [52] and processed using our proposed cleaning pipeline as below:

- 1) **Aspect Ratio Filtering.** We begin by filtering images based on their aspect ratios, removing those with extreme values. This step ensures that the dataset contains images with typical aspect ratios, preventing potential biases during feature extraction. For instance, images that are excessively long or wide may distort the model’s interpretation due to their unusual shapes.
- 2) **Aesthetic Score Filtering.** An aesthetic scoring model is applied to evaluate the visual quality of the images. Based on these scores, images with low aesthetic ratings are discarded. This step eliminates visually poor or poorly composed images, reducing noise and improving the quality of the descriptions generated by the model.
- 3) **Text-Image Similarity Calculation with Coarse Captioning.** The BLIP2 model is used to generate initial captions for images, followed by calculating the text-image similarity

using the CLIP model. Images with low similarity are excluded, as they are likely to contain content that is challenging to describe concisely. This process ensures that the remaining images are both descriptive and interpretable.

4) Visual Feature Clustering. Visual features are extracted using the CLIP vision model, and a k-Nearest-Neighbors (KNN) algorithm is applied for clustering. This method identifies cluster centers in the visual feature space. From each cluster, we select a fixed number of images. This approach ensures diversity within the dataset while maintaining a balanced distribution of semantic categories, improving the model’s ability to generalize across various visual content.

5) Image Re-caption. After filtering and clustering the images, we proceed with detailed re-captioning. Brief captions are generated using InternVL2-8B [31, 53], while the detailed captions are produced with InternVL2-26B [31, 53]. These two types of captions (VL3-Syn7M-short and VL3-Syn7-detailed) are employed at different stages of training to address varying needs.

Through the aforementioned cleaning and re-caption process, we created the VL3-Syn7M dataset, which consists of 7 million image-caption pairs. This high-quality dataset will be a crucial component in training our model, providing a rich and diverse set of images and annotations that support strong performance across a wide range of visual tasks.

3 Training

As illustrated in Figure 3, VideoLLaMA3 consists of four key components: a vision encoder, a video compressor, a projector, and a large language model (LLM). The vision encoder extracts visual tokens and is initialized with the pre-trained SigLIP [54]. To reduce the number of vision tokens representing videos, a video compressor is employed. The projector bridges the features between the vision encoder and the LLM. For the LLM, we utilize Qwen2.5 models [5].

Inspired by previous explorations in MLLMs [7, 25, 29], we develop video understanding capabilities based on strong image understanding foundations. To enable the model with strong image and video understanding capabilities simultaneously, the training of VideoLLaMA3 has four stages: 1) Vision Encoder Adaptation, 2) Vision-Language Alignment, 3) Multi-task Fine-tuning, and 4) Video-centric Fine-tuning. While the first three stages primarily focus on improving image understanding, the final stage is dedicated to video understanding. The details of the training stages are as follows:

1) Vision Encoder Adaptation. In this stage, we fine-tune the vision encoder, which is initialized with the pre-trained SigLIP [54], on a large-scale image dataset. During this stage, the vision encoder is made trainable, while the language decoder remains frozen. This fine-tuning transforms the encoder into a dynamic-resolution processor, enhancing its ability to process images of varying resolutions. Meanwhile, the projector is trained to better align the features of the vision encoder with those of the LLM.

2) Vision-Language Alignment. This stage primarily focuses on introducing multimodal knowledge into the model. During this phase, all parameters are made trainable, enabling both the LLM and the vision encoder to be fine-tuned for integrating multimodal knowledge.

3) Multi-task Fine-tuning. In this stage, we perform instruction fine-tuning using a diverse set of multimodal question-answering data, which includes both image and video-based questions. This step is crucial for improving the model’s ability to follow natural language instructions and enhancing its multimodal understanding. Moreover, this stage lays the foundation for the model’s video understanding capabilities, enabling it to process and analyze temporal information. Also, in this stage, we introduce the video compressor to reduce the number of video tokens.

4) Video-centric Fine-tuning. In this stage, we focus on enhancing the model’s video understanding capabilities. All parameters are unfrozen during this stage. The data used in this stage includes video-text data, image-only data and text-only data.

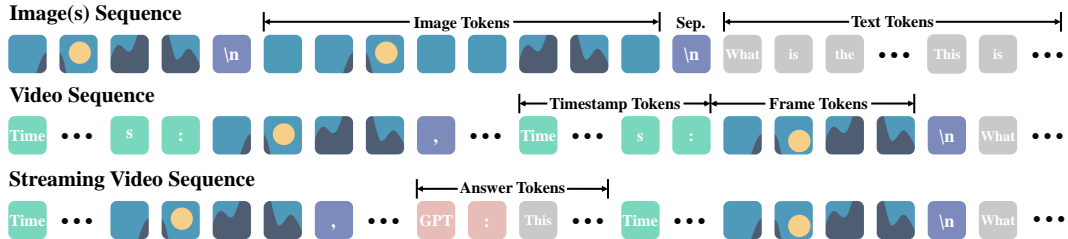


Figure 5: **Data formats for different data types.** ❶ For image sequence, we use “\n” to separate image tokens from different image; ❷ For video sequence, we use “Time: xxs” to indicate timestamps of each frame, “,” to separate different frames, and “\n” to separate tokens from different videos; ❸ For streaming video sequence, videos and texts are organized in an interleaved format.

3.1 Data Format

The data format for images, videos and streaming videos are shown in Figure 5.

Image Sequence. Images are represented as a sequence of tokens, referred to as Image Tokens. The “\n” character is used to separate tokens belonging to different images. Besides, text tokens follow image tokens, separated by “\n”, enabling a mixed representation of image and textual data.

Video Sequence. Frames in a video sequence are represented as Frame Tokens. Before tokens for each frame, a Timestamp Token in the format “Time: xxs” is inserted to denote the time corresponding to that frame. Frames within a video sequence are separated by commas “,”. After the video tokens, “\n” is inserted to separate the video data from any subsequent text tokens, ensuring a clear distinction between the two modalities.

Streaming Video Sequence. For streaming video data, video and text tokens are interleaved in the sequence. Timestamps (i.e., “Time: xxs”) are inserted before the frame tokens, similar to video sequences. To mimic the interactive scenarios of streaming videos, Answer tokens (i.e., “GPT: xxs”) may appear within the sequence to denote contextualized outputs or interactions. The interleaved format ensures a seamless integration of video and textual data streams.

3.2 Data Mixture

Following the principle outlined in LLaVA-OneVision [29], i.e., “quality over quantity”, we conduct rigorous cleaning procedures to guarantee data quality. In this section, we provide a detailed description of the data mixture for each stage, as well as the synthesis and cleaning methods applied to different data subsets.

3.2.1 Vision Encoder Adaptation

Table 1: **Data mixture in vision encoder adaptation stage.**

Task	Dataset	Amount
Scene Image	VL3-Syn7M-short, LLaVA-Pretrain-558k [55], Objects365-Recap [56], SA-1B-Recap [57]	11.84M
Scene Text Image	BLIP3-OCR-Recap [58]	0.93M
Document	pdfa-eng-wds [59], idl-wds [60]	2.80M

The Vision Encoder Adaptation stage is designed to enhance the model’s ability to comprehend a wide range of diverse scenes and improve its feature extraction capacity, with a particular focus on capturing fine-grained information such as objects, regions, and text. As shown in Table 1, the training data in this stage combines scene images and document

recognition images, along with a small portion of scene text images. It should be noted that all data labeled as "Recap" consists of captions generated with InternVL2-8B [31].

For scene images, our data sources include VL3-Syn7M-short, LLaVA-Pretrain-558K [55], Object365 [56], and SA-1B [57]. Notably, the Object365 and SA-1B datasets are included to enhance data diversity, as images in this dataset are mainly complex scenes.

The scene text images are sourced from BLIP3-OCR [58]. Both the brief recaption and the text content within the images are used as captions, and the text content caption following a left-to-right, top-to-bottom pattern across the image.

The document images used in this stage are a subset of pdfa-eng-wds [59] and idl-wds [60]. A total of 2.8 million images were chosen from these two datasets, with the text content of the documents serving as image captions, following the reading order.

3.2.2 Vision-Language Alignment

Table 2: Data mixture in vision-language alignment stage.

Task	Dataset	Amount
Scene Image	VL3-Syn7M-detailed, Objects365-Recap [56], SA-1B-Recap [57], COCO2017-Recap [61], ShareGPT4o [53], TextCaps [62], ShareGPT4V [63], DenseFusion [64], LLaVA-ReCap (LCS-558K) [29]	12.56M
Scene Text Image	Laion-OCR [65], COCO-Text [66], TextOCR [67], BLIP3-OCR-Recap [58], LSVT [68], ReCTS [69]	4.69M
Document	SynthDoG-EN [70], SynthDoG-ZH [70], UReader-TR [71], FUNSD [72], DUDE [73], Vary-600k [74], pdfa-eng-wds [59], idl-wds [60]	2.68M
Chart	Chart-to-Text [75]	0.04M
Fine-grained	Osprey-724K [76], MDVP-Data [77], ADE20K-Recap [78], Object365 [56], Flickr-30K [79], Grand [80]	1.00M
Text-only	Evol-Instruct-143K [81], Infinity-Instruct-code [82], Infinity-Instruct-commonsense [82], Infinity-Instruct-math [82]	6.25M

In this stage, we fine-tune the model using high-quality data. As shown in Table 2, we curate five types of data to cover a wide range of everyday scenarios: scene images, scene text images, documents, charts, and fine-grained data, along with a substantial amount of high-quality text-only data.

For scene images, we include COCO-2017 [66], Object365 [56], SA-1B [57], ShareGPT4o [53], ShareGPT4V [63], DenseFusion [64], and LLaVA-ReCap (LCS-558K) [29]. For Object365, COCO-2017, and SA-1B datasets, we combined the original image annotations with InternVL2-26B [31] to recaption and generate detailed image captions.

The scene text images include a diverse set of Chinese and English scene text recognition datasets. These datasets, such as BLIP3-OCR [58], COCO-Text [66], TextOCR [67], LSVT [68], and ReCTS [69], provide varied examples of text in real-world environments. Furthermore, we filter images from the LAION dataset [65] to include those with clear and readable text, resulting in a collection of 3 million high-quality images, which we term as Laion-OCR dataset. For the Laion-OCR dataset captions, we include both the text content and the corresponding bounding box annotations of the text locations. The caption format is as follows:

{Caption}. The texts in this image are {Text1}<box>[{Bounding Box 1}]</box>, {Text2}<box>[{Bounding Box 2}]</box>, ...

As for document images, we include pdfa-eng-wds [59], idl-wds [60], UReader-TR [71], Vary-600k [74], and SynthDoG [70]. SynthDoG dataset is constructed by generating synthetically accurate document images, avoiding human annotation errors and ensuring precise model

training. Furthermore, we add the handwritten document dataset FUNSD [72] and the complex document dataset DUDE [73]. FUNSD provides annotated handwritten samples for handwriting recognition, while DUDE includes documents with complex layouts, enhancing the model’s ability to handle a variety of document types.

For chart images, since charts share many similarities with documents in terms of content presentation, we only include a limited amount of chart data. These data come from the Chart-to-Text [75] dataset.

For fine-grained images, we construct two types of data: region caption data and grounded caption data. Region caption data describes the content of specific regions within an image. These data are derived and constructed from the Ospery-724K [76], Object365 [56], ADE20K [77], and MDVP-Data [78] datasets. These data help the model to understand the details of the image at the region level. Grounded caption data consist of textual descriptions of objects with corresponding bounding box annotations, primarily constructed from the Flickr-30K [79] and GranD [80] datasets. Both types of data enhance the model’s understanding of images, supporting more accurate object localization and recognition in complex scenes.

3.2.3 Multi-task Fine-tuning

Table 3: Data mixture in massive multi-task fine-tuning stage.

Task	Dataset	Amount
<i>Image & Text Data</i>		
General	LLaVA-SFT-665K [38], LLaVA-OV-SI [29], Cambrian-cleaned [39], Pixmo (docs, cap, points, cap-qa, ask-model-anything) [35]	9.87M
Document	DocVQA [40], Docmatix [41]	1.31M
Chart/Figure	ChartQA [42], MMC_Instruction [83], DVQA [84], LRV_Instruction [85], ChartGemma [86], InfoVQA [87], PlotQA [88]	1.00M
OCR	MultiUI [89], in-house data	0.83M
Grounding	RefCoco [90], VCR [91], in-house data	0.50M
Multi-Image	Demon-Full [92], Contrastive_Caption [93]	0.41M
Text-only	Magpie [94], Magpie-Pro [94], Synthia [95], Infinity-Instruct-subjective [82], NuminaMath [96]	2.21M
<i>Video & Text Data</i>		
General	LLaVA-Video-178K [25], ShareGPT4o-Video [28], FineVideo [97], CinePile [98], ShareGemini-k400 [99], ShareGemini-WebVID [99], VCG-Human [22], VCG-Plus [22], VideoLLaMA2 in-house data, Temporal Grounding in-house data	2.92M

In this stage, we perform instruction tuning with instruction-following data to refine the model’s ability to interpret and follow natural language instructions. This data mixture is designed to cover a wide range of tasks, enabling the model to learn to perform various actions based on instructions across diverse contexts and modalities. Additionally, to activate the model’s video understanding capabilities, we incorporate general video data.

Similar to the vision-language alignment stage, we divide the image data into six distinct groups: general, document, chart/figure, OCR, grounding, and multi-image, as shown in Table 3. Each category targets at a specific aspect of visual understanding, ensuring the model can effectively handle tasks related to different types of visual information. Alongside these visual data categories, we also include a substantial amount of text-only data to improve the model’s ability to handle diverse instruction-following tasks involving both visual and textual inputs.

The general image data includes high-quality datasets, such as LLaVA-SFT-665K [55] and LLaVA-OV-SI [29], which serve as foundational resources for enhancing the model’s scene understanding. We also clean and filter the Cambrian-10M [39] dataset. Furthermore, we

incorporate meaningful data from the Pixmo dataset [35], including tasks such as document analysis, caption generation, and counting. These scene images cover a wide range of tasks, including captioning, counting, document understanding, mathematical reasoning, and *etc.*

For constructing the document and chart/figure datasets, we carefully select high-quality data sources and perform quality cleaning to ensure data reliability. It should be noted that the Docmatix dataset is included as it contains multi-page and diverse documents, crucial for significantly enhancing the model’s ability to understand and long complex document structures and content.

For OCR data, we consider two common cases in real-world scenarios: development scenarios and natural scenarios. For development scenarios, we use the MultiUI dataset [89] to activate the model’s capabilities in understanding and processing text within user interfaces. For natural scenarios, we leverage the Laion-OCR dataset to construct additional instruction-tuning data. The instruction-tuning data for OCR consists of the following five sub-tasks: 1) Text Existence Detection: Determine whether a specific piece of text exists within the image. 2) Text Localization: Locate a specific piece of text within the image and output its bounding box. 3) Text Recognition within a Bounding Box: Given a bounding box, recognize the text contained within it. 4) Text Comparison Between Images: Given two images, determine in which image the specified text appears. 5) Comprehensive Text Detection and Recognition: Detect and recognize all text present in the image.

For grounding images, we select data from established datasets such as RefCOCO [90] and VCR [91], which focus on tasks of grounding visual elements in specific textual descriptions.

For multi-image scenes, we leverage the Demon-Full [92] and Contrastive-Caption [93] datasets. The Demon-Full dataset is particularly valuable as it includes various tasks involving multi-image scenes, such as comparing differences between two images, generating captions for the final image in a comic strip, completing missing text in images with occluded portions, determining whether multiple images belong to the same category, and more. These tasks help the model handle complex scenarios involving multiple images, providing a more comprehensive understanding of how visual information can be interpreted across a series of related images. At the same time, such multi-image data further enhances the model’s video understanding capabilities.

For the video data used in this stage, we incorporate commonly used high-quality video caption datasets, along with a small amount of question-answering data. In addition, we supplement these with high-quality data from VideoLLaMA2 [46] and in-house temporal grounding data. The in-house temporal grounding data specifically focuses on temporal relationships between video frames, enabling the model to grasp the sequence of events and understand the flow of actions across time. These combined data sources contribute to a more robust and nuanced video understanding capability for the model.

3.2.4 Video-centric Fine-tuning

Table 4: Data mixture in video-centric fine-tuning stage.

Task	Dataset	Amount
General Video	LLaVA-Video-178K [25], ShareGPT4o-Video [28], FineVideo [97], CinePile [98], ShareGemini-k400 [99], ShareGemini-WebVID [99], VCG-Human [22], VCG-Plus [22], VideoRefer [100], VideoLLaMA2 in-house data, In-house synthetic data	3.03M
Streaming Video	ActivityNet [101], YouCook2 [102], Ego4D-narration [103], Ego4D-livechat [104]	36.2K
Temporal Grounding	ActivityNet [101], YouCook2 [102], ViTT [105], QuerYD [106], HiREST [107], Charades-STA [108], Moment-10M [109], COIN [110]	0.21M
Image-only	LLaVA-SFT-665K [38], LLaVA-OV-SI [29]	0.88M
Text-only	Magpie [94], Tulu 3 [111]	1.56M

The video-centric fine-tuning stage is designed to tune VideoLLaMA3 to a video specialist and fully unleash its video understanding ability by focusing mainly on large-scale and high-quality video instruction following. We first collect videos with generally annotated caption, question, and answer from multiple open-source datasets including LLaVA-Video [25], ShareGPT-4o [28], FineVideo [97], CinePile [98], ShareGemini [99], VideoGPT+ [22] and VideoRefer [100]. These about 2.7M video-centric conversations eventually form a dataset across various scenes and tasks to serve as examples for teaching the model to understand complex dynamic and static content in videos.

In addition, we further expand the data scale and strengthen the model by synthesizing dense captions and QAs of specific aspects. Specifically, following the pipeline proposed in [25], we first filter 68K dynamic videos from Panda-70M [112] dataset by optical flow, and then employ Qwen2-VL-72B [7] to generate diverse dense captions and QAs for each video from the aspects of temporal understanding, spatial understanding, object description, and time-order understanding. Finally, 242K question-answer pairs are used for training.

Besides general video-centric conversations, we also introduce the feature of streaming video understanding and temporal grounding to extend the application scenarios of our model. For streaming video understanding, we acquire data from ActivityNet [101], YouCook2 [102], and Ego4D [103], and organize video frames and multiple temporal dense captions in an interleaved manner as described in Section 3.1, aiming at enhancing the ability to understand fine-grained events in video and to sustain multi-turn conversations in streaming video. Since these videos are generally long, we cut them into small segments of up to two minutes according to the time interval of dense captions, and remove clips with overly dense and sparse captions. The synthetic streaming conversation from VideoLLM-Online [104] is also involved. For temporal grounding, we collect 205K data from datasets including ActivityNet [101], YouCook2 [102], ViTT [105], QuerYD [106], HiREST [107], Charades-STA [108], Moment-10M [109], and COIN [110], and directly convert the grounding annotation to text format such as "1.0-2.0 s" for training.

Finally, we employ a certain amount of image-only and text-only data from LLaVA [38], LLaVA-OneVision [29], Magpie [94], and Tulu 3 [111] for mitigating the impact of catastrophic forgetting on the model’s capabilities.

3.3 Implementation Details

In this part, we briefly introduce the implementation details of each training stage. For all stages, we adopt the cosine learning rate scheduler. The warm up ratio of the learning rate is set as 0.03. The maximum token length is set as 16384, while the maximum token length for vision tokens is set as 10240. In the stage of Vision Encoder Adaptation, when training VideoLLaMA3-2B, we initialize the vision encoder with the pre-trained weights of SigLIP [54] and the LLM with the pre-trained weights of Qwen2.5-2B [5]. For VideoLLaMA3-7B, the vision encoder is initialized with the fine-tuned SigLIP weights in VideoLLaMA3-2B and the LLM is initialized with Qwen2.5-7B [5]. The projector is implemented as a two-layer MLP with GELU as the activation function. In this stage, we only train the vision encoder and projector, and their learning rates are set as 1.0×10^{-5} and 1.0×10^{-3} , respectively. For the remaining stages, the learning rates for the LLM, the projector, and the vision encoder are set as 1.0×10^{-5} , 1.0×10^{-5} , 2.0×10^{-6} , respectively. The differential frame pruner is applied in the multi-task fine-tuning stage and the video-centric fine-tuning stage where video data is involved. The threshold to discard similar visual tokens is 0.1. To limit context length, the visual tokens of videos are spatially downsampled after vision encoder by a factor of 2 using bilinear interpolation. The visual tokens of images are only downsampled in the video-centric fine-tuning stage to align with video data. For loading video data, we first sample frames at 1 frame per second using FFmpeg. These frames will be further sampled uniformly if the total number of frames is greater than a certain value, which is set to 180 to accommodate most videos that last less than 3 minutes.

Table 5: Evaluation results of 2B models on image benchmarks. * denotes the reproduced results. The best results are in bold and the second best ones are underlined.

Model	🧠 SmolVLM 2B	👁️ InternVL2.5 2B	🦋 Qwen2-VL 2B	VideoLLaMA3 2B
<i>Document/Chart/Scene Text Understanding</i>				
ChartQA	65.3*	<u>79.2</u>	73.5	79.8
DocVQA _{test}	81.6	88.7	<u>90.1</u>	91.9
InfoVQA _{test}	-	60.9	<u>65.5</u>	69.4
OCRBench	622*	804	767*	<u>779</u>
<i>Math</i>				
MathVista _{testmini}	44.6	<u>51.3</u>	43.0	59.2
MathVision _{test}	6.5*	<u>14.7</u>	12.4	15.5
<i>Multi Image</i>				
MMMU-Pro	17.1*	23.7	<u>26.0</u>	28.6
MMMU _{val}	38.8	<u>43.6</u>	41.1	45.3
BLINK _{test}	42.3*	<u>44.0</u>	43.1*	44.2
<i>Knowledge/General QA</i>				
RealWorldQA	48.8*	60.1	<u>62.9</u>	67.3
AI2D	62.1*	<u>74.9</u>	69.9	78.2
GQA	49.2*	59.5*	<u>59.8*</u>	62.7
MME	1600*	2005*	1872	<u>1901</u>

4 Experiment

4.1 Image-based Evaluation

4.1.1 Baselines

To comprehensively evaluate the image performance of VideoLLaMA3, we compare it against a diverse set of baselines. For the 2B version of the model, we select several strong methods, including SmolVLM [37], InternVL2.5-2B [34], and Qwen2VL-2B [7]. For the 7B model, there are more options available. We choose to compare against Molmo-7B-D [35], InternVL2.5-8B [34], LLaVA-OneVision [29], NVILA [36], and Qwen2VL-8B [7].

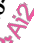
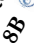
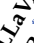
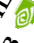
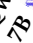
4.1.2 Benchmarks

To evaluate the image recognition and perception capabilities of VideoLLaMA3, we conduct assessments on several representative benchmarks commonly used in Image-LLMs. These benchmarks cover four dimensions: document/chart/scene text understanding, mathematical reasoning, multi-image understanding, and general knowledge QA.

Document/Chart/Scene Text Understanding. To evaluate VideoLLaMA3’s ability to understand various forms of texts in images, including documents, charts, and scene text, we conduct assessments on a range of benchmarks. Specifically, we use: 1) DocVQA [113] for document understanding, which evaluates the model’s ability to process and extract information from text in documents; 2) ChartQA [42] and InfoVQA [113] for chart understanding, assessing the model’s ability to interpret and reason about data presented in graphical forms such as bar charts and line graphs; and 3) OCRBench [114] for scene text image understanding, which tests the model’s capacity to extract and comprehend text from images of real-world scenes.

Mathematical Reasoning. VideoLLaMA3’s mathematical reasoning capabilities are evaluated through the MathVista [115] and MathVision [116] benchmarks. These benchmarks

Table 6: **Evaluation results of 7B models on image benchmarks.** * denotes the reproduced results. † denotes the results retrieved from the official leaderboard. The best results are in **bold** and the second best ones are underlined.

	Molmo-7B-D 7B 	InternVL2.5 8B 	LLaVA-OneVision 7B 	NVILA 8B 	Qwen2-VL 7B 	VideoLLaMA3 7B
<i>Document/Chart/Scene Text Understanding</i>						
ChartQA	84.1	84.8	80.0	<u>86.1</u>	83.0	86.3
DocVQA _{test}	92.2	93.0	87.5	93.7	<u>94.5</u>	94.9
InfoVQA _{test}	72.6	<u>77.6</u>	68.8	70.7	76.5	78.9
OCRBench	-	822	621	676*	845	<u>828</u>
<i>Math</i>						
MathVista _{testmini}	51.6	64.4	63.2	<u>65.4</u>	58.2	67.1
MathVision _{test}	-	<u>19.7</u>	-	11.9*	16.3	26.2
<i>Multi Image</i>						
MMMU-Pro	-	34.3	24.1 [†]	29.5*	31.4*	<u>33.6</u>
MMMU _{val}	45.3	56.0	48.8	<u>49.9</u>	54.1	48.8
BLINK _{test}	-	<u>54.8</u>	48.2	47.0*	43.1*	56.7
<i>Knowledge/General QA</i>						
RealWorldQA	<u>70.7</u>	70.1	66.3	68.6	70.1	72.7
AI2D	93.2	84.5	81.4	<u>92.3</u>	83.0	84.7
GQA	-	-	<u>62.3</u>	-	62.4*	64.9
MME	-	2344	1998	2219	<u>2327</u>	2102

focus on evaluating the model’s ability to reason about and solve mathematical problems presented in visual formats, including text-based mathematical expressions and problem-solving tasks that require visual interpretation.

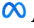
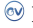
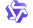
Multi-image Understanding. To assess VideoLLaMA3’s ability to understand and reason about multiple images in conjunction, we evaluate the model on several widely used benchmarks, including MMMU-Pro [117], MMMU [118], and BLINK [119]. These benchmarks test the model’s ability to draw connections between images, handle multiple visual inputs.

General Knowledge QA. Finally, to evaluate VideoLLaMA3’s performance in general question answering, particularly in real-world and complex scenarios, we conduct assessments using several challenging benchmarks. The benchmarks include: 1) RealWorldQA [120], which focuses on answering questions based on realistic images drawn from everyday scenarios, 2) AI2D [121], which evaluates the model’s ability to reason about diagrams and science images, 3) GQA [122], which assesses general question answering with a focus on complex visual reasoning tasks, and 4) MME [123], which includes a wide variety of general knowledge questions that require a deep understanding of visual information.

4.1.3 Evaluation Protocols

When evaluating on benchmarks, we set the temperature as 0.0. The maximum token length is set as the same as the training stage. For benchmarks involving the MCQ, we will give the prompt like “Answer with the option letter from the given choices directly.”. For the benchmarks with short answers, we will give the prompt like “Answer the question with a single word or phrase.”. We follow the original benchmarks to calculate the final

Table 7: **Evaluation results of 2B models on video benchmarks.** * denotes the reproduced results. † denotes the results retrieved from the official leaderboard. The best results are in **bold** and the second best ones are underlined.

Model \ Benchmark	 Apollo 2B	 InternVL2.5 2B	 Qwen2-VL 2B	VideoLLaMA3 2B
<i>General Video Understanding</i>				
VideoMME <i>w/o sub</i>	53.0	51.9	<u>55.6</u>	59.6
VideoMME <i>w/ sub</i>	54.6	54.1	<u>60.4</u>	63.4
MMVU _{val}	-	33.6*	<u>36.5</u> †	39.9
MVBench	-	68.8	63.2	<u>65.5</u>
EgoSchema _{test}	-	<u>58.1</u> *	54.9	58.5
PerceptionTest _{test}	61.0	<u>66.3</u> *	53.9	68.0
ActivityNet-QA	-	<u>54.1</u> *	53.3*	58.2
<i>Long Video Understanding</i>				
MLVU _{dev}	<u>63.3</u>	58.9*	62.7*	65.4
LongVideoBench _{val}	-	<u>52.0</u>	48.7*	57.1
LVBench	-	37.9*	<u>39.4</u> *	41.6
<i>Temporal Reasoning</i>				
TempCompass	60.8	57.7*	<u>62.2</u> *	63.4
NextQA	-	75.6*	<u>77.2</u> *	81.1
Charades-STA	-	-	-	55.5

scores, and we also align our evaluation protocols with other evaluation toolkits, such as Imms-eval [124, 125] and VLMEvalKit [126].

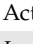
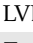
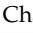


4.1.4 Evaluation Results

We evaluate our VideoLLaMA3 model on the previously mentioned benchmarks. The evaluation results for our 2B model are presented in Table 5. As shown, VideoLLaMA3 demonstrates significant improvements across a range of tasks compared to prior models. For example, in OCR benchmarks such as InfoVQA, VideoLLaMA3 achieves a performance score of 69.4%, compared to the previous best score of 65.5%. In mathematical reasoning tasks, such as MathVista, our 2B model scores 59.2%, surpassing the state-of-the-art method by 7.9%. For multi-image benchmarks like MMMU-Pro, VideoLLaMA3 outperforms the previous top-performing method by 2.6%. In real-world knowledge QA tasks, such as RealWorldQA, VideoLLaMA3 achieves the highest performance with a score of 67.3%, compared to 62.9% from prior methods.

Similarly, we evaluate our larger 7B model on various image benchmarks, with results summarized in Table 6. From the table, it is clear that VideoLLaMA3 consistently outperforms prior models on most benchmarks. Notably, in mathematical reasoning tasks, our 7B model surpasses the previous best by 6.5% on MathVision. In chart understanding tasks, we observe a 1.3% performance improvement over previous methods on InfoVQA. Additionally, in general reasoning tasks like RealWorldQA, VideoLLaMA3 outperforms prior models by 2.0%.

Overall, the results confirm that VideoLLaMA3 provides consistent advancements across a broad range of benchmarks, demonstrating its efficacy and versatility in handling complex tasks, including OCR, mathematical reasoning, and general knowledge. These improvements position VideoLLaMA3 as a powerful tool for real-world applications, advancing the field of multi-modal learning.

Table 8: **Evaluation results of 7B models on video benchmarks.** * denotes the reproduced results. † denotes the results retrieved from the official leaderboard. The best results are in **bold** and the second best ones are underlined.

	Qwen2-VL 7B 	InternVL2.5 8B 	LLaVA-Video 7B 	NVILA 8B 	Apollo 7B 	VideoLLaMA 2.1-7B	VideoLLaMA 3-7B
<i>General Video Understanding</i>							
VideoMME <i>w/o sub</i>	63.3	<u>64.2</u>	63.3	<u>64.2</u>	61.3	54.9	66.2
VideoMME <i>w/ sub</i>	69.0	66.9	69.7	<u>70.0</u>	63.3	56.4	70.3
MMVU _{val}	42.1 [†]	41.1 [†]	42.4*	<u>43.7*</u>	-	39.5 [†]	44.1
MVBench	67.0	72.0	58.6	68.1	-	57.3	<u>69.7</u>
EgoSchema _{test}	66.7	<u>66.2*</u>	57.3	54.3*	-	53.1	63.3
PerceptionTest _{test}	62.3	68.9*	<u>67.9*</u>	65.4*	-	54.9	72.8
ActivityNet-QA	57.4*	58.9*	56.5	<u>60.9</u>	-	53.0	61.3
<i>Long Video Understanding</i>							
MLVU _{dev}	69.8*	69.0*	70.8*	70.6*	<u>70.9</u>	57.4	73.0
LongVideoBench _{val}	55.6 [†]	60.0	58.2	57.7	58.5	-	<u>59.8</u>
LVBench	<u>44.7*</u>	43.2*	41.5*	44.0*	-	36.2	45.3
<i>Temporal Reasoning</i>							
TempCompass	67.9 [†]	<u>68.3*</u>	65.4	69.7*	64.9	56.8	68.1
NextQA	81.2*	85.0*	83.2	82.2	-	75.6	<u>84.5</u>
Charades-STA	-	-	-	-	-	-	60.7

4.2 Video-based Evaluation

4.2.1 Baselines

To comprehensively evaluate the video performance of VideoLLaMA3, we compare it with a diverse set of baseline models. Similar to image evaluation, there are few available models with a 2B parameter size in the community. We select several strong baselines, including Apollo-2B [15], InternVL2.5-2B [34], and Qwen2VL-2B [7]. For the 7B model, we compare it with generalist models such as Qwen2VL-7B [7], InternVL2.5-8B [34], and NVILA [36], as well as specialist models like LLaVA-Video [25], Apollo-7B [15], and our previous generation model, VideoLLaMA2 [46].

4.2.2 Benchmarks

The video understanding capabilities of VideoLLaMA3 are systematically evaluated across three core dimensions: general understanding, temporal reasoning, and long-form video comprehension.

General Video Understanding. We assess VideoLLaMA3’s general video understanding capabilities through established benchmarks: (1) Multi-Choice Video Question Answering (MC-VQA) tasks, including MVBench [27], VideoMME [127], EgoSchema [128], and Perception-Test [129]. (2) Open-Ended Video Question Answering (OE-VQA) tasks, including ActivityNet-QA [130] and VCGBench [26]. This evaluation suite follows the protocol of VideoLLaMA2 [46]. We also run evaluations on MMVU [131] which includes both the task types mentioned above.

Long Video Understanding. To further examine the capacity of VideoLLaMA3 to process and comprehend long-form video content, we assess performance on three long-video

understanding (LVU) benchmarks: (1) MLVU [132]: diverse long-video understanding tasks for videos ranging from 3 minutes to more than 2 hours, (2) LongVideoBench [133]: video reasoning over the referred context within long video-language interleaved inputs, and (3) LVBench [134]: extreme long video understanding.

Video Temporal Reasoning. To assess the temporal awareness and reasoning capabilities of VideoLLaMA3, we conduct evaluations on the following tasks: (1) Temporal Perception and Reasoning tasks, including TempCompass [135] and NextQA [136]; and (2) Temporal Sentence Grounding task on Charades-STA [108] benchmark, with mean Intersection over Union (mIoU) metric.

4.2.3 Evaluation Protocols

We expand the max number of visual tokens to 16K when evaluating our models on video-based benchmarks, ensuring that each frame corresponds to a reasonable number of tokens and the total context length is within the maximum range of the base LLM. The maximum number of frames is set to 180, which is the same as training. For reproducibility, we keep these hyperparameters the same on all benchmarks and disable sampling when decoding.

For general multi-choice question answering evaluation, we follow the official setting to construct the instruction using provided questions and options. An addition prompt like "Answer with the option's letter from the given choices directly" is added to control the model output. In addition, we apply CoT prompt on MMVU benchmark following the official evaluation protocol. For temporal grounding evaluation, we add an extra prompt "Please output the start and end timestamps in seconds" after the question. The numbers in the model response are extracted by regular expression, and then treated as one or multiple time intervals. Based on this strategy, we finally report the mIoU between the ground-truth intervals and the predicted intervals.

4.2.4 Evaluation Results.

Table 7 evaluates the performance of Video Understanding models with 2B model size. VideoLLaMA3 consistently demonstrates competitive results and outperforms baseline methods. In General Video Understanding, VideoLLaMA3 achieves the highest scores on VideoMME w/o sub (59.6%), VideoMME w/ sub (63.4%), ActivityNet-QA (58.2%), PerceptionTest-test (68.0%), MVBench (65.5%), and MMVU (37.6%). On MVBench, it ranks second (65.5%), slightly behind InternVL2.5 2B (68.8%). For Long Video Understanding, VideoLLaMA3 achieves the best performance on all benchmarks: MLVU-dev (65.4%), LongVideoBench-val (57.1%), and LVBench (40.4%), showcasing its superior ability to handle long video content. In Temporal Reasoning, VideoLLaMA3 leads on TempCompass (63.4%), and NextQA (81.1%), and Charades-STA (55.5%). Compared to Apollo-2B, InternVL2.5-2B, and Qwen2-VL-2B, VideoLLaMA3 not only secures the top position in most benchmarks but also demonstrates consistent superiority in tasks requiring comprehensive and long-term video understanding, reinforcing its strong capability across diverse video-related tasks.

As for the VideoLLaMA3-7B model, the results are shown in Table 8. On 7B model size, VideoLLaMA3-7B still exhibits competitive results. For general video understanding, it leads on 5 out of 7 benchmarks, including VideoMME w/o sub, VideoMME w/ sub, PerceptionTest-test, and ActivityNet-QA. On MVBench, it also achieves comparable results to InternVL2.5-8B. For long video understanding, VideoLLaMA3-7B scores the highest on MLVU-dev, and achieves the second best results on LongVideoBench-val and LVBench.

4.3 Case Study

Chart Image Understanding. In Figure 6, we show two cases for chart image understanding. In the first case, VideoLLaMA3 can analyze stock trends and offer some reasonable suggestions for investment. As for the second case, the model can compare the performance of MLLMs and know the tradeoff between the number of parameters and performances.

OCR and Document Understanding. Figure 7 shows two cases for images with texts. In this first example, the model can successfully parse the words in the design image, and offer

some suggestions to make the poster better. In the second image, we ask VideoLLaMA3 to perform OCR task on the given document image. VideoLLaMA3 can successfully recognize the words in the document image, demonstrating the strong performance of VideoLLaMA3 in understanding dense information in images.

Multi-Image Understanding. Figure 8 gives three examples on multi-image understanding tasks. In the first example, VideoLLaMA3 can tell the differences between two types of birds. The second example demonstrates that VideoLLaMA3 is able to locate answers from long documents (even with multiple images) rather than simply parsing words. It is an advanced capability beyond OCR. While in the last example, VideoLLaMA3 can understand storylines from comic strips.

General Image and Video Understanding. Figure 9 demonstrates VideoLLaMA3’s capability in understanding general images, including VQA tasks, answering questions using knowledges and providing videos with captions. Also in Figure 10, we give five cases for video understanding. VideoLLaMA3 can comprehend video content through temporal dimensions, rather than relying solely on inferences from static content.

Long video understanding, temporal grounding, and video-image joint understanding. In Figure 11, we present several cases involving more complex video tasks, including long video grounding, video temporal grounding, and video-image joint understanding. Our VideoLLaMA3 model demonstrates the ability to perform complex long video question-answering tasks. For tasks requiring temporal grounding, our model accurately identifies the specified time. Additionally, for video-image joint understanding, the model effectively captures the relationships between videos and images, enabling it to tackle more intricate tasks.

4.4 Ablation Study

Table 9: Ablation Study on Vision Encoders.

Model	GQA	AI2D	ChartQA	DocVQA _{val}	MME
clip-vit-large-patch14-336 [137]	61.50	56.28	18.32	24.86	1668.41
dfn5B-clip-vit-h-14-378 [138]	62.70	56.87	16.40	23.09	1665.35
siglip-so400m-patch14-384 [54]	62.92	57.12	22.44	31.32	1667.92

In MLLMs, the embeddings of pre-trained vision encoder should be trained to align with embeddings of LLMs. Therefore, the representation performance of vision encoder is crucial to the final performance of MLLMs. In this work, we study the impact of different vision encoders. Specifically, we compare three pre-trained transformer-based vision encoders: CLIP [137], DFN [138], and SigLIP [54]. Due to the computation limitation, we perform the study on the subset of the whole dataset. Also, to investigate the performance of the original pre-trained weights, we fix the weights of vision encoders and keep the visual inputs as the fixed resolution, which is the same as the pretrained resolution of the vision encoder (336×336 for CLIP, 378×378 for DFN, and 384×384 for SigLIP). The training has three stages: 1) Training projector with LLaVA-Pretrain-558K [55]; 2) Tuning all parameters with our captioned COYO data; 3) SFT with LLaVA-SFT-665K [38]. The comparison results are shown in Table. 9. SigLIP outperforms the other two vision encoders, especially in fine-grained understanding tasks involving texts. Based on this ablation study, we choose the pretrained SigLIP as our base vision encoder, and then adapt it to taking dynamic resolutions as inputs.

5 Related Work

Multimodal LLMs for Native Video Understanding. Early video MLLMs primarily relied on sparsely sampled frames and simple connectors, such as MLPs [12, 13, 139], discrete visual tokenizers [140], and Q-formers [141, 142], to link visual encoders with large language



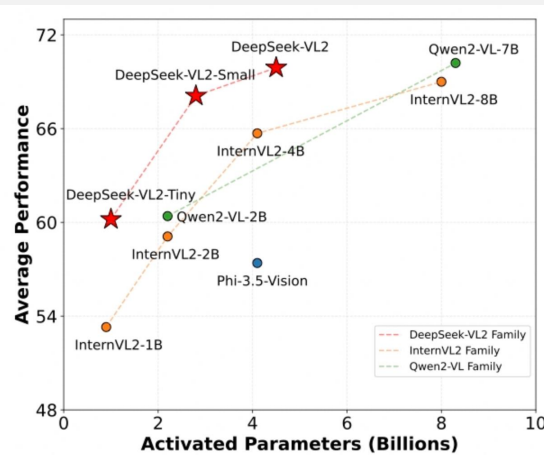
Question: What do you think of this stock? Is it worth holding? Why?



VideoLLaMA 3: The stock shown in this chart *appears to be a volatile and potentially risky investment*. The price has been *fluctuating significantly, with sharp peaks and drops*. This kind of volatility can be indicative of a stock that is not stable or reliable. *While the MACD indicator shows some positive momentum*, suggesting potential for future growth, the overall pattern of the stock price is concerning. *The large red candlesticks at the bottom of the chart indicate significant selling pressure, which could be a sign of investor uncertainty or negative news affecting the stock*. Without more context about the company's fundamentals, industry trends, or market conditions, it's difficult to make a definitive recommendation. However, based solely on the price action shown in this chart, I would advise caution. The stock seems to be experiencing high volatility, which increases risk for investors. If you're considering holding this stock, I'd suggest doing thorough research and possibly diversifying your portfolio to mitigate risk. *It might be better to wait for more stable price action or consider alternative investments that appear less volatile*.



Question: Which model do you think is the strongest? Why?



VideoLLaMA 3: The strongest model in the graph is Qwen2-VL-7B, as it *has the highest average performance and the most activated parameters*. This suggests that it has a *more powerful architecture* and is better at processing visual information.



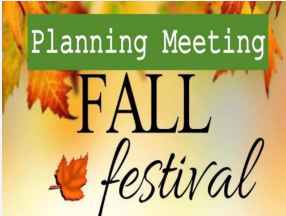
Question: Which model do you think is the optimal? Why?



VideoLLaMA 3: The optimal model in the graph is DeepSeek-VL2, *as it has a good balance between average performance and activated parameters*. It performs well without being overly complex, *making it a good choice for applications that require high performance but also need to be efficient*.

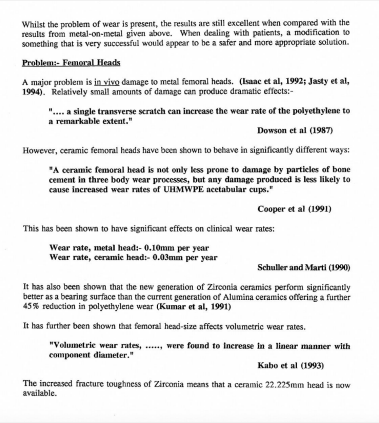
Figure 6: Case study of chart images understanding.

Question: What is this image designed for? What do you think are the shortcomings of the design of this poster? How can it be improved?



VideoLLaMA 3: This image appears to be designed as a poster for a Fall Festival planning meeting. The main text “Planning Meeting” is prominently displayed, followed by “Fall Festival” in a larger, *more decorative font*. The design has some shortcomings: 1. *Lack of visual hierarchy*: The text doesn’t clearly differentiate between the meeting and the festival. 2. *Insufficient information*: It doesn’t provide details about the date, time, location, or purpose of the meeting. 3. *Limited color palette*: The use of green and orange leaves on a white background might not be very eye-catching or memorable. *To improve this design*: 1. Add a clear visual hierarchy by using different font sizes and styles for “Planning Meeting” and “Fall Festival.” 2. Include essential details like date, time, and location below the main text. 3. Use a more vibrant color scheme that stands out, perhaps incorporating autumn colors like reds, yellows, and oranges. 4. *Consider adding a small icon or graphic related to fall (e.g., a pumpkin or acorn)* to make it more visually appealing. 5. Ensure the text is legible and easy to read, especially from a distance. By addressing these issues, the poster would become more effective at communicating important information and attracting attention to the Fall Festival planning meeting.

Question: Please parse all the text in the document, from left to right, from top to bottom.



VideoLLaMA 3: Here is the text from the document: *Whilst the problem of wear is present, the results are still excellent when compared with the results from metal-on-metal given above. When dealing with patients, a modification to something that is very successful would appear to be a safer and more appropriate solution.*

Problem:- Femoral Heads

A major problem is in vivo damage to metal femoral heads. (Isaac et al, 1992; Jasty et al, 1994). Relatively small amounts of damage can produce dramatic effects:-

"... a single transverse scratch can increase the wear rate of the polyethylene to a remarkable extent." Dowson et al (1987)

However, ceramic femoral heads have been shown to behave in significantly different ways:

"A ceramic femoral head is not only less prone to damage by particles of bone cement in three body wear processes, but any damage produced is less likely to cause increased wear rates of UHMWPE acetabular cups." Cooper et al (1991)

This has been shown to have significant effects on clinical wear rates:

Wear rate, metal head:- 0.10mm per year
Wear rate, ceramic head:- 0.03mm per year Schuller and Marti (1990)

It has also been shown that the new generation of Zirconia ceramics perform significantly better as a bearing surface than the current generation of Alumina ceramics offering a further 45% reduction in polyethylene wear (Kumar et al, 1991)

It has further been shown that femoral head-size affects volumetric wear rates.

"Volumetric wear rates,, were found to increase in a linear manner with component diameter." Kabo et al (1993)

The increased fracture toughness of Zirconia means that a ceramic 22.225mm head is now available.

Whilst the problem of wear is present, the results are still excellent when compared with the results from metal-on-metal given above. When dealing with patients, a modification to something that is very successful would appear to be a safer and more appropriate solution.

Problem:- Femoral Heads

A major problem is in vivo damage to metal femoral heads. (Isaac et al, 1992; Jasty et al, 1994). Relatively small amounts of damage can produce dramatic effects:

"... a single transverse scratch can increase the wear rate of the polyethylene to a remarkable extent."

Dowson et al (1987)

However, ceramic femoral heads have been shown to behave in significantly different ways:

"A ceramic femoral head is not only less prone to damage by particles of bone cement in three body wear processes, but any damage produced is less likely to cause increased wear rates of UHMWPE acetabular cups."

Cooper et al (1991)

This has been shown to have significant effects on clinical wear rates:

Wear rate, metal head:- 0.10mm per year

Wear rate, ceramic head:- 0.03mm per year

Schuller and Marti (1990)

It has also been shown that the new generation of Zirconia ceramics perform significantly better as a bearing surface than the current generation of Alumina ceramics offering a further 45% reduction in polyethylene wear (Kumar et al, 1991)

It has further been shown that femoral head-size affects volumetric wear rates...

Figure 7: Case study of OCR and document images.



Question: Please elaborate on the difference between these two animals.

Animal1:



Animal2:



VideoLaMA 3: Animal1 is brown with a white neck while animal2 is yellow with black stripes and red on the back.



Question: What are the requirements for team grading and when will it take place?

MARIST COLLEGE CANBERRA CRICKET CLUB

2021/2022 SEASON

The Marist College Canberra Cricket Club aims to provide an effective cricket support infrastructure for Marist students to play cricket and to engage their families and interest in the game.

Goals

- To provide an excellent standard of player coaching and player development.
- To provide a high level of participation of students who allow us to play cricket.
- To provide good fun cricket for the college cricket teams.
- To provide all players with the opportunity to be involved in all levels of the game.
- To engage Marist students in the enjoyment of the sport and ensure that students and players play the game year-round.
- To develop all teams, and especially the Jnr XI, into successful contributors where all boys are able to contribute to sports success.
- To keep the College and general community informed of the progress of the cricket teams.
- To ensure the level of participation matches game scheduling criteria so they will want to play the game.
- To ensure all players are safe and enjoy the game.
- To ensure all players are able to participate in the game.
- To ensure that funding is available to address all other goals.

How to join:

Students are required to be registered members of Marist College. Please check the appropriate box you complete for on-line registration form (you can also tick, exempt, exempt to under 10 years). You must also complete the on-line registration form for the season. Please note that it is a mandatory requirement for all returning players to be working with SafeSport. They can:

- Be exempt from the 2021-22 season will open in July. Registration (and payment) may be completed on-line through MyCric.

What are the costs for the 2021-22?

Based on prices in 2020 - 2021 (a Year 12): Depending on date of birth (see table below), you can expect to pay for your cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10). If you have any equipment from last year, you will need to bring it with you. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year.

Age Group	Days of each week
Under 7	Monday - Friday
Under 8	Monday - Friday
Under 9	Monday - Friday
Under 10	Monday - Friday
Under 11	Monday - Friday
Under 12	Monday - Friday
Under 13	Monday - Friday
Under 14	Monday - Friday
Under 15	Monday - Friday
Under 16	Monday - Friday
Under 17	Monday - Friday
Under 18	Monday - Friday
Under 19	Monday - Friday
Under 20	Monday - Friday
Under 21	Monday - Friday

Join us at Marist or Marist 1st or Marist 2nd girls' members (for more you are encouraged to play for Marist 2nd) contact us on: 0262 520 1010.

What about trying for teams?

Players of all grades 10 to Year 12 and only those players who wish to try out for Division 1 teams can attend the trials. For more information, please contact the club on 0262 520 1010.

The season will commence in mid-October and will conclude in mid-March. The trial will take place in late August for the Canberra holidays and commencing at the end of January 2022. The season will finish in March 2022.

What about teaming?

With a scheduled playing time on six hours and fifteen to the start of the season, an indication is provided below. Please note that this is an indication and does not guarantee a playing time.

Age Group	Friday	Saturday	Sunday
Under 7 (Division 1)	8.30 - 10.00		
Under 8 (Division 1)	8.30 - 10.00		
Under 9 (Division 1)	8.30 - 10.00		
Under 10 (Division 1)	8.30 - 10.00		
Under 11 (Division 1)	8.30 - 10.00		
Under 12 (Division 1)	8.30 - 10.00		
Under 13 (Division 1)	8.30 - 10.00		
Under 14 (Division 1)	8.30 - 10.00		
Under 15 (Division 1)	8.30 - 10.00		
Under 16 (Division 1)	8.30 - 10.00		
Under 17 (Division 1)	8.30 - 10.00		
Under 18 (Division 1)	8.30 - 10.00		
Under 19 (Division 1)	8.30 - 10.00		
Under 20 (Division 1)	8.30 - 10.00		
Under 21 (Division 1)	8.30 - 10.00		
Under 22 (Division 1)	8.30 - 10.00		

What are the requirements for teaming?

All players are required to be registered members of Marist College. To attend playing time, you must attend the College on your weekend sessions. The club's playing time is being shared and used by 157. Please check with the club's management for more information. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year.

What are the requirements for teaming?

All players are required to be registered members of Marist College. To attend playing time, you must attend the College on your weekend sessions. The club's playing time is being shared and used by 157. Please check with the club's management for more information. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year.

What are the requirements for teaming?

All players are required to be registered members of Marist College. To attend playing time, you must attend the College on your weekend sessions. The club's playing time is being shared and used by 157. Please check with the club's management for more information. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year.

What are the requirements for teaming?

All players are required to be registered members of Marist College. To attend playing time, you must attend the College on your weekend sessions. The club's playing time is being shared and used by 157. Please check with the club's management for more information. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year. Please note that you will need to bring your own cricket gear (Year 12: \$150, Year 11: \$100, Year 10: \$80, Year 9: \$60, Year 8: \$40, Year 7: \$20, Year 6: \$10) if you do not have any equipment from last year.



VideoLaMA 3: The team grading process will take place in August/September 2021. Players are not graded in Under 10s, and only those players who wish to try out for Division 1 teams need to attend the trials.



Question: With the narratives paired with the initial images, how would you conclude the story using the last picture?



Caption#1: Jack and Lizzy enjoying their holiday at the park with other kids around. But they made mistake by not informing their parents about their whereabouts.



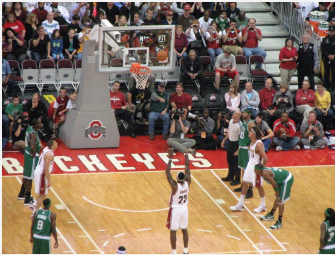
Caption#2: Their parents Martha and Jeff were really tensed about their kids. Jeff was really angry on Martha about not following the kids properly. The argument leads to some bad misunderstanding. They decided to search their kids outside.



VideoLaMA 3: Jack and Lizzy were playing in the park. They were enjoying their time. Martha and Jeff found them and they were very happy. They were relieved that they are safe.

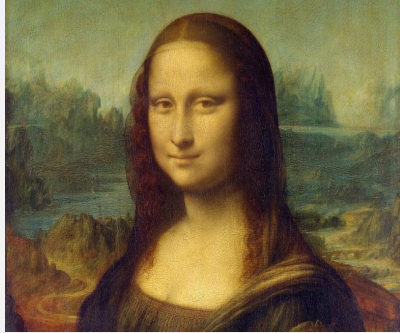
Figure 8: Case study of multi-image understanding.

Question: What part of a basketball game is this?



VideoLLaMA 3: This appears to be a *free throw situation* in a basketball game. *The player wearing number 23 is standing at the free throw line*, which is typically located 15 feet from the basket. He's holding the ball and preparing to take a shot, with his arms raised above his head in the classic free throw stance. The presence of other players on the court, including those wearing green jerseys, suggests this is during an actual game rather than a practice session. The crowd visible in the background further supports that this is taking place during an official match. Free throws are a crucial part of basketball, often occurring after fouls or technical violations. This moment captures the tension and focus as the player attempts to score points for his team.

Question: Discuss the historical impact and the significance of this painting in the art world.



VideoLLaMA 3: *The Mona Lisa, painted by Leonardo da Vinci in the early 16th century*, is one of the most iconic and influential paintings in the history of art. It represents a pinnacle of Renaissance art, characterized by its lifelike representation, subtle use of light and shadow, and the enigmatic smile of the subject. *The painting's historical impact lies in its technical mastery and its enduring mystery*, which has captivated viewers for centuries. The Mona Lisa has been a subject of fascination, inspiring countless imitations, parodies, and references in popular culture. Its significance extends beyond its artistic value; *it symbolizes the Renaissance's humanistic ideals and the pursuit of knowledge and beauty*. The painting's fame has also contributed to the establishment of museums as cultural institutions, with the Louvre Museum in Paris being one of the world's most visited due in part to the Mona Lisa.

Question: Please describe this video in detail.



VideoLLaMA 3: The video begins with *a view of space*, showcasing *a spacecraft* flying through the *vastness of the cosmos*. The scene transitions to *a group of puppies dressed in astronaut suits, sitting inside the spacecraft*. They appear to be on a mission, looking out into the unknown. The camera then *captures a stunning view of Earth from space*, with the puppies still visible inside the spacecraft. The video wraps up with a *close-up of the puppies, highlighting their adorable and curious expressions* as they explore the wonders of space.

Figure 9: Case study of images with general knowledge.

Question: What objects appear in the video? What is the order and position of objects?



VideoLLaMA 3: The video shows a *backyard with a wooden planter box, a patio area with white outdoor furniture, and a house with large glass sliding doors*. The camera pans to the *right*, revealing more of the backyard and patio area. Inside the house, there is a *living room with a couch, a TV mounted on the wall, and a kitchen area* with white cabinets and stainless steel appliances.

Question: Which keyboard key is the last to disappear?



VideoLLaMA 3: The key with the number '*1*' is the last to disappear.

Question: What is unusual in the video?



VideoLLaMA 3: The unusual aspect of the video is that a large bear and two smaller bears are sitting at a wooden table, eating sushi in a serene forest setting.

Question: Describe this video in detail.



VideoLLaMA 3: In the video, there is a *unicorn made out of bubbles*. The unicorn is *blowing bubbles out of its mouth*. The bubbles are *floating in the air*. There is also a *yellow rubber duck* nearby.

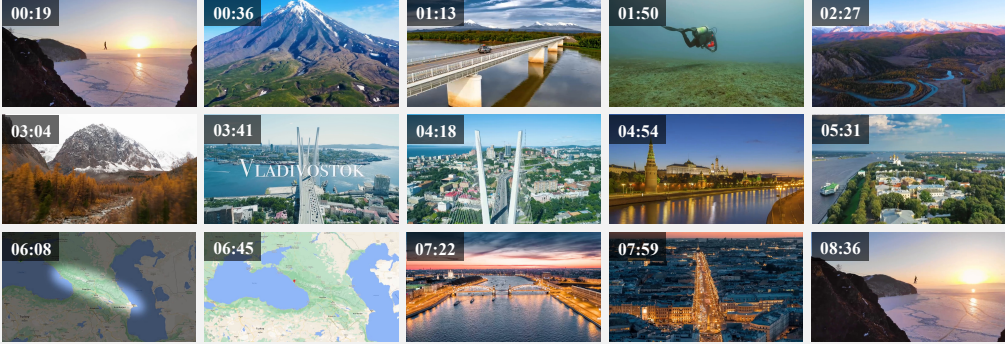
Question: Who won the competition?



VideoLLaMA 3: The athlete *with the number 184 on his jersey* won the competition.

Figure 10: Case study of video understanding.

Question: Please describe the video in detail.

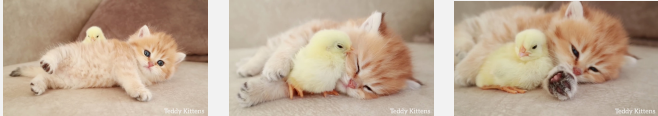
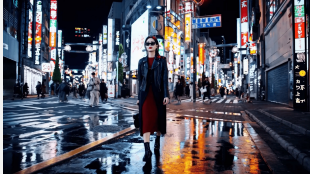


VideoLLaMA 3: The video shows a montage of various locations in Russia. It starts with *a person walking along a tightrope over icy water*, followed by *scenes of landscapes and cityscapes*. The video also includes an advertisement for a *Russian travel destination*, showcasing the country's natural beauty. The montage is accompanied by music, creating a visually appealing experience.

Question: When did the man pour the cola into the cup? Please output the start and end timestamps.



VideoLLaMA 3: The man poured the cola into the cup from 23.1 to 36.8.

Question: <video> What is the connection between the video and the following image? <image>

VideoLLaMA 3: The video and the image *do not seem to be related*. The video features *a kitten and a chick cuddling together*, while the image shows *a woman walking down a city street at night*.

Figure 11: Case study of long video understanding, temporal grounding, and video-image joint understanding.

models. Subsequent models propose various methods to overcome token limitations and support long-form video understanding. For example, [14] directly extends the context window of LLMs to achieve long video understanding, while others [25, 143–150] introduce video token compression techniques that perform pooling across spatial, temporal, or both dimensions. While most approaches utilize image-based encoders [15–19, 28, 151–153], some incorporate video-specific encoders to better capture temporal dependencies [20–23].

More recent works [23, 24, 154–161] extend beyond visual inputs by incorporating audio, using separate encoders for each modality and integrating them through an LLM decoder. These models leverage joint instruction tuning on video-audio datasets [162–165] to capture interactions between visual and auditory information. Additionally, recent advances in streaming video understanding focus on real-time processing [166–170], employing techniques like adaptive memory and incremental processing for tasks such as live event detection and real-time captioning. Previous works [15, 20, 23, 25, 143, 150] typically follow a training recipe that involves an alignment phase, followed by supervised fine-tuning, with instruction-tuning datasets [19, 25–28] often being video dominant. However, we propose a vision-centric training paradigm to enhance video understanding capabilities by focusing on large-scale image understanding pre-training. This approach leverages high-quality image-text datasets to build robust vision encoders that are then adapted for video tasks.

Multimodal LLMs for General Vision Understanding. Recently, a growing number of general MLLMs have been developed to process both images and videos. While, in principle, models designed to handle multiple images are inherently capable of processing video data, achieving optimal performance requires dedicated training on video-specific datasets.

Previous studies [29, 30, 93, 171] have demonstrated that general MLLMs with robust image understanding capabilities can achieve remarkable performance on video understanding tasks, even with minimal or no dedicated video training data. These findings highlight the effectiveness of task transfer from images to videos, showcasing the models’ strong video comprehension and cross-scenario adaptability.

Furthermore, Qwen2-VL [7] adopts a unified framework for processing both images and videos, enhancing the model’s visual perception capabilities. Models such as Qwen2-VL, InternVL-2 [31], and InternVL-2.5 [34], which scale both model sizes (ranging from 1 billion to 78 billion parameters) and the volume of training data, have achieved highly competitive performance in both image and video understanding tasks. To address the challenges of processing longer video inputs, recent studies [32, 33] have proposed solutions such as adapting model architectures by incorporating a hybrid design of Mamba and Transformer blocks or training with extensive long video datasets to support extended input and output sequences.

Recent studies [4, 172, 173] have integrated text, image, and video modalities with audio and speech modalities to improve models’ video understanding and cross-scenario performance. Additionally, Aria [8] utilizes a fine-grained mixture-of-experts decoder, which enables more efficient training and inference compared to dense decoders when handling multimodal inputs.

6 Discussion, Limitations, and Future Work

6.1 Discussion

The introduction of VideoLLaMA3 marks a significant advancement in the realm of MLLMs, particularly in bridging the gap between image and video understanding. By adopting a vision-centric training paradigm, VideoLLaMA3 leverages the robustness of image-centric data to enhance video comprehension, effectively mitigating the challenge associated with temporal dynamics and the complexity of video data. This approach underscores the inherent value of high-quality image-text datasets, which are more readily available and easier to curate compared to their video-text counterparts. The success of VideoLLaMA3 on diverse benchmarks, including VideoMME, PerceptionTest, MLVU, DocVQA, and MathVista, demonstrates its versatility and efficacy across various multimodal tasks.

The model’s ability to maintain strong performance in both image and video domains highlights the effectiveness of our vision-centric framework designs. Specifically, the dynamic resolution adaptation and vision token compression strategies facilitate a more flexible and efficient representation of visual inputs, enabling the model to handle a wide range of image and video formats with minimal information loss. This flexibility is crucial for real-world applications where visual data can vary significantly in resolution and aspect ratio.

Furthermore, the multi-task fine-tuning stage of our training paradigm contributes to the model’s robust generalization capabilities. By exposing VideoLLaMA3 to a variety of downstream tasks, including interactive question answering and video captioning, the model develops a comprehensive understanding of both static and dynamic visual information. This comprehensive training enables VideoLLaMA3 to excel not only in standard benchmarks but also in specialized tasks that require nuanced comprehension of visual content.

6.2 Limitations

Despite the impressive performance of VideoLLaMA3, several limitations must be acknowledged.

Video Data Quality and Diversity. While leveraging large-scale image-text datasets has proven beneficial, the quality and diversity of video-text datasets remain a constraint. Video data often suffer from lower annotation quality and limited diversity, which can impede the model’s ability to generalize across different video domains and genres.

Real-time Processing. The current model architecture may not be optimized for real-time video processing tasks, which is essential for applications such as autonomous driving and live video analytics. The computational overhead associated with processing high-resolution and lengthy video inputs can hinder real-time performance.

Generalization to Unseen Modalities. While VideoLLaMA3 excels in image and video understanding, its capability to generalize to other modalities, such as audio or speech data, remains unexplored. Integrating additional modalities could further enhance the model’s multimodal comprehension but poses significant challenges in terms of architecture and training.

6.3 Future Work

Building on the foundations laid by VideoLLaMA3, several avenues for future research are proposed to address the identified limitations and further enhance the model’s capabilities.

Enhanced Video-Text Datasets. Investing in the creation and curation of higher quality and more diverse video-text datasets will be crucial. Incorporating annotations that capture nuanced temporal and contextual information can significantly improve the model’s temporal understanding and generalization across different video domains.

Real-time Inference Optimization. Optimizing the model architecture for real-time inference by reducing latency and improving processing speed is essential for applications requiring immediate responses. Techniques such as model acceleration, parallel processing, and efficient tokenization strategies can contribute to achieving real-time performance.

Multimodal Expansion. Extending VideoLLaMA3 to incorporate additional modalities like audio, speech, and sensor data can create a more holistic understanding of multimodal inputs. Research into unified architectures that seamlessly integrate multiple data types will be pivotal in achieving comprehensive multimodal intelligence.

Advanced Post-Training Techniques. Implementing more sophisticated post-training methodologies, such as scaling RL techniques for MLLMs, can further refine VideoLLaMA3’s performance. RLHF and other RL-based approaches can be employed to better align the model’s outputs with human preferences and task-specific requirements. Scaling these RL techniques to accommodate the complexities of multimodal data will enhance the model’s ability to generate more accurate, contextually appropriate, and user-aligned responses, thereby advancing its overall multimodal intelligence.

In summary, while VideoLLaMA3 represents a significant step forward in multimodal AI, addressing its current limitations through targeted research and development will pave the way for even more powerful and versatile models in the future.

References

- [1] OpenAI. Gpt-4o system card, 2024. [1](#)
- [2] Anthropic. The claude 3 model family: Opus, sonnet, haiku.
- [3] Gemini Team Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [4] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. [24](#)
- [5] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. [6](#), [11](#)
- [6] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. [1](#)
- [7] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. [2](#), [4](#), [6](#), [11](#), [12](#), [15](#), [24](#)
- [8] Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Fan Zhou, Chengen Huang, Yanpeng Li, Chongyan Zhu, Xiaoyi Ren, Chao Li, Yifan Ye, Peng Liu, Lihuan Zhang, Hanshu Yan, Guoyin Wang, Bei Chen, and Junnan Li. Aria: An open multimodal native mixture-of-experts model, 2025. [24](#)
- [9] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024.
- [10] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. [2](#)
- [11] KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. [2](#)
- [12] Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023. [17](#)
- [13] Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and Mohamed Elhoseiny. Minigpt4-video: Advancing multimodal llms for video understanding with interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024. [17](#)
- [14] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*, 2024. [24](#)
- [15] Orr Zohar, Xiaohan Wang, Yann Dubois, Nikhil Mehta, Tong Xiao, Philippe Hansen-Estruch, Licheng Yu, Xiaofang Wang, Felix Juefei-Xu, Ning Zhang, et al. Apollo: An exploration of video understanding in large multimodal models. *arXiv preprint arXiv:2412.10360*, 2024. [15](#), [24](#)

- [16] Enxin Song, Wenhao Chai, Tian Ye, Jenq-Neng Hwang, Xi Li, and Gaoang Wang. Moviechat+: Question-aware sparse memory for long video question answering. *arXiv preprint arXiv:2404.17176*, 2024.
- [17] Jiajun Fei, Dian Li, Zhidong Deng, Zekun Wang, Gang Liu, and Hui Wang. Videoccam: Enhancing video-language understanding with causal cross-attention masks for short and long videos. *arXiv preprint arXiv:2408.14023*, 2024.
- [18] Jiajun Liu, Yibing Wang, Hanghang Ma, Xiaoping Wu, Xiaoqi Ma, Xiaoming Wei, Jianbin Jiao, Enhua Wu, and Jie Hu. Kangaroo: A powerful video-language model supporting long-context video input. *arXiv preprint arXiv:2408.15542*, 2024.
- [19] Zuyan Liu, Yuhao Dong, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. Oryx mllm: On-demand spatial-temporal understanding at arbitrary resolution. *arXiv preprint arXiv:2409.12961*, 2024. 24
- [20] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, et al. Internvideo2: Scaling video foundation models for multimodal video understanding. *arXiv preprint arXiv:2403.15377*, 2024. 24
- [21] Raehyuk Jung, Hyojun Go, Jaehyuk Yi, Jiho Jang, Daniel Kim, Jay Suh, Aiden Lee, Cooper Han, Jae Lee, Jeff Kim, et al. Pegasus-v1 technical report. *arXiv preprint arXiv:2404.14687*, 2024.
- [22] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Videogpt+: Integrating image and video encoders for enhanced video understanding. *arxiv*, 2024. 9, 10, 11
- [23] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 2, 24
- [24] Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, Yuxuan Wang, and Chao Zhang. video-salmonn: Speech-enhanced audio-visual large language models. *arXiv preprint arXiv:2406.15704*, 2024. 24
- [25] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 2, 6, 9, 10, 11, 15, 24
- [26] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*, 2024. 15
- [27] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. *arXiv preprint arXiv:2311.17005*, 2023. 15
- [28] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Bin Lin, Zhenyu Tang, et al. Sharegpt4video: Improving video understanding and generation with better captions. *arXiv preprint arXiv:2406.04325*, 2024. 2, 9, 10, 11, 24
- [29] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 2, 3, 4, 6, 7, 8, 9, 10, 11, 12, 24
- [30] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models, 2024. 24

- [31] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023. 2, 3, 4, 6, 8, 24
- [32] Xidong Wang, Dingjie Song, Shunian Chen, Chen Zhang, and Benyou Wang. Longlava: Scaling multi-modal llms to 1000 images efficiently via a hybrid architecture, 2024. 24
- [33] Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, Songyang Zhang, Wenwei Zhang, Yining Li, Yang Gao, Peng Sun, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Hang Yan, Conghui He, Xingcheng Zhang, Kai Chen, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output, 2024. 2, 24
- [34] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling, 2025. 2, 3, 4, 12, 15, 24
- [35] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024. 9, 10, 12
- [36] Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, et al. Nvila: Efficient frontier visual language models. *arXiv preprint arXiv:2412.04468*, 2024. 12, 15
- [37] Hugging Face Team. Smolvlm - small yet mighty vision language model. <https://huggingface.co/blog/smolvlm>, 2023. Accessed: 2025-01-19. 2, 12
- [38] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 9, 10, 11, 17
- [39] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024. 9
- [40] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021. 9
- [41] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding vision-language models: insights and future directions., 2024. 9

- [42] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022. 9, 12
- [43] Kung-Hsiang Huang, Hou Pong Chan, Yi R. Fung, Haoyi Qiu, Mingyang Zhou, Shafiq Joty, Shih-Fu Chang, and Heng Ji. From pixels to insights: A survey on automatic chart understanding in the era of large foundation models, 2024. 2
- [44] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 2
- [45] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. In Yansong Feng and Els Lefever, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 - System Demonstrations, Singapore, December 6-10, 2023*, pages 543–553. Association for Computational Linguistics, 2023. 2
- [46] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 2, 10, 15
- [47] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3, 4
- [48] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 3, 4
- [49] Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Patch n’pack: Navit, a vision transformer for any aspect ratio and resolution. *arXiv preprint arXiv:2307.06304*, 2023. 4
- [50] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 5
- [51] Rohan Choudhury, Guanglei Zhu, Sihan Liu, Koichiro Niinuma, Kris M Kitani, and László Jeni. Don’t look twice: Faster video transformers with run-length tokenization. *arXiv preprint arXiv:2411.05222*, 2024. 5
- [52] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Sae-hoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022. 5
- [53] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiao wen Dong, Hang Yan, Hewei Guo, Conghui He, Zhenjiang Jin, Chaochao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, and Yu Qiao. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *ArXiv*, abs/2404.16821, 2024. 6, 8
- [54] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11941–11952, 2023. 6, 11, 17
- [55] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 7, 8, 9, 17
- [56] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8429–8438, 2019. 7, 8, 9

- [57] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 7, 8
- [58] Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S. Ryoo, Shrikant B. Kendre, Jieyu Zhang, Can Qin, Shu Zhen Zhang, Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Awalganekar, Shelby Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caiming Xiong, and Ran Xu. xgen-mm (blip-3): A family of open large multimodal models. *ArXiv*, abs/2408.08872, 2024. 7, 8
- [59] pdfa-eng-wds. <https://huggingface.co/datasets/pixparse/pdfa-eng-wds>. 7, 8
- [60] idl-wds. <https://huggingface.co/datasets/pixparse/idl-wds>. 7, 8
- [61] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 8
- [62] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 742–758. Springer, 2020. 8
- [63] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 8
- [64] Xiaotong Li, Fan Zhang, Haiwen Diao, Yueze Wang, Xinlong Wang, and Ling-Yu Duan. Densefusion-1m: Merging vision experts for comprehensive multimodal perception. *arXiv preprint arXiv:2407.08303*, 2024. 8
- [65] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 8
- [66] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. Coco-text: Dataset and benchmark for text detection and recognition in natural images. *arXiv preprint arXiv:1601.07140*, 2016. 8
- [67] Amanpreet Singh, Guan Pang, Mandy Toh, Jing Huang, Wojciech Galuba, and Tal Hassner. Textocr: Towards large-scale end-to-end reasoning for arbitrary-shaped scene text. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8802–8812, 2021. 8
- [68] Yipeng Sun, Zihan Ni, Chee-Kheng Chng, Yuliang Liu, Canjie Luo, Chun Chet Ng, Junyu Han, Errui Ding, Jingtuo Liu, Dimosthenis Karatzas, Chee Seng Chan, and Lianwen Jin. Icdar 2019 competition on large-scale street view text with partial labeling - rrc-lsvt. *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1557–1562, 2019. 8
- [69] Xi Liu, Rui Zhang, Yongsheng Zhou, Qianyi Jiang, Qi Song, Nan Li, Kai Zhou, Lei Wang, Dong Wang, Minghui Liao, et al. Icdar 2019 robust reading challenge on reading chinese text on signboard. *arXiv preprint arXiv:1912.09641*, 2019. 8
- [70] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pages 498–517. Springer, 2022. 8

- [71] Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Mingshi Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Feiyan Huang. Ureader: Universal ocr-free visually-situated language understanding with multi-modal large language model. In *Conference on Empirical Methods in Natural Language Processing*, 2023. 8
- [72] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, volume 2, pages 1–6, 2019. 8, 9
- [73] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Dawid Jurkiewicz, Rafał Powalski, Paweł Józiać, Sanket Biswas, Mickaël Coustaty, and Tomasz Stanisławek. Icdar 2023 competition on document understanding of everything (dude). In *International Conference on Document Analysis and Recognition*, pages 420–434, 2023. 8, 9
- [74] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language models. *arXiv preprint arXiv:2312.06109*, 2023. 8
- [75] Shankar Kantharaj, Rixie Tiffany Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. Chart-to-text: A large-scale benchmark for chart summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022. 8, 9
- [76] Yuqian Yuan, Wentong Li, Jian Liu, Dongqi Tang, Xinjie Luo, Chi Qin, Lei Zhang, and Jianke Zhu. Osprey: Pixel understanding with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28202–28211, 2024. 8, 9
- [77] Weifeng Lin, Xinyu Wei, Ruichuan An, Peng Gao, Bocheng Zou, Yulin Luo, Siyuan Huang, Shanghang Zhang, and Hongsheng Li. Draw-and-understand: Leveraging visual prompts to enable mllms to comprehend what you want. *arXiv preprint arXiv:2403.20271*, 2024. 8, 9
- [78] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. 8, 9
- [79] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. 8, 9
- [80] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 8, 9
- [81] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model, 2024. 8
- [82] Beijing Academy of Artificial Intelligence (BAAI). Infinity instruct. *GitHub repository, HuggingFace repository*, 2024. 8, 9
- [83] Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoub, and Dong Yu. Mmc: Advancing multimodal chart understanding with large-scale instruction tuning. *arXiv preprint arXiv:2311.10774*, 2023. 9

- [84] Kushal Kafle, Scott Cohen, Brian Price, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *CVPR*, 2018. 9
- [85] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 9
- [86] Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty. Chartgemma: Visual instruction-tuning for chart reasoning in the wild, 2024. 9
- [87] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706, 2022. 9
- [88] Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. Plotqa: Reasoning over scientific plots. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1527–1536, 2020. 9
- [89] Junpeng Liu, Tianyue Ou, Yifan Song, Yuxiao Qu, Wai Lam, Chenyan Xiong, Wenhui Chen, Graham Neubig, and Xiang Yue. Harnessing webpage uis for text-rich visual understanding, 2024. 9, 10
- [90] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 9, 10
- [91] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 9, 10
- [92] Juncheng Li, Kaihang Pan, Zhiqi Ge, Minghe Gao, Wei Ji, Wenqiao Zhang, Tat-Seng Chua, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. Fine-tuning multimodal llms to follow zero-shot demonstrative instructions. In *The Twelfth International Conference on Learning Representations*, 2024. 9, 10
- [93] Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. Mantis: Interleaved multi-image instruction tuning, 2024. 9, 10, 24
- [94] Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *ArXiv*, abs/2406.08464, 2024. 9, 10, 11
- [95] Migel Tissera. Synthia-70b-v1.2: Synthetic intelligent agent. <https://huggingface.co/migtissera/Synthia-13B>, 2023. 9
- [96] Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*, 13, 2024. 9
- [97] Miquel Farré, Andi Marafioti, Lewis Tunstall, Leandro Von Werra, and Thomas Wolf. Finevideo. <https://huggingface.co/datasets/HuggingFaceFV/finevideo>, 2024. 9, 10, 11
- [98] Ruchit Rawal, Khalid Saifullah, Miquel Farré, Ronen Basri, David Jacobs, Gowthami Somepalli, and Tom Goldstein. Cinepile: A long video question answering dataset and benchmark. *arXiv preprint arXiv:2405.08813*, 2024. 9, 10, 11
- [99] Share. Sharegemini: Scaling up video caption data for multimodal large language models, June 2024. 9, 10, 11

- [100] Yuqian Yuan, Hang Zhang, Wentong Li, Zesen Cheng, Boqiang Zhang, Long Li, Xin Li, Deli Zhao, Wenqiao Zhang, Yueting Zhuang, et al. Videorefer suite: Advancing spatial-temporal object understanding with video llm. *arXiv preprint arXiv:2501.00599*, 2024. 10, 11
- [101] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 706–715, 2017. 10, 11
- [102] Luwei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, number 1, 2018. 10, 11
- [103] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 10, 11
- [104] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18407–18418, 2024. 10, 11
- [105] Gabriel Huang, Bo Pang, Zhenhai Zhu, Clara Rivera, and Radu Soricut. Multimodal pretraining for dense video captioning. *arXiv preprint arXiv:2011.11760*, 2020. 10, 11
- [106] Andreea-Maria Oncescu, Joao F Henriques, Yang Liu, Andrew Zisserman, and Samuel Albanie. Queryd: A video dataset with high-quality text and audio narrations. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2265–2269. IEEE, 2021. 10, 11
- [107] Abhay Zala, Jaemin Cho, Satwik Kottur, Xilun Chen, Barlas Oguz, Yashar Mehdad, and Mohit Bansal. Hierarchical video-moment retrieval and step-captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23056–23065, 2023. 10, 11
- [108] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, pages 5267–5275, 2017. 10, 11, 16
- [109] Long Qian, Juncheng Li, Yu Wu, Yaobo Ye, Hao Fei, Tat-Seng Chua, Yueting Zhuang, and Siliang Tang. Momentor: Advancing video large language model with fine-grained temporal reasoning. *arXiv preprint arXiv:2402.11435*, 2024. 10, 11
- [110] Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019. 10, 11
- [111] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training. 2024. 10, 11
- [112] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiangwei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. *arXiv preprint arXiv:2402.19479*, 2024. 11

- [113] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021. 12
- [114] Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023. 12
- [115] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 12
- [116] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*, 2024. 12
- [117] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*, 2024. 13
- [118] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 13
- [119] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166, 2025. 13
- [120] xai. Realworldqa benchmark. <https://huggingface.co/datasets/xai-org/RealworldQA>, 2024. 13
- [121] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 235–251, 2016. 13
- [122] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 13
- [123] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 13
- [124] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models, 2024. 14
- [125] Li* Bo, Zhang* Peiyuan, Zhang* Kaichen, Pu* Fanyi, Du Xinrun, Dong Yuhao, Liu Haotian, Zhang Yuanhan, Zhang Ge, Li Chunyuan, and Ziwei Liu. Lmms-eval: Accelerating the development of large multimodal models, March 2024. 14
- [126] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 11198–11201, 2024. 14

- [127] Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024. 15
- [128] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36, 2024. 15
- [129] Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, et al. Perception test: A diagnostic benchmark for multimodal video models. *Advances in Neural Information Processing Systems*, 36, 2024. 15
- [130] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9127–9134, 2019. 15
- [131] Yilun Zhao, Lujing Xie, Haowei Zhang, Guo Gan, Yitao Long, Zhiyuan Hu, Tongyan Hu, Weiyuan Chen, Chuhan Li, Junyang Song, Zhijian Xu, Chengye Wang, Weifeng Pan, Ziyao Shangguan, Xiangru Tang, Zhenwen Liang, Yixin Liu, Chen Zhao, and Arman Cohan. Mmvu: Measuring expert-level multi-discipline video understanding, 2025. 15
- [132] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*, 2024. 16
- [133] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding, 2024. 16
- [134] Weihang Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu Huang, Bin Xu, Yuxiao Dong, et al. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024. 16
- [135] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024. 16
- [136] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786, 2021. 16
- [137] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763, 2021. 17
- [138] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023. 17
- [139] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023. 17
- [140] Yang Jin, Zhicheng Sun, Kun Xu, Liwei Chen, Hao Jiang, Quzhe Huang, Chengru Song, Yuliang Liu, Di Zhang, Yang Song, et al. Video-lavit: Unified video-language pre-training with decoupled visual-motional tokenization. *arXiv preprint arXiv:2402.03161*, 2024. 17

- [141] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 17
- [142] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024. 17
- [143] Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, et al. Longvu: Spatiotemporal adaptive compression for long video-language understanding. *arXiv preprint arXiv:2410.17434*, 2024. 24
- [144] Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. Pllava: Parameter-free llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*, 2024.
- [145] Peng Jin, Ryuichi Takano, Wancai Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13700–13710, 2024.
- [146] Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin Dehghan. Slowfast-llava: A strong training-free baseline for video large language models. *arXiv preprint arXiv:2407.15841*, 2024.
- [147] Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. Longvlm: Efficient long video understanding via large language models. In *European Conference on Computer Vision*, pages 453–470. Springer, 2025.
- [148] Zhende Song, Chenchen Wang, Jiamu Sheng, Chi Zhang, Gang Yu, Jiayuan Fan, and Tao Chen. Moviellm: Enhancing long video understanding with ai-generated movies. *arXiv preprint arXiv:2403.01422*, 2024.
- [149] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In *European Conference on Computer Vision*, pages 323–340. Springer, 2025.
- [150] Reuben Tan, Ximeng Sun, Ping Hu, Jui-hsien Wang, Hanieh Deilamsalehy, Bryan A Plummer, Bryan Russell, and Kate Saenko. Koala: Key frame-conditioned long video-llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13581–13591, 2024. 24
- [151] Jiawei Wang, Liping Yuan, and Yuchen Zhang. Tarsier: Recipes for training and evaluating large video description models, 2024. 24
- [152] Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024.
- [153] Shimin Chen, Xiaohan Lan, Yitian Yuan, Zequn Jie, and Lin Ma. Timemarker: A versatile video-llm for long and short video understanding with superior temporal localization ability. *arXiv preprint arXiv:2411.18211*, 2024. 24
- [154] Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. In *European Conference on Computer Vision*, pages 146–164. Springer, 2025. 24
- [155] Fangxun Shu, Lei Zhang, Hao Jiang, and Cihang Xie. Audio-visual llm for video understanding. *arXiv preprint arXiv:2312.06720*, 2023.

- [156] Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- [157] Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*, 2023.
- [158] Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li, Dongxu Li, Shafiq Joty, Ran Xu, Silvio Savarese, Caiming Xiong, and Juan Carlos Niebles. X-instructblip: A framework for aligning x-modal instruction-aware representations to llms and emergent cross-modal reasoning. *arXiv preprint arXiv:2311.18799*, 2023.
- [159] Jiaming Han, Kaixiong Gong, Yiyuan Zhang, Jiaqi Wang, Kaipeng Zhang, Dahua Lin, Yu Qiao, Peng Gao, and Xiangyu Yue. Onellm: One framework to align all modalities with language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26584–26595, 2024.
- [160] Shoubin Yu, Jaehong Yoon, and Mohit Bansal. Crema: Multimodal compositional video reasoning via efficient modular adaptation and fusion. *arXiv preprint arXiv:2402.05889*, 2024.
- [161] Yunlong Tang, Daiki Shimada, Jing Bi, and Chenliang Xu. Avicuna: Audio-visual llm with interleaver and context-boundary alignment for temporal referential dialogue. *arXiv preprint arXiv:2403.16276*, 2024. 24
- [162] Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. Avqa: A dataset for audio-visual question answering on videos. *Proceedings of the 30th ACM International Conference on Multimedia*, 2022. 24
- [163] Huda AlAmri, Vincent Cartillier, Abhishek Das, Jue Wang, Stefan Lee, Peter Anderson, Irfan Essa, Devi Parikh, Dhruv Batra, Anoop Cherian, Tim K. Marks, and Chiori Hori. Audio visual scene-aware dialog. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7550–7559, 2019.
- [164] Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: Enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios, 2024.
- [165] Sicong Leng, Yun Xing, Zesen Cheng, Yang Zhou, Hang Zhang, Xin Li, Deli Zhao, Shijian Lu, Chunyan Miao, and Lidong Bing. The curse of multi-modalities: Evaluating hallucinations of large multimodal models across language, visual, and audio. *arXiv preprint arXiv:2410.12787*, 2024. 24
- [166] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-based real-time understanding for long video streams. *arXiv preprint arXiv:2406.08085*, 2024. 24
- [167] Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video understanding with large language models. *arXiv preprint arXiv:2405.16009*, 2024.
- [168] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18407–18418, 2024.
- [169] Pan Zhang, Xiao wen Dong, Yuhang Cao, Yuhang Zang, Rui Qian, Xilin Wei, Lin Chen, Yifei Li, Junbo Niu, Shuangrui Ding, Qipeng Guo, Haodong Duan, Xin Chen, Han Lv, Zheng Nie, Min Zhang, Bin Wang, Wenwei Zhang, Xinyue Zhang, Jiaye Ge, Wei Li, Jingwen Li, Zhongying Tu, Conghui He, Xingcheng Zhang, Kai Chen, Yuanbo Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2.5-omnilive: A comprehensive multimodal system for long-term streaming video and audio interactions. 2024.

- [170] Jihao Liu, Zhiding Yu, Shiyi Lan, Shihao Wang, Rongyao Fang, Jan Kautz, Hongsheng Li, and Jose M. Alvarez. Streamchat: Chatting with streaming video. 2024. 24
- [171] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models, 2024. 24
- [172] Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, Ran He, Rongrong Ji, Yunsheng Wu, Caifeng Shan, and Xing Sun. Vita: Towards open-source interactive omni multimodal llm, 2024. 24
- [173] Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Yangze Li, Zuwei Long, Heting Gao, Ke Li, Xiawu Zheng, Rongrong Ji, Xing Sun, Caifeng Shan, and Ran He. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction, 2025. 24