

InternLM-XComposer2.5-OmniLive: A Comprehensive Multimodal System for Long-term Streaming Video and Audio Interactions

Pan Zhang^{*1}, Xiaoyi Dong^{*1,2}, Yuhang Cao^{*1}, Yuhang Zang^{*1}, Rui Qian^{*1,2†}, Xilin Wei^{1,3†}, Lin Chen^{1,4†}, Yifei Li^{1,5†}, Junbo Niu^{1,6†}, Shuangrui Ding^{1,2†}, Qipeng Guo¹, Haodong Duan¹, Xin Chen¹, Han Lv¹, Zheng Nie¹, Min Zhang¹, Bin Wang¹, Wenwei Zhang¹, Xinyue Zhang¹, Jiaye Ge¹, Wei Li¹, Jingwen Li¹, Zhongying Tu¹, Conghui He⁷, Xingcheng Zhang⁷, Kai Chen¹, Yu Qiao¹, Dahua Lin^{1,2}, Jiaqi Wang^{1,✉}

¹Shanghai Artificial Intelligence Laboratory, ²The Chinese University of Hong Kong,

³Fudan University, ⁴University of Science and Technology of China,

⁵Tsinghua University, ⁶Beihang University, ⁷SenseTime Group

internlm@pjlab.org.cn

Abstract

Creating AI systems that can interact with environments over long periods, similar to human cognition, has been a longstanding research goal. Recent advancements in multimodal large language models (MLLMs) have made significant strides in open-world understanding. However, the challenge of continuous and simultaneous streaming perception, memory, and reasoning remains largely unexplored. Current MLLMs are constrained by their sequence-to-sequence architecture, which limits their ability to process inputs and generate responses simultaneously, akin to being unable to think while perceiving. Furthermore, relying on long contexts to store historical data is impractical for long-term interactions, as retaining all information becomes costly and inefficient. Therefore, rather than relying on a single foundation model to perform all functions, this project draws inspiration from the concept of the **Specialized Generalist AI** and introduces disentangled streaming perception, reasoning, and memory mechanisms, enabling real-time interaction with streaming video and audio input. The proposed framework **InternLM-XComposer2.5-OmniLive (IXC2.5-OL)** consists of three key modules: (1) **Streaming Perception Module**: Processes multimodal information in real-time, storing key details in memory and triggering reasoning in response to user queries. (2) **Multi-modal Long Memory Module**: Integrates short-term and long-term memory, compressing short-term memories into long-term ones for efficient retrieval and improved accuracy. (3) **Reasoning Module**: Responds to queries and executes reasoning tasks, coordinating with the perception and memory modules. This project simulates human-like cogni-

* indicates equal contribution. † indicates interns at IXCLab, Shanghai AI Laboratory

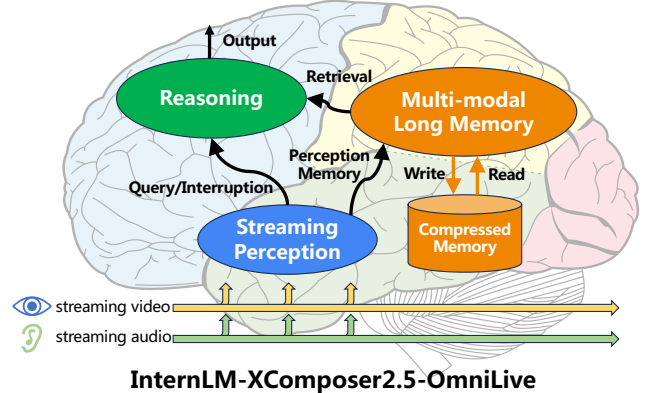


Figure 1. Inspired by human-like cognition and Specialized Generalist AI, we introduce InternLM-XComposer2.5-OmniLive (IXC2.5-OL), a system that facilitates real-time interaction with: (1) a **streaming perception** module supports streaming video and audio inputs; (2) a **multi-modal long memory** module that compresses short-term memory into long-term memory; and (3) a **reasoning** module that answers queries based on retrieved memories.

tion, enabling multimodal large language models to provide continuous and adaptive service over time. All code and models of **InternLM-XComposer2.5-OmniLive (IXC2.5-OL)** are publicly available at <https://github.com/InternLM/InternLM-XComposer/tree/main/InternLM-XComposer-2.5-OmniLive>.

1. Introduction

The goal of developing AI systems [55] that can understand and interact with environments over long periods, akin to human cognition, has been a central focus of research for

decades. The rise of large-scale data corpora [54, 69, 95, 112] and multimodal large language models [83, 84, 107] has driven significant advances in free-form multimodal question answering. Recent developments, such as Mini-Omni [123], VideoLLM-Online [12], and VITA [38], have made notable strides toward enabling more natural and immersive online interactions. However, challenges persist in creating systems capable of continuous interaction due to the intrinsic limitations of a single decoder-only large language model architecture.

Existing architectures [12, 38, 123, 149] encounter significant limitations in real-time and long-term streaming perception, reasoning, and memory. The sequence-to-sequence decoder-only architecture used in current MLLMs forces a switch between perception (e.g., seeing and hearing) and thinking, limiting the simultaneous processing of inputs and outputs. Additionally, existing works [33, 118, 145] rely on the integration of multimodal memories within context windows. The reliance on long contexts to store historical information proves impractical for long-term use, especially in scenarios requiring continuous AI assistance. Multimodal data, like video streams, can quickly accumulate millions of tokens within a few hours, making it impractical to maintain context over multiple days of service. The cost and inefficiency of storing all historical clues within the context further limit the system’s capacity to provide continuous and long-term service. In contrast, the human brain can effortlessly integrate perception and cognition, preserving long-term multimodal memories. This is believed to be closely related to the functional partitioning design of the human brain cortex, where different areas of the cortex are responsible for distinct tasks, such as perception, memory, and cognition.

Inspired by the paradigm of Specialized Generalist AI [146], we propose a system **InternLM-XComposer2.5-OmniLive (IXC2.5-OL)** composed of fused specialized generalist models for streaming perception, reasoning, and memory, respectively. The system is designed to enable AI models to engage continuously with environments while retaining observations over time. By integrating short-term and long-term multimodal memory, our approach attempts to emulate human-like cognition, enabling more dynamic and sustained interactions.

As shown in Figure 1, the IXC2.5-OL system consists of three key modules: (1) **Streaming Perception Module**: This module processes the multimodal information stream on-the-fly. To ensure perception accuracy and efficiency, the video and audio streams are handled separately. A live video perception model processes the video stream, encoding the information and storing key details in memory. Meanwhile, an audio model recognizes the contents of human speech and other sounds, *e.g.*, barking, knocking, or whistling. It triggers the reasoning process when human

queries occur. (2) **Multi-modal Long Memory Module**: This component integrates both long-term and short-term memory, enabling the retrieval of detailed short-term information as well as long-term historical cues. It continuously compresses short-term memories into more information-rich long-term memories to enhance retrieval efficiency and accuracy. (3) **Reasoning Module**: The reasoning module, activated by the perception module, handles queries and performs reasoning tasks. As the component with the most model parameters, it serves as the core of the system’s deep cognitive processes.

The proposed system empowers AI with the ability to perceive, think, and memorize simultaneously. By overcoming the limitations of alternating perception and reasoning, IXC2.5-OL seeks to provide continuous, adaptive service, and long-term AI service. The proposed system will not only enhance the performance of AI assistants but will also contribute to the broader AI applications capable of continuously interacting and adapting to dynamic environments.

The **IXC2.5-OL** demonstrates strong performance across both audio and video benchmarks. Among the open-source models, IXC2.5-OL achieves competitive results on audio recognition (ASR) benchmarks such as WenetSpeech [140] for Chinese and LibriSpeech [87] for English. For video understanding benchmarks, IXC2.5-OL achieves state-of-the-art results among models with less than 10B parameters, obtaining an M-Avg of 66.2% on MLVU [155] and an overall accuracy of 68.7% on MVBench [62]. Additionally, it demonstrates competitive performance on Video-MME [37] (60.6%) and MMBench-Video [34] (1.42). On recent streaming video bench StreamingBench [67], IXC2.5-OL achieves new SOTA results on open-source models (73.79%), highlighting its exceptional capabilities for real-time video interactions.

To foster the development of the multimodal streaming interaction community, alongside the model parameters, the inference and deployment source code, encompassing both the web frontend and backend code, has also been released. All code and models of IXC2.5-OL are publicly available at <https://github.com/InternLM/InternLM-XComposer/tree/main/InternLM-XComposer-2.5-OmniLive>.

2. Related Works

MLLMs for Text-Image Conversation. Large Language Models (LLMs) [5, 7, 9, 24, 46, 51, 81, 86, 90, 108–110, 136] have garnered significant attention for their remarkable capabilities in language comprehension and generation. Building on this success, Large Vision-Language Models (LVLMs) [3, 6, 17–19, 28, 30, 31, 36, 56, 68, 82, 88, 132, 147, 147, 156] have been developed by integrating LLMs with vision encoders [4, 10, 14, 21, 22, 29, 70, 74,

75, 85, 91, 104, 115, 135, 138, 141, 150], extending their ability to comprehend visual content and enabling applications like text-image conversations. Earlier LVLMs were primarily designed for single-image, multi-round conversations, whereas recent advancements [1, 4, 30, 48, 58, 68, 103, 148, 153] have expanded their capabilities to process and understand multi-image inputs.

MLLMs for Video Understanding. In addition to advancements in image understanding, the field of MLLMs has seen growing efforts in video analysis [32, 34, 61, 73, 80, 98, 100, 113, 127]. To address the complexity of video inputs, existing approaches leverage techniques such as sparse sampling or temporal pooling [44, 66, 77, 79, 133], compressed video tokens [16, 49, 60, 63, 94, 119, 144], and memory banks [33, 43, 89, 98, 100, 118, 145]. Additionally, some methods utilize language as a bridge for video understanding [45, 50, 142]. Beyond these video-specific strategies, video analysis can also be framed as interpreting a high-resolution composite image generated from sampled video frames [52, 126, 149]. Recent advancements [12, 117, 120, 145] have increasingly focused on online video understanding, aiming to simulate real-world scenarios where AI processes video streams in real-time to comprehend the environment on-the-fly. However, existing solutions still lack the capability to simultaneously perform perception, memory, and reasoning, limiting their applicability for consistent and long-term human-AI interactions.

MLLMs for Audio Understanding. Audio understanding can be effectively modeled as a sequence-to-sequence (Seq2Seq) task [93], which enables powerful integration with large language models by incorporating audio tokenizers and encoders [25, 105, 137, 143]. In addition to receiving the audio input, recent research investigates streaming duplex speech models [78, 114, 116, 134] that allow speakers to interrupt freely. Beyond audio-text models, emerging research delves into audio-visual models [59, 96] and unified architectures that process audio, visual, and text modalities [38, 64, 139].

MLLMs for Omni-Modal Understanding. Integrating multiple modalities into a single omni-modal foundation model represents a promising research direction. Existing works [13, 38, 42, 64, 102, 121, 124, 139] explore models capable of processing omni-modal inputs, typically combining video and audio, to produce outputs in various formats. These outputs include text [38, 42, 64], audio [13, 102, 124], and omni-modal contents [121, 139]. In the current design of IXC2.5-OL, we handle the audio and video modalities separately to mitigate potential influence during joint training. In future versions, our model will incorporate joint training across all modalities, enabling seamless omni-modality integration.

Table 1. Overview of datasets used in pretraining and supervised fine-tuning (SFT) for the Audio Translation Module. The pre-training stage focuses solely on the automatic speech recognition (ASR) task, utilizing the GigaSpeech and WenetSpeech datasets. The SFT stage includes both ASR and audio classification (CLS) tasks, leveraging diverse datasets. For CommonVoice, we only use the English and Chinese splits. Additionally, 475 self-constructed “Silence” samples are used for CLS tasks.

Stage	Task	Dataset	Data Num
Pretrain	ASR	GigaSpeech [11]	8,282,987
		WenetSpeech [140]	17,821,017
SFT	ASR	LibriSpeech [87]	281,241
		VCTK [111]	44,070
		AISHELL-1 [8]	120,098
		AISHELL-4 [39]	102,254
		MD-RAMC [129]	219,325
		ASCEND [76]	12,314
		KeSpeech [106]	888,428
		DASR [27]	190,732
	CLS	CommonVoice [2]	2,813,852
		FSD50K [35]	40,966
		AudioSet [53]	18,683
		Silence	475

3. Method

As we briefly introduced in Sec.1, the IXC2.5-OL has three disentangled modules: 1) the Streaming Perception Module for on-the-fly visual and audio information processing, 2) the Multi-modal Long Memory Module for memory integration and retrieval, and 3) the Reasoning Module collect information from the perception and memory module, and handles queries and performs reasoning tasks. All the modules work simultaneously and interact asynchronously.

3.1. Streaming Perception Module

Besides nature language, the IXC2.5-OL could handle video and audio natively. To realize this, the Streaming Perception Module contains an Audio Translation Module and a Video Perception Module.

Audio Translation Module contains an audio encoder, an audio projector, and a Small Language Model (SLM). The audio encoder encodes the input audio sample into high-dimension features, and the audio projector further maps the feature to the input space of the SLM. The SLM outputs both the class (e.g. laughing, clapping, or raining) of the audio and the natural language within the audio (i.e. the automatic speech recognition). In practice, we use the Whisper [92] model as the audio encoder and a Qwen2-1.8B [128] as the SLM. The training contains two stages and we list the training data in Table 1.

Video Perception Module provides coarse-grained visual

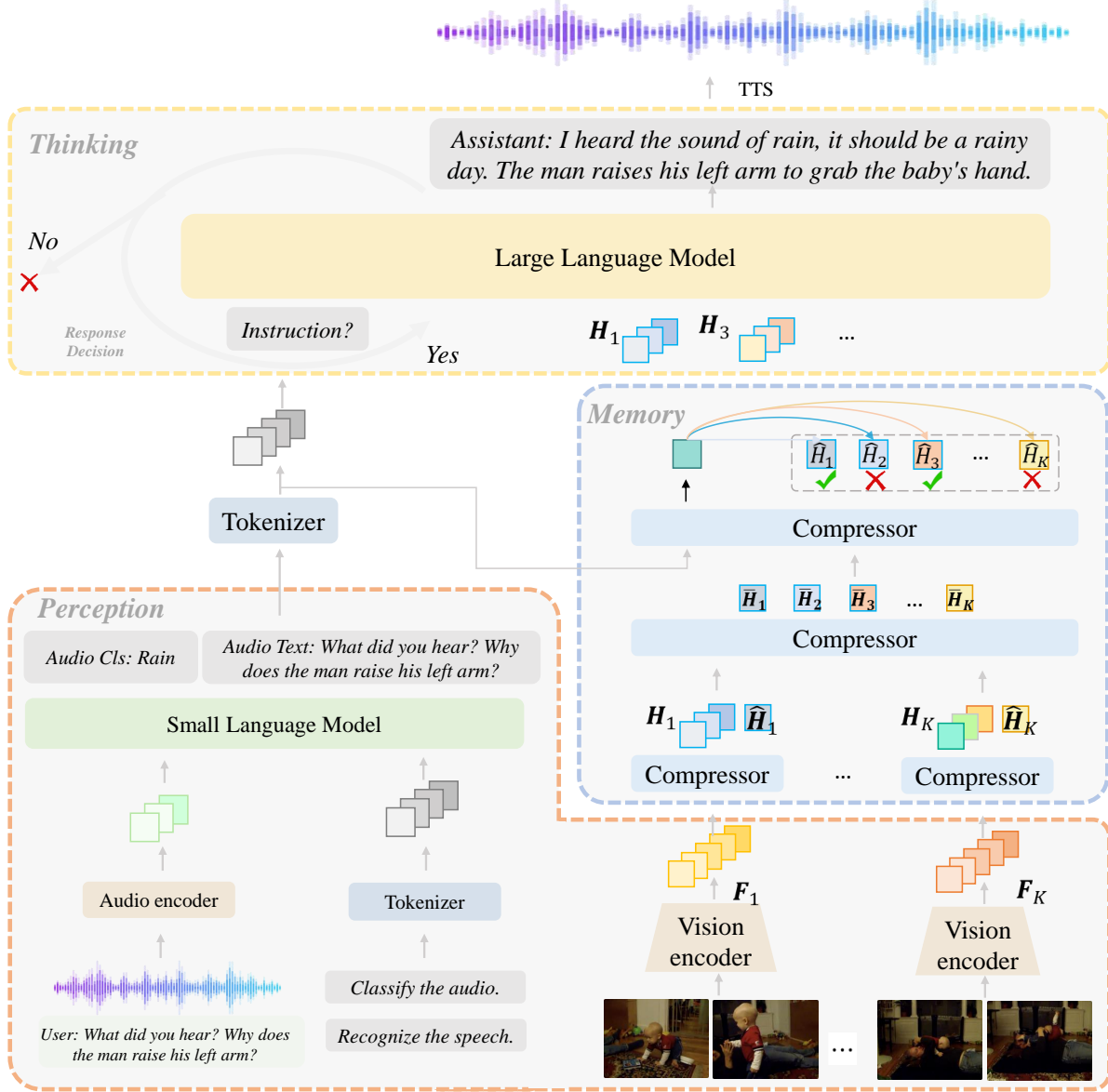


Figure 2. **Pipeline of the InternLM-XComposer2.5-OmniLive (IXC2.5-OL).** The IXC2.5-OL is a real-time interacting system that is constructed by three simultaneous modules: 1) the Streaming Perception Module, 2) the Multi-modal Long Memory Module, and 3) the Reasoning Module.

information to the Multi-modal Long Memory Module. It processes the real-time video input stream and encodes each frame into semantic features. For efficiency, we use the OpenAI CLIP-L/14 [91] In practice.

3.2. Multi-modal Long Memory Module

The Multi-modal Long Memory Module is the core design to handle extremely long video input and helps the Reasoning Module to get rid of millions of tokens from its context window. It shares a similar idea from the VideoStream-

ing [89] that encodes video clips into short-term memories and integrates them into long-term memory. With the given questions, it retrieved the most related video clips for the Reasoning Module. Formally, the Multi-modal Long Memory Module is trained with three tasks:

Video Clip Compression. With features of k_{th} video clip extracted from the Perception Module $F_k \in \mathbb{R}^{TN \times C}$, we initialize its short-term memory $H_k \in \mathbb{R}^{TP \times C}$ by the spatial down-sampling and its global memory $\hat{H}_k \in \mathbb{R}^{1 \times C}$. We realize the compression by the auto-regressive and fea-

ture aggregation nature of LLMs:

$$\mathbf{H}_k, \hat{\mathbf{H}}_k = \text{Compressor}([\mathbf{F}_k \circ \mathbf{H}_k \circ \hat{\mathbf{H}}_k]).$$

Memory Integration. Short-term memory represents the detailed information of each short video clip while the model still lacks a macro view of the video. To this end, with the short-term and global memory of a list of video clips, we integrate them into long-term memory by the Compressor in the following format:

$$\bar{\mathbf{H}}_1, \bar{\mathbf{H}}_2, \dots, \bar{\mathbf{H}}_k = \text{Compressor}([\mathbf{H}_1 \circ \mathbf{H}_2 \dots \circ \mathbf{H}_k \circ \hat{\mathbf{H}}_1 \circ \hat{\mathbf{H}}_2 \dots \circ \hat{\mathbf{H}}_k]).$$

the $\bar{\mathbf{H}} = [\bar{\mathbf{H}}_1, \bar{\mathbf{H}}_2, \dots, \bar{\mathbf{H}}_k] \in \mathbb{R}^{k \times C}$ represents the video in a high-compressed way and we denote it as the long-term memory.

Video Clip Retrieval. When users raise questions, the Multi-modal Long Memory Module retrieves the question-related video clips and provides both the video clips and their short-term memory to the Reasoning Module. In practice, we first encode the question to the feature space of the memory. We concatenate the long-term memory with the tokenized question as the Compressor input, and we view the last token of the output features as the memory-space-aligned question feature. Then we calculate the similarity between the question feature and each video’s global memory, and select the most related clips for the Reasoning Module.

Implementation Detail. We use Qwen2-1.8B [128] as the LLMs and construct several kinds of training data for the three aforementioned tasks. As shown in Table. 2, we train the Video Clip Compression task with short video captioning data from multiple sources, using the same prefix captioning task designed in VideoStreaming [89]. For the Memory Integration task and Video Clip Retrieval task, besides the off-the-shelf video grounding data, we also construct data for two unique tasks: ‘Semantics Implicit Question’ and ‘Reference Implicit Question’.

The ‘Semantics Implicit Question’ means the question does not point to some object directly, but mentions the usage or meaning of the object, and the model should find out the object by understanding the implicit question. For example, when the user asks ‘How about the weather today?’, the model should find out some weather-related object in the past video stream, such as an umbrella, a sun-glass, or something. Another example could be ‘I’m hungry, where can I heat my sandwiches?’, the model should find the microwave oven it has seen before.

The ‘Reference Implicit Question’ means the question uses pronouns rather than nouns. For example, ‘What is this’ means the models should retrieve the current frames, although it does not mention any exact objects.

Model	Dataset
Memory Module	ShareGPT4Video [15], Ego4D[41]
	ActivityNet [32]
	Semantics Implicit QA
	Reference Implicit QA
IXC2.5	ShareGPT4Video [15], ActivityNet [32]
	FunQA [122], TrafficQA [125]
	VideoChat2-IT[61], LLaVA-Video [152]

Table 2. **Video Datasets used in IXC2.5-OL.**

Both kinds of implicit questions are commonly used in real-world communication while current models failed to handle them, so we construct corresponding training data to empower the model with these capabilities.

3.3. Reasoning Module

The Reasoning Module is initialized by an improved version of InternLM-XComposer2.5 (IXC2.5 in the following for simplified statement) and we add a memory projector to align the memory feature with IXC-2.5. For a given questions and both visual and memory information provided by the Memory Module, we formulate the input as:

Question: < |Que| >,
 Here is the question related video clip < |Img| >;
 Here is the question related memory < |Mem| >

In real-world usage, there exists some noisy input that should not be answered (e.g., the user says ‘enn...’ or ‘ok...’), the model should keep salient and wait for the next question. To realize this, we add an additional ‘Instruction Prediction’ process for each question to decide it should be answered or not.

3.4. System Pipeline

As illustrated in Figure 3, the system comprises the Frontend, SRS Server, and Backend Server.

Frontend. The frontend application, developed with JavaScript, enables the camera and microphone to capture video and audio stream inputs, which are then pushed to the SRS server. Concurrently, it establishes a WebSocket connection with the backend to listen for audio outputs and interrupt signals. When audio data is received, the frontend plays it. Upon receiving an interrupt signal, the frontend suspends the audio playback and discards the pending audio.

SRS Server. SRS (Simple Realtime Server) is a straightforward and efficient real-time video server, adept at supporting a multitude of real-time streaming protocols such as RTMP, WebRTC, HLS, HTTP-FLV, SRT, and others. It is renowned for its ability to reliably receive and deliver audio and video streams.

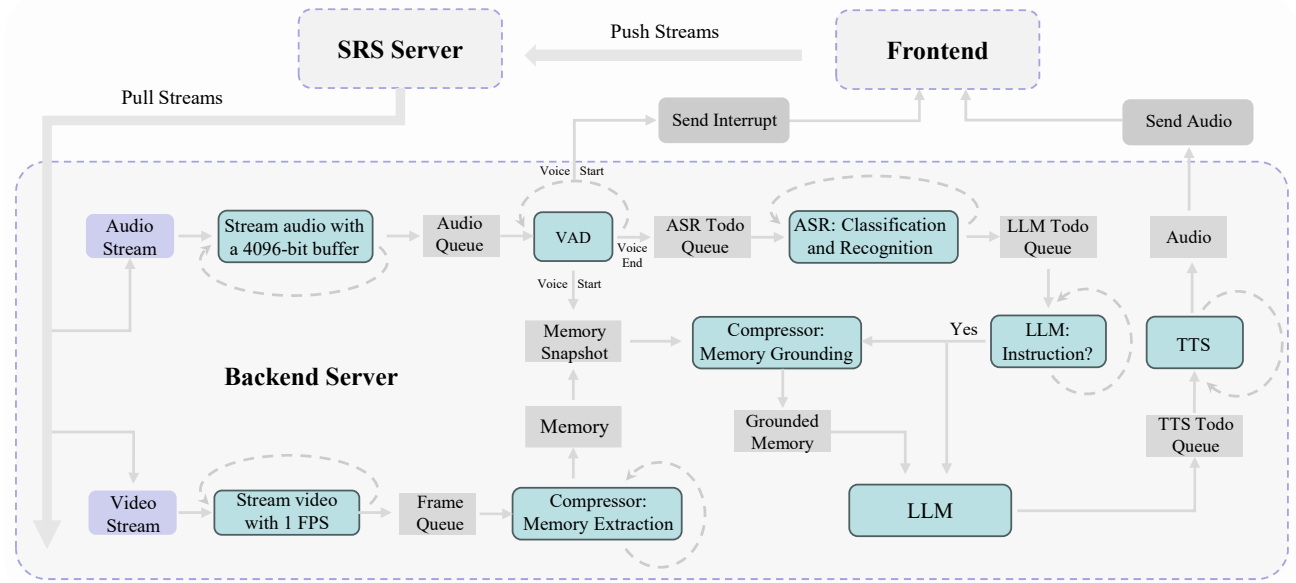


Figure 3. **System pipeline of the IXC2.5-OL.** The system comprises the Frontend, SRS Server, and Backend Server. The Frontend is utilized for capturing video and audio streams and for playing audio from the Backend Server. The SRS Server is employed for managing live streams. The Backend Server is responsible for reading audio and video, extracting memory, and answering questions. The green boxes in the figure represent a thread or a process.

Backend Server. After establishing a WebSocket connection with the frontend, the backend will pull streaming from the SRS Server and initiate separate threads to read audio and video.

The audio reading thread will segment the audio stream into 4096-bit chunks and enqueue them into the *Audio Queue*. The Voice Activity Detection (VAD) [40] thread continuously reads data from *Audio Queue* and detects the start and end of voice activity. Upon detecting the start of voice activity, the backend sends an interrupt signal to the frontend to pause the currently playing audio, and at the same time, dispatches a backup signal to the video process, directing it to save the current memory state. When detecting the end of voice activity, the entire voice segment will be enqueued into *ASR Todo Queue*. The ASR thread continuously reads audio segments from *ASR Todo Queue*, performs background noise classification and voice recognition on them, and then enqueues the results into *LLM Todo Queue* for use by the LLM.

The video reading thread reads video frames at a rate of 1 frame per second and enqueues them into *Frame Queue*. The compressor process reads video frames from the queue, recognizes them, extracts relevant memory, and stores it. Upon receiving a backup signal from the VAD thread, the compressor process will save the current memory state for later retrieval.

The LLM process reads text from the *LLM Todo Queue* and determines whether it is an instruction that requires a re-

sponse from the model. For texts identified as instructions, the compressor process will use the current instruction and the backed-up memory to perform memory grounding, in order to retrieve memories related to the instruction. The LLM process will then generate a response based on the retrieved memories and the instruction, and enqueue the resulting output into *TTS Todo Queue*. An additional TTS thread (e.g., F5-TTS [20], MeloTTS [154]) will convert the text from the *TTS Todo Queue* into audio and send it to the frontend.

4. Experiments

In this section, we validate the benchmark performance of our InternLM-XComposer2.5-OmniLive (IXC2.5-OL), including both audio and video benchmarks.

4.1. Audio Benchmarks

We evaluate our audio models on two prominent automatic speech recognition (ASR) benchmarks: WenetSpeech [140] for Chinese (CN) and LibriSpeech [87] for English (EN). WenetSpeech includes two test sets: Test_Net, which represents high-quality and relatively clean Chinese speech, and Test_Meeting, which captures more challenging conversational scenarios. LibriSpeech consists of four splits: Dev_clean and Test_clean, which contain clean, high-quality English speech, and Dev_other and Test_other, which include noisier, more complex utterances.

As shown in Table 3, our IXC2.5-OL demonstrates supe-

Table 3. **Evaluation results on ASR tasks:** "CN" refers to Chinese speech, while "ENG" refers to English speech. The performance is measured using WER ↓ (Word Error Rate).

Method	LLM	Wenetspeech (CN)		Librispeech (ENG)			
		Test_Net ↓	Test_Meeting ↓	Dev_clean ↓	Dev_other ↓	Test_clean ↓	Test_other ↓
Qwen2-Audio [26]	Qwen2-7B [128]	7.8	8.4	1.3	3.4	1.6	3.6
Mini-Omni [123]	Qwen2-0.5B [128]	-	-	4.5	9.7	4.6	9.2
VITA [38]	Mixtral-8x7B [47]	12.2	16.5	7.6	16.6	8.1	18.4
IXC2.5-OL	Qwen2-1.5B [128]	9.0	9.2	2.5	5.7	2.6	5.8

Table 4. **Evaluation results on MLVU benchmark.** IXC2.5-OL has demonstrated excellent performance, surpassing both open-source models and closed-source APIs, achieving SOTA at the 7B model scale.

Method	Params	Topic Rea.	Anomaly Recog.	Needle QA	Ego Rea.	Plot QA	Action Or.	Action Co.	M-Avg
<i>Closed-source APIs.</i>									
Claude-3-Opus	-	67.2	43.5	21.6	40.2	47.8	18.2	16.7	36.5
Qwen-VL-Max	-	67.4	63.5	40.3	40.9	43.3	25.0	14.8	42.2
GPT-4 Turbo	-	79.5	68.0	45.9	47.4	60.6	26.5	16.1	49.2
GPT-4o	-	87.4	74.5	64.8	57.1	65.1	56.7	46.3	64.6
<i>Open-source models.</i>									
MovieChat [99]	7B	29.5	25.0	24.2	24.7	25.8	28.6	22.8	25.8
LLaMA-VID [65]	7B	50.8	34.5	30.1	32.7	32.5	23.9	27.8	33.2
LLaVA-1.6 [71]	7B	60.6	41.0	43.1	38.4	41.0	25.5	25.7	39.3
ShareGPT4Video [15]	7B	75.8	51.5	47.6	43.2	48.4	34.0	23.3	46.4
VideoLLaMA2 [23]	7B	74.6	64.5	49.9	43.8	45.1	34.0	27.4	48.5
LongVA [149]	7B	83.3	58.5	69.3	50.0	67.2	38.6	27.2	56.3
IXC2.5 [148]	7B	-	-	-	-	-	-	-	58.8
InternVL2 [22]	8B	-	-	-	-	-	-	-	64.0
LLaVA-OneVision [57]	7B	-	-	-	-	-	-	-	64.7
Video-XL [97]	7B	-	-	-	-	-	-	-	64.9
IXC2.5-OL	7B	84.1	68.5	76.6	60.8	75.1	57.1	41.3	66.2

Table 5. **Evaluation results on Video-MME benchmark.** IXC2.5-OL demonstrates performance close to that of the open-source SOTA.

Method	Params	Short	Medium	Long	Overall
<i>Closed-source APIs.</i>					
GPT-4V	-	70.5	55.8	53.5	59.9
Claude 3.5 Sonnet	-	71.0	57.4	51.2	60.0
GPT-4o mini	-	72.5	63.1	58.6	64.8
GPT-4o	-	80.0	70.3	65.3	71.9
Gemini 1.5 Pro	-	81.7	74.3	67.4	75.0
<i>Open-source models.</i>					
ShareGPT4Video [15]	7B	48.3	36.3	35.0	39.9
VideoLLaMA2 [23]	7B	-	-	-	47.9
LongVA [149]	7B	61.1	50.4	46.2	52.6
Video-XL [97]	7B	64.0	53.2	49.2	55.5
VITA [38]	8×7B	65.9	52.9	48.6	55.8
IXC2.5 [148]	7B	-	-	-	55.8
InternVL2 [22]	8B	-	-	-	56.3
LLaVA-OneVision [57]	7B	-	-	-	58.2
mPLUG-Owl3 [131]	7B	70.0	57.7	50.1	59.3
MiniCPM-V 2.6 [130]	8B	-	-	-	60.9
IXC2.5-OL	7B	72.7	58.2	50.8	60.6

rior performance compared to recent streaming audio LLMs such as VITA and Mini-Omni, particularly achieving lower

Word Error Rates (WER) across both CN and EN benchmarks with merely a lightweight 1.5B LLM.

4.2. Video Benchmarks

In Tables 4, 5, 7 and 8, we compare IXC2.5-OL with both closed-source APIs and open-source models on conventional video understanding benchmarks, including MLVU [155], Video-MME [37], MMBench-Video [34] and MVBench [62]. Furthermore, we also assess the performance of different models on the recently proposed StreamingBench [67], which is designed to better evaluate performance for real-time video interactions. The results of this comparison are presented in Table 6. For the video benchmarks, the base model utilizes 64 sampled frames for each video during evaluation.

MLVU MLVU is a comprehensive benchmark designed for evaluating Multimodal Large Language Models in Long Video Understanding tasks. The videos range from 3 minutes to 2 hours and include nine distinct evaluation tasks. Here, we evaluate seven multi-choice tasks, including Topic Reasoning, Anomaly Recognition, Needle QA, Ego Reasoning, Plot QA, Action Order, and Action Count. The de-

Table 6. **Evaluation results on StreamingBench** for Real-Time Visual Understanding. Metrics include Object Perception (OP), Causal Reasoning (CR), Clips Summarization (CS), Attribute Perception (ATP), Event Understanding (EU), Text-Rich Understanding (TR), Prospective Reasoning (PR), Spatial Understanding (SU), Action Perception (ACP), and Counting (CT). IXC2.5-OL excels among all open-source models, and falling just short of the closed-source API, Gemini 1.5 Pro.

Method	Params	Real-Time Visual Understanding										Overall
		OP	CR	CS	ATP	EU	TR	PR	SU	ACP	CT	
Human	-	89.47	92.00	93.60	91.47	95.65	92.52	88.00	88.75	89.74	91.30	91.46
<i>Closed-source APIs.</i>												
Claude 3.5 Sonnet	-	80.49	77.34	82.02	81.73	72.33	75.39	61.11	61.79	69.32	43.09	72.44
GPT-4o	-	77.11	80.47	83.91	76.47	70.19	83.80	66.67	62.19	69.12	49.22	73.28
Gemini 1.5 Pro	-	79.02	80.47	83.54	79.67	80.00	84.74	77.78	64.23	71.95	48.70	75.69
<i>Open-source models.</i>												
VideoLLM-online [12]	8B	39.07	40.06	34.49	31.05	45.96	32.40	31.48	34.16	42.49	27.89	35.99
VideoLLaMA2 [23]	7B	55.86	55.47	57.41	58.17	52.80	43.61	39.21	42.68	45.61	35.23	49.52
VILA-1.5 [68]	8B	53.68	49.22	70.98	56.86	53.42	53.89	54.63	48.78	50.14	17.62	52.32
LongVA [149]	7B	70.03	63.28	61.20	70.92	62.73	59.50	61.11	53.66	54.67	34.72	59.96
InternVL2 [22]	8B	68.12	60.94	69.40	77.12	67.70	62.93	59.26	53.25	54.96	56.48	63.72
Kangaroo [72]	7B	71.12	84.38	70.66	73.20	67.08	61.68	56.48	55.69	62.04	38.86	64.60
MiniCPM-V 2.6 [130]	8B	71.93	71.09	77.92	75.82	64.60	65.73	70.37	56.10	62.32	53.37	67.44
Qwen2-VL [113]	7B	75.20	82.81	73.19	77.45	68.32	71.03	72.22	61.19	69.04	46.11	69.04
LLaVA-OneVision [57]	7B	80.38	74.22	76.03	80.72	72.67	71.65	67.59	65.45	65.72	45.08	71.12
IXC2.5-OL	7B	82.83	73.77	78.66	82.95	72.50	76.01	61.11	60.67	71.59	58.85	73.79

Table 7. **Evaluation results on MMBench-Video**. Tasks include Coarse Perception (CP), Single-Instance Finegrained Perception (FP-S), Cross-Instance Finegrained Perception (FP-C), Hallucination (HL), Logic Reasoning (LR), Attribute Reasoning (AR), Relation Reasoning (RR), Commonsense Reasoning (CSR), and Temporal Reasoning (TP).

Method	Params	Perception					Reasoning						Overall
		CP	FP-S	FP-C	HL	Mean	LR	AR	RR	CSR	TP	Mean	
Closed-source APIs.													
Claude 3.5 Sonnet	-	1.57	1.39	1.07	1.40	1.38	1.13	1.70	1.48	1.54	1.04	1.35	1.38
Gemini 1.0 Pro	-	1.61	1.56	1.30	0.65	1.50	1.15	1.57	1.55	1.36	1.33	1.39	1.48
Gemini 1.5 Pro	-	1.99	2.04	1.70	1.90	1.98	1.98	2.02	1.92	1.78	1.63	1.86	1.94
GPT-4V	-	1.83	1.65	1.40	1.76	1.66	1.45	1.91	1.86	1.83	1.53	1.69	1.68
GPT-4o	-	2.23	2.24	2.01	1.90	2.19	2.11	2.12	2.17	1.94	1.97	2.08	2.15
Open-source models.													
MovieLLM [101]	7B	0.95	0.82	0.70	0.15	0.81	0.52	1.12	1.22	0.54	1.05	0.97	0.87
LLaVA-OneVision [57]	72B	1.22	1.07	0.90	0.21	1.03	0.76	0.96	0.55	0.81	0.48	0.70	0.94
PLLaVA [126]	7B	1.08	1.06	0.86	0.52	1.02	0.64	1.25	1.17	0.98	1.01	1.03	1.03
ShareGPT4Video [15]	7B	1.20	1.05	1.00	0.32	1.04	0.89	1.06	1.19	1.01	0.99	1.03	1.05
VideoStreaming [89]	7B	1.38	1.13	0.8	0.32	1.13	0.77	1.27	1.11	1.01	1.10	1.09	1.12
LLaVA-NeXT-Video [151]	7B	1.35	1.15	0.97	0.58	1.14	0.64	1.38	1.30	1.27	1.03	1.13	1.14
VILA1.5 [68]	13B	1.51	1.45	1.26	0.24	1.39	0.80	1.52	1.30	1.40	1.28	1.28	1.36
InternVL2 [22]	8B	1.41	1.37	1.15	0.19	1.30	0.90	1.34	1.38	1.14	1.00	1.16	1.26
Qwen2-VL [113]	7B	1.63	1.51	1.19	0.55	1.46	1.16	1.56	1.49	1.37	1.21	1.35	1.44
IXC2.5-OL	7B	1.53	1.61	1.20	0.15	1.49	0.93	1.44	1.57	1.30	1.08	1.25	1.42

Table 8. **Evaluation results on MVBench.** Tasks include Action Sequence (AS), Action Prediction (AP), Action Antonym (AA), Fine-grained Action (FA), Unexpected Action (UA), Object Existence (OE), Object Interaction (OI), Object Shuffle (OS), Moving Direction (MD), Action Localization (AL), Scene Transition (ST), Action Count (AC), Moving Count (MC), Moving Attribute (MA), State Change (SC), Fine-grained Pose (FP), Character Order (CO), Egocentric Navigation (EN), Episodic Reasoning (ER), and Counterfactual Inference (CI).

Method	Params	AS	AP	AA	FA	UA	OE	OI	OS	MD	AL	ST	AC	MC	MA	SC	FP	CO	EN	ER	CI	Avg
<i>Closed-source APIs.</i>																						
GPT-4V	-	55.5	63.5	72.0	46.5	73.5	18.5	59.0	29.5	12.0	40.5	83.5	39.0	12.0	22.5	45.0	47.5	52.0	31.0	59.0	11.0	43.5
GPT-4o	-	61.5	56.5	72.0	54.0	82.0	62.5	66.5	44.0	36.5	33.5	93.0	54.5	33.5	54.5	53.5	74.5	71.5	32.5	71.0	42.5	57.5
<i>Open-source models.</i>																						
VideoLLaMA [144]	7B	27.5	25.5	51.0	29.0	39.0	48.0	40.5	38.0	22.5	22.5	43.0	34.0	22.5	32.5	45.5	32.5	40.0	30.0	21.0	37.0	34.1
VideoChat [60]	7B	33.5	26.5	56.0	33.5	40.5	53.0	40.5	30.0	25.5	27.0	48.5	35.0	20.5	42.5	46.0	26.5	41.0	23.5	23.5	36.0	35.5
MiniCPM-V 2.6 [130]	7B	38.0	43.0	63.0	35.5	67.5	55.5	46.0	35.5	25.5	33.0	77.5	48.0	37.0	54.0	42.5	40.0	31.0	38.0	43.0	40.5	44.7
VideoChat2 [62]	7B	66.0	47.5	83.5	49.5	60.0	58.0	71.5	42.5	23.0	23.0	88.5	39.0	42.0	58.5	44.0	49.0	36.5	35.0	40.5	65.5	51.1
Qwen2-VL [113]	7B	51.0	58.0	77.5	47.0	64.0	63.0	65.5	40.0	25.5	35.5	77.0	43.5	47.0	62.0	42.0	61.5	49.5	41.5	47.5	41.5	52.0
PLLaVA [126]	34B	65.0	53.0	83.5	45.0	77.5	70.0	64.5	38.5	37.5	49.0	89.5	41.5	43.5	70.0	53.0	52.5	65.0	39.5	60.5	58.0	57.8
LLaVA-OneVision [57]	72B	63.0	58.0	84.5	46.5	85.5	64.0	73.5	41.5	37.0	69.0	95.0	47.5	47.5	75.5	53.5	52.0	70.5	34.0	64.0	54.5	60.8
InternVL2 [22]	8B	75.0	62.0	83.5	40.5	69.5	96.0	72.0	29.5	58.0	53.0	88.5	39.5	83.0	97.0	51.0	78.5	65.0	33.0	48.0	67.0	64.5
IXC2.5-OL	7B	84.5	81.0	75.0	46.0	81.0	92.0	79.5	36.5	83.0	47.0	90.0	60.5	75.0	93.0	58.0	60.5	74.0	42.0	53.0	62.0	68.7

tailed comparisons are given in Table 4. The IXC2.5-OL exhibits state-of-the-art (SOTA) performance among closed-source APIs, and open-source models with parameters less than 10 billion, surpassing the previous SOTA by 1.3% for Video-XL, 1.6% for GPT-4o.

Video-MME Video-MME is a high-quality video benchmark. The videos are collected from 6 primary visual domains with 30 subfields to ensure broad scenario generalizability, encompassing both short-, medium-, and long-term videos, ranging from 11 seconds to 1 hour. As demonstrated in Table 5, the IXC2.5-OL exhibits competitive performance on this benchmark, comparable to previous SOTA MiniCPM-V 2.6.

StreamingBench StreamingBench is a streaming video benchmark designed for real-time video evaluation. It comprises 18 tasks, showcasing 900 videos and 4,500 human-curated QA pairs. In this context, we focus on assessing visual understanding in real-time. Table 6 illustrates the comparative analysis, demonstrating that IXC2.5-OL excels among all open-source models, achieving a 2.67% improvement over the previous state-of-the-art model, LLaVA-OneVision, and falling just short of the closed-source API, Gemini 1.5 Pro. This performance solidifies IXC2.5-OL’s remarkable prowess in real-time video interaction.

MMBench-Video MMBench-Video is a free-form QA video benchmark consisting of 600 videos and 2000 QA pairs. The duration of each video varies from 30 seconds to 6 minutes. Given the open-ended nature of the answers, the benchmark utilizes GPT-4-based evaluation to enhance

quality in terms of accuracy, consistency, and alignment with human judgment. The results are presented in Table 7. IXC2.5-OL demonstrates state-of-the-art performance on perception tasks and comparable performance on overall evaluations.

MVBench MVBench is a video benchmark that emphasizes temporal understanding. It encompasses 20 challenging video tasks that cannot be effectively addressed using a single frame. As shown in Table 8, IXC2.5-OL, despite having a smaller 7B parameter size, has outperformed both the GPT-4 series and the 72B open-source model LLaVA-OneVision, demonstrating its strong capability in understanding video temporal dynamics.

5. Conclusion

We have presented IXC2.5-OL, a real-time streaming model that advances multi-modal text, audio, and visual capabilities with long-term memory. IXC2.5-OL empowers users to engage in dynamic and interactive experiences. Our model’s real-time processing enables fluid and responsive interactions, allowing users to engage with ever-changing environments of multimodal data seamlessly, providing a more intuitive and efficient user experience. Our future work will focus on reducing system latency to provide a seamless user experience.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35: 23716–23736, 2022. 3
- [2] Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670*, 2019. 3
- [3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. OpenFlamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 2
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A frontier large vision-language model with versatile abilities. *arXiv.org*, 2023. 2, 3
- [5] Baichuan. Baichuan 2: Open large-scale language models. *arXiv.org*, 2023. 2
- [6] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Saġnak Taşlılar. Introducing our multimodal models, 2023. 2
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:1877–1901, 2020. 2
- [8] Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline. In *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*, pages 1–5. IEEE, 2017. 3
- [9] Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024. 2
- [10] Yuhang Cao, Pan Zhang, Xiaoyi Dong, Dahua Lin, and Jiaqi Wang. DualFocus: Integrating macro and micro perspectives in multi-modal large language models. *arXiv preprint arXiv:2402.14767*, 2024. 2
- [11] Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, et al. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv:2106.06909*, 2021. 3
- [12] 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. 2, 3, 8
- [13] Kai Chen, Yunhao Gou, Runhui Huang, Zhili Liu, Daxin Tan, Jing Xu, Chunwei Wang, Yi Zhu, Yihan Zeng, Kuo Yang, et al. EMOVA: Empowering language models to see, hear and speak with vivid emotions. *arXiv preprint arXiv:2409.18042*, 2024. 3
- [14] 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. 2
- [15] 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. 5, 7, 8
- [16] 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, 2024. 3
- [17] Xi Chen, Josip Djolonga, Piotr Padlewski, et al. PaLI-X: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023. 2
- [18] Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, Jialin Wu, Paul Voigtlaender, et al. Pali-3 vision language models: Smaller, faster, stronger, 2023.
- [19] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model, 2023. 2
- [20] Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng, Chunhui Wang, Jian Zhao, Kai Yu, and Xie Chen. F5-tts: A fairy-taler that fakes fluent and faithful speech with flow matching. *arXiv preprint arXiv:2410.06885*, 2024. 6
- [21] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023. 2
- [22] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites, 2024. 2, 7, 8, 9
- [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. 7, 8
- [24] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv.org*, 2022. 2
- [25] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-Audio: Advancing universal audio understanding via uni-

- fied large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023. 3
- [26] Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024. 7
- [27] Samuele Cornell, Taejin Park, Steve Huang, Christoph Boeddeker, Xuankai Chang, Matthew Maciejewski, Matthew Wiesner, Paola Garcia, and Shinji Watanabe. The chime-8 dasr challenge for generalizable and array agnostic distant automatic speech recognition and diarization. *arXiv preprint arXiv:2407.16447*, 2024. 3
- [28] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 2
- [29] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024. 2
- [30] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024. 2, 3
- [31] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, et al. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023. 2
- [32] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. ActivityNet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 3, 5
- [33] Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. Videoagent: A memory-augmented multimodal agent for video understanding, 2024. 2, 3
- [34] Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao, Yining Li, Dahua Lin, and Kai Chen. MMBench-Video: A long-form multi-shot benchmark for holistic video understanding. *arXiv preprint arXiv:2406.14515*, 2024. 2, 3, 7
- [35] E Fonseca, X Favory, J Pons, F Font, and X Serra. Fsd50k: an open dataset of human-labeled sound events, in *arxiv. arXiv preprint arXiv:2010.00475*, 2020. 3
- [36] Chaoyou Fu, Renrui Zhang, Zihan Wang, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, Yunhang Shen, Mengdan Zhang, Peixian Chen, Sirui Zhao, Shao-hui Lin, Deqiang Jiang, Di Yin, Peng Gao, Ke Li, Hongsheng Li, and Xing Sun. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*, 2023. 2
- [37] Chaoyou Fu, Yuhang 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. 2, 7
- [38] 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. 2, 3, 7
- [39] Yihui Fu, Luyao Cheng, Shubo Lv, Yukai Jv, Yuxiang Kong, Zhuo Chen, Yanxin Hu, Lei Xie, Jian Wu, Hui Bu, et al. Aishell-4: An open source dataset for speech enhancement, separation, recognition and speaker diarization in conference scenario. *arXiv preprint arXiv:2104.03603*, 2021. 3
- [40] Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Zhangyu Xiao, and Shiliang Zhang. Funasr: A fundamental end-to-end speech recognition toolkit. In *INTERSPEECH*, 2023. 6
- [41] 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. 5
- [42] 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, 2023. 3
- [43] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. Ma-Imm: Memory-augmented large multimodal model for long-term video understanding. *arXiv preprint arXiv:2404.05726*, 2024. 3
- [44] Suyuan Huang, Haoxin Zhang, Yan Gao, Yao Hu, and Zengchang Qin. From image to video, what do we need in multimodal llms? *arXiv preprint arXiv:2404.11865*, 2024. 3
- [45] Md Mohaiminul Islam, Ngan Ho, Xitong Yang, Tushar Nagarajan, Lorenzo Torresani, and Gedas Bertasius. Video recap: Recursive captioning of hour-long videos. *arXiv preprint arXiv:2402.13250*, 2024. 3
- [46] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, et al. Mistral 7b, 2023. 2
- [47] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 7

- [48] Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. Mantis: Interleaved multi-image instruction tuning, 2024. 3
- [49] Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual representation empowers large language models with image and video understanding. *arXiv preprint arXiv:2311.08046*, 2023. 3
- [50] Kumara Kahatapitiya, Kanchana Ranasinghe, Jongwoo Park, and Michael S Ryoo. Language repository for long video understanding. *arXiv preprint arXiv:2403.14622*, 2024. 3
- [51] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 2
- [52] Wonkyun Kim, Changin Choi, Wonseok Lee, and Wonjong Rhee. An image grid can be worth a video: Zero-shot video question answering using a vlm, 2024. 3
- [53] Qiuqiang Kong, Yong Xu, Wenwu Wang, and Mark D Plumbley. Audio set classification with attention model: A probabilistic perspective. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 316–320. IEEE, 2018. 3
- [54] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4. *International Journal of Computer Vision (IJCV)*, 128(7):1956–1981, 2020. 2
- [55] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 2022. 1
- [56] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv.org*, 2023. 2
- [57] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 7, 8, 9
- [58] Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Guoyin Wang, Bei Chen, and Junnan Li. Aria: An open multimodal native mixture-of-experts model, 2024. 3
- [59] Jungang Li, Sicheng Tao, Yibo Yan, Xiaojie Gu, Haodong Xu, Xu Zheng, Yuanhuiyi Lyu, Linfeng Zhang, and Xuming Hu. SAVEn-Vid: Synergistic audio-visual integration for enhanced understanding in long video context. *arXiv preprint arXiv:2411.16213*, 2024. 3
- [60] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. 3, 9
- [61] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, Limin Wang, and Yu Qiao. Mvbench: A comprehensive multi-modal video understanding benchmark, 2023. 3, 5
- [62] 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 (CVPR)*, 2024. 2, 7, 9
- [63] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. *arXiv preprint arXiv:2311.17043*, 2023. 3
- [64] Yadong Li, Haoze Sun, Mingan Lin, Tianpeng Li, Guosheng Dong, Tao Zhang, Bowen Ding, Wei Song, Zhenglin Cheng, Yuqi Huo, Song Chen, Xu Li, Da Pan, Shusen Zhang, Xin Wu, Zheng Liang, Jun Liu, Tao Zhang, Keer Lu, Yaqi Zhao, Yanjun Shen, Fan Yang, Kaicheng Yu, Tao Lin, Jianhua Xu, Zenan Zhou, and Weipeng Chen. Ocean-omni: To understand the world with omni-modality, 2024. 3
- [65] 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. 7
- [66] 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. 3
- [67] Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen Luo, Peng Li, Yang Liu, and Maosong Sun. Streaming-bench: Assessing the gap for mllms to achieve streaming video understanding. *arXiv preprint arXiv:2411.03628*, 2024. 2, 7
- [68] Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2024. 2, 3, 8
- [69] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 740–755, 2014. 2
- [70] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023. 2
- [71] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 7
- [72] 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. 8
- [73] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Temp-Compass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024. 3
- [74] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s, 2022. 2

- [75] Ziyu Liu, Zeyi Sun, Yuhang Zang, Wei Li, Pan Zhang, Xiaoyi Dong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. RAR: Retrieving and ranking augmented mllms for visual recognition. *arXiv preprint arXiv:2403.13805*, 2024. 3
- [76] Holy Lovenia, Samuel Cahyawijaya, Genta Indra Winata, Peng Xu, Xu Yan, Zihan Liu, Rita Frieske, Tiezheng Yu, Wenliang Dai, Elham J Barezi, et al. Ascend: A spontaneous chinese-english dataset for code-switching in multi-turn conversation. *arXiv preprint arXiv:2112.06223*, 2021. 3
- [77] Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023. 3
- [78] Ziyang Ma, Yakun Song, Chenpeng Du, Jian Cong, Zhuo Chen, Yuping Wang, Yuxuan Wang, and Xie Chen. Language model can listen while speaking. *arXiv preprint arXiv:2408.02622*, 2024. 3
- [79] 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. 3
- [80] Munan Ning, Bin Zhu, Yujia Xie, Bin Lin, Jiayi Cui, Lu Yuan, Dongdong Chen, and Li Yuan. Video-bench: A comprehensive benchmark and toolkit for evaluating video-based large language models. *arXiv preprint arXiv:2311.16103*, 2023. 3
- [81] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>, 2022. 2
- [82] OpenAI. Gpt-4 technical report, 2023. 2
- [83] OpenAI. Gpt-4v(ision) system card, 2023. 2
- [84] OpenAI. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2
- [85] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. 3
- [86] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [87] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE, 2015. 2, 3, 6
- [88] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv.org*, 2023. 2
- [89] Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video understanding with large language models, 2024. 3, 4, 5, 8
- [90] Qwen. Introducing Qwen-7B: Open foundation and human-aligned models (of the state-of-the-arts), 2023. 2
- [91] 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 *Proceedings of the International Conference on Machine learning (ICML)*, 2021. 3, 4
- [92] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. 3
- [93] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *ICML*, 2023. 3
- [94] Michael S. Ryoo, Honglu Zhou, Shrikant Kendre, Can Qin, Le Xue, Manli Shu, Silvio Savarese, Ran Xu, Caiming Xiong, and Juan Carlos Niebles. xgen-mm-vid (blip-3-video): You only need 32 tokens to represent a video even in vlms, 2024. 3
- [95] 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. 2022. 2
- [96] Fangxun Shu, Lei Zhang, Hao Jiang, and Cihang Xie. Audio-visual llm for video understanding. *arXiv preprint arXiv:2312.06720*, 2023. 3
- [97] Yan Shu, Peitian Zhang, Zheng Liu, Minghao Qin, Junjie Zhou, Tiejun Huang, and Bo Zhao. Video-xl: Extra-long vision language model for hour-scale video understanding. *arXiv preprint arXiv:2409.14485*, 2024. 7
- [98] Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. Moviechat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 3
- [99] Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024. 7
- [100] 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. 3
- [101] 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. 8
- [102] 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, 2024. 3
- [103] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Zhengxiong Luo, Yuezhe Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners, 2024. 3
- [104] Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Alpha-CLIP: A clip model focusing on wherever you want. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3
- [105] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. *arXiv preprint arXiv:2310.13289*, 2023. 3
- [106] Zhiyuan Tang, Dong Wang, Yanguang Xu, Jianwei Sun, Xiaoning Lei, Shuaijiang Zhao, Cheng Wen, Xingjun Tan, Chuandong Xie, Shuran Zhou, et al. Kespeech: An open source speech dataset of mandarin and its eight subdialects. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. 3
- [107] Gemini Team. Gemini: A family of highly capable multimodal models, 2023. 2
- [108] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. <https://github.com/InternLM/InternLM>, 2023. 2
- [109] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv.org*, 2023.
- [110] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models, 2023. 2
- [111] Christophe Veaux, Junichi Yamagishi, and Kirsten Macdonald. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit. 2017. 3
- [112] Jiaqi Wang, Pan Zhang, Tao Chu, Yuhang Cao, Yujie Zhou, Tong Wu, Bin Wang, Conghui He, and Dahua Lin. V3det: Vast vocabulary visual detection dataset. In *The IEEE International Conference on Computer Vision (ICCV)*, 2023. 2
- [113] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution, 2024. 3, 8, 9
- [114] Peng Wang, Songshuo Lu, Yaohua Tang, Sijie Yan, Wei Xia, and Yuanjun Xiong. A full-duplex speech dialogue scheme based on large language models. *arXiv preprint arXiv:2405.19487*, 2024. 3
- [115] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023. 3
- [116] Xiong Wang, Yangze Li, Chaoyou Fu, Lei Xie, Ke Li, Xing Sun, and Long Ma. Freeze-Omni: A smart and low latency speech-to-speech dialogue model with frozen llm. *arXiv preprint arXiv:2411.00774*, 2024. 3
- [117] Yueqian Wang, Xiaojun Meng, Yuxuan Wang, Jianxin Liang, Jiansheng Wei, Huishuai Zhang, and Dongyan Zhao. Videollm knows when to speak: Enhancing time-sensitive video comprehension with video-text duet interaction format, 2024. 3
- [118] Yuxuan Wang, Cihang Xie, Yang Liu, and Zilong Zheng. Videollamb: Long-context video understanding with recurrent memory bridges, 2024. 2, 3
- [119] Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. Longvlm: Efficient long video understanding via large language models. *arXiv preprint arXiv:2404.03384*, 2024. 3
- [120] Shiwei Wu, Joya Chen, Kevin Qinghong Lin, Qimeng Wang, Yan Gao, Qianli Xu, Tong Xu, Yao Hu, Enhong Chen, and Mike Zheng Shou. Videollm-mod: Efficient video-language streaming with mixture-of-depths vision computation, 2024. 3
- [121] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm, 2024. 3
- [122] Binzhu Xie, Sicheng Zhang, Zitang Zhou, Bo Li, Yuanhan Zhang, Jack Hessel, Jingkang Yang, and Ziwei Liu. Funqa: Towards surprising video comprehension. In *European Conference on Computer Vision*, pages 39–57. Springer, 2025. 5
- [123] Zhifei Xie and Changqiao Wu. Mini-omni: Language models can hear, talk while thinking in streaming. *arXiv preprint arXiv:2408.16725*, 2024. 2, 7
- [124] Zhifei Xie and Changqiao Wu. Mini-omni2: Towards open-source gpt-4o with vision, speech and duplex capabilities, 2024. 3
- [125] Li Xu, He Huang, and Jun Liu. Sutd-trafficqa: A question answering benchmark and an efficient network for video reasoning over traffic events. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9878–9888, 2021. 5
- [126] 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, 2024. 3, 8, 9
- [127] Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian Tang, Shang Yang, Zhijian Liu, Ethan He, Hongxu Yin, Pavlo Molchanov, Jan Kautz, Linxi Fan, Yuke Zhu, Yao Lu, and Song Han. Longvila: Scaling long-context visual language models for long videos, 2024. 3
- [128] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng

- Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 3, 5, 7
- [129] Zehui Yang, Yifan Chen, Lei Luo, Runyan Yang, Lingxuan Ye, Gaofeng Cheng, Ji Xu, Yaohui Jin, Qingqing Zhang, Pengyuan Zhang, et al. Open source magicdata-ramc: A rich annotated mandarin conversational (ramc) speech dataset. *arXiv preprint arXiv:2203.16844*, 2022. 3
- [130] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024. 7, 8, 9
- [131] 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. *arXiv preprint arXiv:2408.04840*, 2024. 7
- [132] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv.org*, 2023. 2
- [133] Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 3
- [134] Wenyi Yu, Siyin Wang, Xiaoyu Yang, Xianzhao Chen, Xiaohai Tian, Jun Zhang, Guangzhi Sun, Lu Lu, Yuxuan Wang, and Chao Zhang. SALMONN-omni: A codec-free llm for full-duplex speech understanding and generation. *arXiv preprint arXiv:2411.18138*, 2024. 3
- [135] Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and Chen Change Loy. Contextual object detection with multimodal large language models. *arXiv preprint arXiv:2305.18279*, 2023. 3
- [136] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. GLM-130b: An open bilingual pre-trained model. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. 2
- [137] Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. GLM-4-Voice: Towards intelligent and human-like end-to-end spoken chatbot. *arXiv preprint arXiv:2412.02612*, 2024. 3
- [138] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. 3
- [139] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yugang Jiang, and Xipeng Qiu. Anygpt: Unified multimodal llm with discrete sequence modeling, 2024. 3
- [140] Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, et al. Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6182–6186. IEEE, 2022. 2, 3, 6
- [141] Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-CLIP: Unlocking the long-text capability of clip. *arXiv preprint arXiv:2403.15378*, 2024. 3
- [142] Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. A simple llm framework for long-range video question-answering. *arXiv preprint arXiv:2312.17235*, 2023. 3
- [143] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. SpeechGPT: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv preprint arXiv:2305.11000*, 2023. 3
- [144] 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. 3, 9
- [145] 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, 2024. 2, 3
- [146] Kaiyan Zhang, Biqing Qi, and Bowen Zhou. Towards building specialized generalist ai with system 1 and system 2 fusion. *arXiv preprint arXiv:2407.08642*, 2024. 2
- [147] Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*, 2023. 2
- [148] 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. *arXiv preprint arXiv:2407.03320*, 2024. 3, 7
- [149] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkan 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. 2, 3, 7, 8
- [150] Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Chen Change Loy, and Shuicheng Yan. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding, 2024. 3
- [151] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-next: A strong zero-shot video understanding model, 2024. 8
- [152] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data, 2024. 5

- [153] Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. Mmicl: Empowering vision-language model with multi-modal in-context learning. *arXiv.org*, 2023. [3](#)
- [154] Wenliang Zhao, Xumin Yu, and Zengyi Qin. Melotts: High-quality multi-lingual multi-accent text-to-speech, 2023. [6](#)
- [155] 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. [2](#), [7](#)
- [156] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv.org*, 2023. [2](#)