# *Ola*: Pushing the Frontiers of Omni-Modal Language Model with Progressive Modality Alignment

Zuyan Liu[1,2*]    Yuhao Dong[3,2*]    Jiahui Wang[1]
Ziwei Liu[3]    Winston Hu[2]    Jiwen Lu[1†]    Yongming Rao[2,1†]

[1] Tsinghua University    [2] Tencent Hunyuan Research    [3] S-Lab, NTU
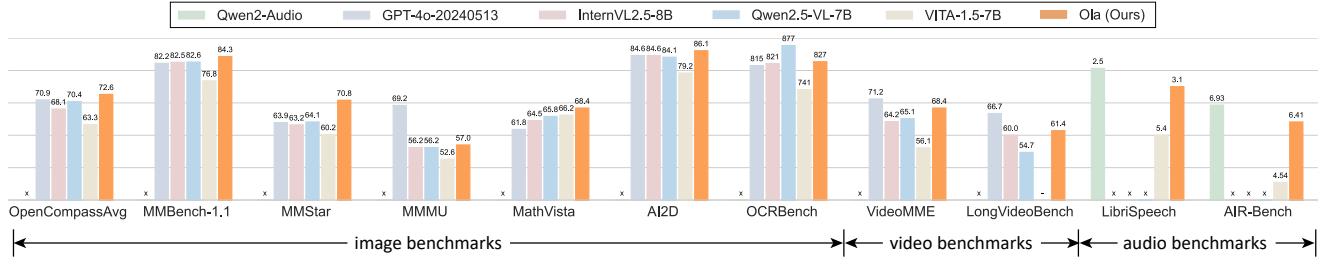
https://ola-omni.github.io/

Figure 1. ***Ola* pushes the frontiers of the omni-modal language model across image, video and audio understanding benchmarks.** We compare *Ola* with existing state-of-the-art open-sourced multimodal models and GPT-4o on their abilities in mainstream image, video, and audio benchmarks. For fair comparisons, we select around 7B versions of existing MLLMs. *Ola* can achieve outperforming performance against omni-modal and specialized MLLMs in all modalities thanks to our progressive alignment strategy. "×" indicates that the model is not capable of the task and "−" indicates the result is lacking. The score for LibriSpeech is inverted as lower is better for the WER metric.

## Abstract

*Recent advances in large language models, particularly following GPT-4o, have sparked increasing interest in developing omni-modal models capable of understanding more modalities. While some open-source alternatives have emerged, there is still a notable lag behind specialized single-modality models in performance. In this paper, we present **Ola**, an **O**mni-modal **la**nguage model that achieves competitive performance across image, video, and audio understanding compared to specialized counterparts. The core design of Ola lies in its progressive modality alignment strategy that extends the supporting modality of the language model progressively. Our training pipeline begins with the most distinct modalities: image and text, then gradually expands the skill sets of the model using speech data that connects language and audio knowledge, and video data that connects all modalities. The progressive learning pipeline also enables us to maintain a relatively small size of the cross-modal alignment data, making developing omni-modal from existing vision-language models easy and less costly. Moreover, to unlock an advanced interactive experience like GPT-4o, we further design a sentence-wise decoding solution for streaming speech generation. Extensive experiments demonstrate that Ola surpasses existing open omni-modal LLMs across all modalities while achieving highly competitive performance compared to state-of-the-art specialized models of similar sizes. We aim to make Ola a fully open omni-modal understanding solution to advance future research in this emerging field. Model weights, code, and data are open-sourced at https://github.com/Ola-Omni/Ola.*

## 1. Introduction

Multi-Modal Large Language Models are drawing increasing attention owing to their strong instruction-following capabilities and abundant knowledge of handling complex inputs including texts, images, videos, and audio. Based on the strong performance of open-sourced large language models [56, 76],

---

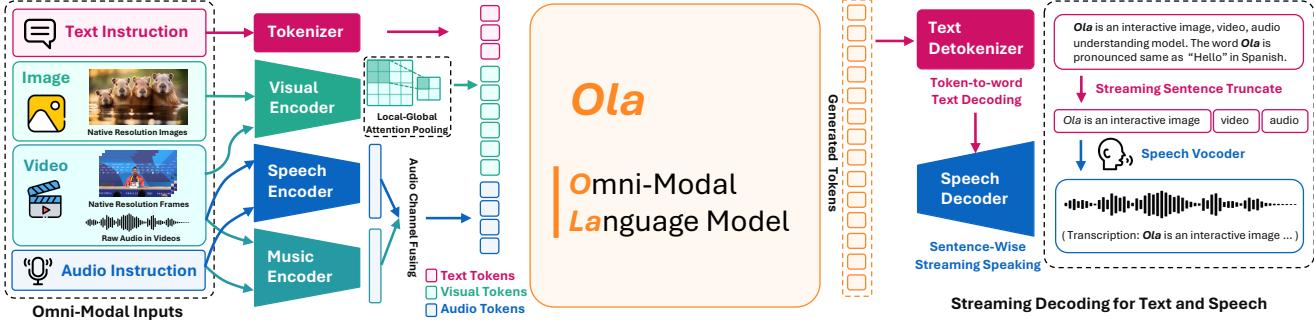*  Equal Contribution.   † Corresponding authors.

Figure 2. **Ola Architecture.** *Ola* supports omni-modal inputs including text, image, video, and audio, capable of processing the inputs simultaneously with competitive performance on understanding tasks for all these modalities. Meanwhile, *Ola* supports user-friendly real-time streaming decoding for texts and speeches thanks to the text detokenizer and the speech decoder.
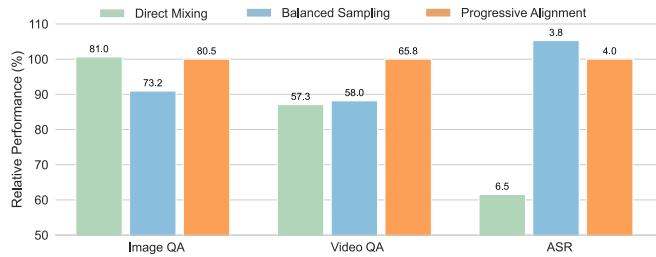


Figure 3. **Progressive modality alignment helps to learn better omni-modal models.** We compare our progressive alignment strategy with two baseline training pipelines on Image QA(MMBench [40]), Video QA(VideoMME [21]), and ASR(LibriSpeech [54]): 1) direct mixing where all instruction tuning data is merged and trained in a single stage, and 2) balanced sampling where we upsample certain sources to make the training data more balanced among modalities. The experiment is conducted on a subsampled training set for efficiency and we train models for the same number of steps for fair comparisons. The score is normalized based on the score of progressive alignment to calculate the relative score and the ASR score is inverted as lower is better for the WER metric.

extensive research has been done on connecting specific modalities with language responses [10, 32, 36, 55, 57, 67]. Recently, the success of GPT-4o [53] and Gemini [24] aiming at supporting more modalities in Large Language Models inspires researchers to take one important steps towards omni models that understand all the inputs in one model.

The core challenges in training omni-modal Large Language Models lie in the modeling of modalities in various distributions, and the design of an effective learning pipeline to achieve competitive, balanced performance on all the supported tasks. Several attempts have been made to overcome the difficulty of omni-modal models [17, 22, 72], where we illustrate the mainstream works and the state-of-the-art MLLMs [11, 32, 67, 84] in specific domains in Fig. 1. Impressive performance and modality breadth are contradicted

in previous works, while existing open-sourced omni-modal solutions still have a large performance gap between state-of-the-art specialized LLMs, making a strong barrier between the concept of omni-modal and real-world applications. Moreover, the lack of capability in specific domains or tasks, the mass data demand, the delay for user interaction, and the inadequate alignment across modalities show the suboptimal for existing omni-modal models.

In this paper, we propose the *Ola* model, exploring the solution for training an omni-modal Large Language Model with comparable performance with state-of-the-art specific LLMs, real-time interaction, and high efficiency on alignment data. The core design of the *Ola* model is the progressive modality alignment strategy. To build connections between language and vision, we start from two fundamental and separated modalities, image and text, to build the basic knowledge for omni-modal models. Subsequently, we gradually expand the training sets to equip the model with an extended ability including video frames that strengthen the visual understanding capability, speech data that connects the language and audio knowledge, and the video with audio that mixes up the information from language, video, and audio comprehensively. The progressively learning strategy makes omni-modal learning easier by disassembling the complex training procedure into small steps, therefore maintaining a small size of cross-modal alignment data and making it easier to start from existing achievements in vision-language models. As shown in Fig. 3, the performance of Ola largely benefits from our progressive training pipeline, leading to more balanced and competitive results on all modalities.

To cooperate with the training strategy, important improvements have been made to the architecture and data domains. 1) The *Ola* architecture supports omni-modal inputs and streaming text and speech generation with extendable and concise architecture design. We design the joint alignment module for vision and audio, fusing the visual inputs with a local-global attention pooling layer, and make free combinations for visual, audio, and text tokens.

Moreover, we integrate the sentence-wise streaming decoding module for high-quality voice synthesis. 2) Beyond the collected fine-tuning data in vision and audio aspects, we deeply excavate the relationships between video and the corresponding audio to construct the bridge between visual and audio modality. Specifically, we collect raw videos from academic and open-ended web sources, design separated cleaning pipelines, and then utilize vision-language models to generate question-answering pairs based on the subtitles and video content.

We evaluate *Ola* under the complete omni-modal benchmarks including image, video, and audio aspects. With only 7B parameters, *Ola* achieves competitive performance across mainstream multi-modal benchmarks. On Image Benchmarks, *Ola* excels at general understanding and specific-task understanding, with an overall mean accuracy of 72.6% on the challenging OpenCompass benchmark [16], 84.3% average scores on MMBench-1.1 [40], 57.0% average scores on MMMU [77], etc. On the challenging VideoMME [21] multiple-choice benchmark ranging from videos within 30 seconds to 1 hour, *Ola* achieves the impressive accuracy of 68.4% with video and audio inputs. *Ola* also excels at audio understanding tasks such as audio-speech recognition and chat evaluation, achieving 3.1 mean WER on LibriSpeech [54] and 6.41 GPT-eval score on AIR-Bench [74]. Results on the benchmarks show a giant promotion compared with existing omni-modal LLMs and even outperforming the performance of state-of-the-art specialized LLMs.

## 2. Related Works

**Large Vision-Language Models.** Inspired by the success of AI assistants and large language models [50–52], research has increasingly focused on vision-language multi-modal large language models. Significant advancements have been made in architecture design [3, 13, 33, 39, 44], training strategies [10, 43], model scaling [32], and data curation [31, 43, 45, 73]. Furthermore, models are evolving beyond static images to support video [7, 36, 42, 47], 3D [27, 37], and mixed visual inputs [57, 59]. However, extending these models to effectively integrate audio modalities while maintaining balanced and robust performance remains an area that has not been fully explored.

**Large Audio-Text Models.** Large Language Models, mainly focused on text inputs and outputs, have a foundational link to speech, with pioneering efforts integrating speech inputs through adapter-based modifications [4, 19, 70]. The challenge of LLM-based speech generation has been addressed with the development of speech decoders [61, 80], marking a significant step towards omni-modal models. Beyond speech, research is expanding into audio-based LLMs that encompass the understanding of music, events, and more. Notable examples include AudioGPT [28] and SALMONN [63],

which explore these audio dimensions, while models like Qwen2-Audio [11] demonstrate enhanced understanding capabilities.

**Towards Large Omni-Modal Models.** Recent advancements in large language models [24, 53] have spurred interest in developing omni-modal models that can handle multiple modalities simultaneously. Notable examples include SpeechGPT [80] and LLaMA-Omni [17], which integrate audio-text understanding with speech generation. The Mini-Omni [72] series addresses challenges in speech streaming generation through parallel decoding techniques. VITA [22] extends this capability by unifying audio, image, video, and text understanding. While these models excel in understanding tasks, efforts are also being made to tackle both understanding and generation tasks [68, 71]. However, existing omni-modal models often fall short in managing the full spectrum of input modalities and output formats, or they suffer from significantly poorer performance. Ola aims to address these limitations by enhancing the capability and efficiency of omni-modal models with better architecture, training strategy, and targeting data preparation.

## 3. *Ola*: Omni-Modal Understanding

We put main efforts into three aspects to obtain an omni-modal understanding for *Ola*, capable of reading, hearing, seeing, typing, and speaking arbitrarily. 1) The *Ola* architecture introduced in Sec. 3.1 supports omni-modal inputs and streaming outputs for both text and speech. 2) We design the progressive training strategy in Sec. 3.2 to bridge the modality gaps between language and vision from primary to periphery. 3) Effective omni-modal training data in Sec. 3.3 provide strong performance across all the benchmarks, especially with our proposed cross-modal video data that stresses to learn audio-relevant information from videos.

### 3.1. *Ola* Architecture

A general view of *Ola* architecture is illustrated in Fig. 2. The encoding part accepts omni-modal inputs in text, image, video, and audio formats with modal-specific encoders or embeddings. Subsequently, the joint alignment operations process all the inputs in a unified manner, fusing and concatenating all the sequences into flattened tokens for the core *Ola* Large Language Model. The LLM generates text tokens in serial, and we adopt the speech decoder for streaming speech decoding.

**Omni-Modal Inputs Encoding.** We encode visual, audio, and text inputs separately based on the successful practice of the previous text-to-one-modal Large Language Models. For visual inputs that include images $I$, videos $V_{1,2,\cdots,n}$ with $n$ frames, we follow vision-language models [33, 39] to use multi-modal visual encoder $\mathcal{E}_{I,V}(I, V_{f_1,f_2,\dots})$ for encoding. Note that we preserve the original aspect ratio of
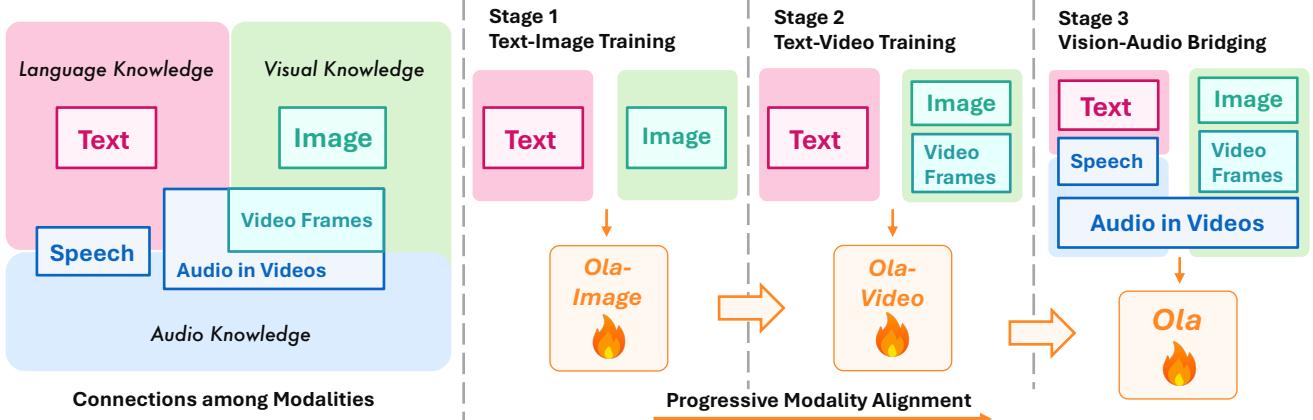
Figure 4. **Illustrations of the *Ola* Progressive Modality Alignment.** We visualize the relationships among modalities in the left part. Speech acts as the connection between language and audio knowledge, while video constructs the bridge with highly relevant visual and audio information. Therefore, we design the progressive alignment training strategy from primary to periphery. Furthermore, we design the cross-modality video-audio data to better capture the relationships among modalities.

each image or frame for the arbitrary resolution vision encoder OryxViT [42] initialized from SigLIP-400M [79], as OryxViT performs a more natural solution for visual inputs. We obtain the image features $f_I$, and the video features for each frame $[f_{V_1}, f_{V_2}, \cdots, f_{V_n}]$ based on the image patches.

For audio encoding, we propose the dual encoder approach for the audio input. Specifically, we use Whisper-v3 [58] as the speech encoder and BEATs [8] as the music encoder for better alignment between audios and texts, providing richer audio information. The music encoder takes the original wav audio $A$ as inputs, while the speech encoder takes the wav transformed into Mel spectrogram representation $A_{(mel)}$ as inputs. Note that the Whisper encoder only supports a certain length of audio inputs, therefore we fix the sample rate as 16000Hz and cut the overlong audio into segments of 30 seconds $A_1, A_2, \cdots, A_n$ and conduct the encoder operation in batches $[f_{A_1}, f_{A_2}, \cdots, f_{A_n}] = \mathcal{E}_A([A_1, A_2, \cdots, A_n])$. The embedding features of the speech and music encoders are concatenated across channel dimensions for the comprehensive audio features $f_A$.

For the text inputs, we use the carefully designed tokenizer and the well-trained embedding layers from the pretrained Large Language Model for the text tokens $t_T$ directly.

**Joint Alignment for Vision and Audio.** The alignment modules act as the converter from specific-modal spaces to the text embedding space, which is an essential part of omni-modal Large Language Models. To reduce the token length of visual features for higher efficiency, we obtain one step forward based on the motivation of structural downsampling in previous works [35], and propose the *Local-Global Attention Pooling* layer for better downsampled features with less information loss. Specifically, for image or video frame feature in spatial shape $H \times W$ and channel $C$, we adopt

the bilinear interpolation for 2x downsampling to obtain $f^{\text{global}}$, which contains the global information of the downsampled region. We combine the original and global features for local-global embeddings and use Softmax to predict the importance $\pi$ of each downsampled spatial region:

$$f = \text{Concat}[f, f^{\text{global}}], \quad \pi = \text{Softmax}(\text{MLP}(f)) \quad (1)$$

The downsampled feature $f^{\text{global}}$ determines the weight of each previous region based on the score $\pi$ with the Hadamard product.

We apply the simple yet effective two-layer non-linear MLP connectors $\text{MLP}_A, \text{MLP}_V$ following the previous works [38, 39] to project the specific modal features $[f_I, f_V, f_A]$ into unified tokens $[t_I, t_V, t_A]$. We define visual and audio start, separate, newline, and end tokens to indicate special positions for inputs. The omni-modal tokens $[t_I, t_V, t_A]$ are concatenated with text tokens $t_T$ in free combination for LLM decoding.

**Streaming Speech Generation.** We adopt CosyVoice [15] as a high-quality speech decoder for speech generation. To support user-friendly streaming decoding, we detect the generated text tokens in real time and truncate the sentence once punctuation is encountered. Afterward, the previous sentence is fed into the speech decoder for audio synthesis. Therefore, *Ola* does not need to wait for the whole sentence to finish while supporting streaming decoding. Though several attempts [17, 72] have been made to train the speech generation module end-to-end, the external text-to-speech decoder is a more efficient, high-quality, and training-free solution for omni-modal models.

4

## 3.2. Progressive Omni-Modal Alignment

**Rethinking Modality Gaps among Language, Vision and Audio.** From our exploration, we recognize two critical issues in omni-modal training. 1) *Modal balancing.* As illustrated in Fig. 3, directly combining data from all modalities negatively affects benchmark performance. Therefore, we propose a rational training procedure that progressively equips the sense organ to the *Ola* model. We assert that texts and images are the core modalities in omni-modal learning, while speeches and videos are variants of texts and images, respectively. Learning to recognize texts and images ensures the model's basic cross-modal ability, so we prioritize these harder cases. Subsequently, we gradually incorporate video, audio, and speech into the training for the omni-modal LLM. 2) *Connections between audio and vision.* Another problem lies in building connections between audio and vision, which has been overlooked by previous works. However, jointly learning audio and vision data can yield surprising results in omni-modal learning by providing a more comprehensive view across different modalities. For the *Ola* model, we consider video as the bridge between audio and vision, as videos contain natural, abundant, and highly relevant information between frames and the accompanying audio. We test our hypothesis by optimizing the training pipeline and preparing targeted training data, as introduced below.

**Stage 1: Text-Image Training.** The *Ola* training starts from a pre-trained Large Language Model, where we use Qwen2.5-7B [64] in our implementation for better trade-offs for model sizes and performance. The *Ola* text-image training includes MLP alignment, large-scale pre-training, and supervised fine-tuning following common practice in large-scale multi-modal learning [32, 67]. We initialize the visual MLP adapter and freeze other parameters in MLP alignment with the image captioning task. Subsequently, we unfreeze all the parameters including the vision encoder in the pre-training and supervised fine-tuning phase. The downsampling module is well-trained in the text-image training stage to hold the 2x compression for visual data including images and videos.

**Stage 2: Continuous Training for Images and Videos.** Based on a strong text-image multi-modal LLM, we continuously extend the capability for *Ola* with video data. We keep most of the experimental setting for supervised fine-tuning while freezing the vision encoder in this stage as the encoder is already fully trained beforehand. We mix the previous image data and the video data to preserve the original text-image performance.

**Stage 3: Bridging Vision and Audio with Videos.** The audio-relevant training is included in stage 3. We follow the training strategy for the visual MLP adapter while initializing the audio MLP adapter with a basic audio-speech recognition (ASR) task. Then we mix up the text & speech under-

standing, text & music understanding, audio & video joint comprehension, and the foremost text-image multi-modal tasks together for the formal training. *Ola* concentrates on learning audio recognition and identifying the relationships between vision and audio in this stage, resulting in a comprehensive image, video, and audio understanding model after training.

## 3.3. Data

The training data of *Ola* includes the general supervised fine-tuning data collected from open-source academic datasets in image, video, and audio categories. Additionally, we design a pipeline to generate cross-modal video-audio data for omni-modal alignment.

**Image Data.** We follow the simple setting in [38] for image MLP alignment. The MLP alignment data includes 800k image captioning pairs from the LAION dataset [62]. For the large-scale pre-training phase, we collect around 20M text-image data pairs from open-sourced and in-house data to build the basic capability of the model. For text-image supervised fine-tuning data, we collect abundant data from various tasks including captions, conversations, OCR, etc. The source of the training data involves the mixture of LLaVA-OneVision [32], Cauldron [31], Cambrian-1 [67], MAmmoTH-VL [26], PixMo [12], etc., resulting in 7.3M image training data in total.

**Video Data.** For text-video training data, we collect useful video datasets from LLaVA-Video-178K [84], VideoChatGPT-Plus [47], LLaVA-Hound [82], and Cinepile [60], with 1.9M video conversation pieces in total. We randomly sample 2/3 video-language data pairs from LLaVA-Video-178K, resulting in 1.2M high-quality training data, and we use the full set of other data sources. In stage 2 for multi-image training, we randomly sample 0.8M image data from stage 1 and mix it with the video datasets for continuous training to maintain the basic performance.

**Audio Data.** We prepare audio training data for comprehensive speech and music understanding. For text-speech understanding, we design multiple tasks including ASR from LibriSpeech [54] and GigaSpeech [5] datasets, audio captioning from AudioCaps [30] and Clotho [14] datasets, speech question answering from LibriSpeech [54] datasets, audio question answering from WavCaps [49] and Audio-Caps [30] datasets. For text-music understanding, we collect tasks including music captioning from MusicCaps [1], music question answering from MillionSong [48] and MusicNet [66]. The overall audio training data involves 1.1M samples. The relevant text question-answering representations are collected from SALMONN [63].

**Generating Cross-Modal Video Data.** Most existing video training data are annotated or synthesized solely from frame inputs, often overlooking the valuable information in accom-

Table 1. **Main Results across Image, Video, and Audio Understanding Benchmarks.** We select representative benchmarks among image, video, and audio benchmarks, and select mainstream state-of-the-art open-source large language models in each modality. We also include open-source omni-modal LLMs for comparison. In the table, "−" indicates the model is capable of solving the tasks theoretically, while the result is lacking. "✗" indicates that the model is not capable of the task. ↓ indicates that lower score is better. * LLaMA-Omni is not optimized for ASR and thus cannot produce reasonable results on this task.

| Model | Size | Image Benchmarks | | | | | | | Video Benchmarks | | | Audio Benchmarks | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | MMBench-1.1 | MMStar | MMMU | MathVista | HalluBench | AI2D | OCRBench | VideoMME | LongVideoBench | MVBench | LibriSpeech↓ | AIR-Bench |
| *Image LLMs* | | | | | | | | | | | | | |
| Cambrian-1 [67] | 8B | 68.2 | 50.7 | 41.8 | 48.1 | 30.6 | 74.6 | 614 | ✗ | ✗ | ✗ | ✗ | ✗ |
| Pixtral [2] | 12B | 72.7 | 54.5 | 51.1 | 56.3 | 47.0 | 79.0 | 685 | ✗ | ✗ | ✗ | ✗ | ✗ |
| *Video LLMs* | | | | | | | | | | | | | |
| VideoCCAM [20] | 9B | – | – | – | – | – | – | – | 53.9 | – | 64.6 | ✗ | ✗ |
| LLaVA-Video [84] | 7B | – | – | – | – | – | – | – | 63.3 | 58.2 | 58.6 | ✗ | ✗ |
| *Vision Comprehensive LLMs* | | | | | | | | | | | | | |
| LLaVA-OneVision [32] | 7B | 80.9 | 61.9 | 47.9 | 62.6 | 31.6 | 82.4 | 622 | 58.2 | 61.3 | 59.4 | ✗ | ✗ |
| MiniCPM-V 2.6 [75] | 8B | 78.0 | 57.5 | 49.8 | 60.8 | 48.1 | 82.1 | 852 | 60.9 | – | – | ✗ | ✗ |
| InternVL2.5 [9] | 8B | 82.5 | 63.2 | 56.2 | 64.5 | 49.0 | 84.6 | 821 | 64.2 | 60.0 | **72.0** | ✗ | ✗ |
| Qwen2.5-VL [65] | 7B | 82.6 | 64.1 | 56.2 | 65.8 | **56.3** | 84.1 | **877** | 65.1 | 54.7 | 69.6 | ✗ | ✗ |
| *Audio LLMs* | | | | | | | | | | | | | |
| SALMONN [63] | 13B | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | 3.5 | 6.12 |
| Qwen2-Audio [11] | 7B | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | **2.5** | **6.93** |
| *Omni-Modal LLMs* | | | | | | | | | | | | | |
| LLaMA-Omni [17] | 8B | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | 120.4* | 4.70 |
| Mini-Omni2 [72] | 0.5B | 32.1 | – | 24.9 | – | – | – | 6 | – | – | – | 7.2 | 3.20 |
| VITA-1.5 [23] | 7B | 76.8 | 60.2 | 52.6 | 66.2 | 44.6 | 79.2 | 741 | 56.1 | – | 55.4 | 5.4 | 4.54 |
| IXC2.5-OmniLive [81] | 8B | 79.4 | 59.9 | 42.9 | 64.0 | 43.1 | 81.6 | 686 | 60.6 | – | 68.7 | 4.4 | 1.67 |
| *Ola* | 7B | **84.3** | **70.8** | **57.0** | **68.4** | 53.5 | **86.1** | 827 | **68.4** | **61.4** | 66.3 | 3.1 | 6.41 |

panying audio. To address this, we designed a pipeline to generate cross-modal video data, aiming to uncover the intrinsic relationships between video and audio. This guides an omni-modal large language model in learning cross-modality information. Specifically, we developed two tasks for cross-modal learning: video-audio question answering and video speech recognition. We collected videos from the academic dataset LLaVA-Video-178k [84] and the open-ended video datasets from FineVideo [18]. Due to the lack of subtitles in the academic datasets, we used Whisper-v3 [58] to generate subtitles from the video audio and conducted a language-based cleaning procedure. We then employed a large language model to assess whether the subtitles were complete and informative. We gathered 41k pure videos from LLaVA-Video-178k, along with the original 42k videos from FineVideo. Subsequently, we used Qwen2-VL-72B [57] to generate questions and answers based on the video and corresponding subtitles. The model was instructed to focus on the subtitle inputs while using the videos as supplementary information. We created three question-answer pairs for each video, resulting in 243k cross-modal video-audio data points. Additionally, we included the original video subtitling tasks with 83k training data to help the model maintain its ASR ability in noisy environments. During training, the models processed multiple image frames, audio, and text inputs

simultaneously, significantly enhancing their cross-modal capabilities.

The prepared 324k cross-modal video data is mixed with audio data implemented above for stage 3 training. We mix all the 1.1M pure text-audio training data and 324k cross-modal video-audio data for the comprehensive training stage. Additionally, we sample 400k image data from stage 1 to maintain the basic ability and create 200k image data with voice instructions to equip the model with interaction capability.

## 4. Experiments

We conduct all-sided benchmarking in Sec. 4.2 to evaluate the all-powerful *Ola* model, including the representative benchmarks in image, video, and audio understanding. Subsequently, we conduct detailed results on critical benchmarks in Sec. 4.3 to demonstrate the effectiveness of our design on motivation, training, and data preparations.

### 4.1. Implementation Details

The *Ola* model builds upon the Qwen-2.5-7B [64] framework, incorporating OryxViT [42] as the vision encoder initialized from SigLIP-400M [79], Whisper-V3-Large [58] as the speech encoder, and BEATs-AS2M(cpt2) [8] as the

music encoder. Initially, we employ a relatively high learning rate of 1e-3 for MLP adapter pre-training. During supervised fine-tuning, the learning rate is gradually reduced from 2e-5 for text-image and multi-image training to 1e-5 for video-audio training. We utilize a batch size of 256 for fine-tuning, leveraging 64 NVIDIA A800 GPUs to conduct our training. We adopt $10\times$ downsampled rate for audio features to reduce the token length, resulting in 300 tokens per minute. During training and inference, the maximum token length is set to 16384 and the maximum number of audio trunks is set to 25.

## 4.2. Results on Omni Understanding

**Benchmarks.** We conduct extensive comparisons across image, video, and audio understanding benchmarks to demonstrate the omni-modal capabilities of the *Ola* model. For image benchmarks, we utilize comprehensive understanding datasets including MMBench-1.1 [40], MMMU [77], MM-Star [6], MathVista [46], Hallusion Bench [25], AI2D [29], and OCRBench [41]. In the video domain, we evaluate using VideoMME [21], which involves multiple-choice questions on videos of varying lengths, and LongVideoBench [69] for assessing performance on extremely long video content, MVBench [34] for the general recognition ability. For audio benchmarks, we focus on two primary tasks relevant to audio LLMs. Librispeech [54] serves as a traditional audio-speech recognition (ASR) dataset, testing the model's ability to accurately transcribe spoken language. AIR-Bench [74] provides a comprehensive evaluation of audio question-answering capabilities, incorporating speech, sound, and music inputs. The responses are evaluated with a GPT-based [52] scorer against ground truth answers. We report the mainstream evaluation metric on image and video benchmarks and report the mean metric in Tab. 2 for ASR and audio understanding tasks for simplicity.

**Baselines.** We selected a range of state-of-the-art multi-modal large language models across different modalities for comparison and reference. We categorized vision-language models into three groups: image-centric LLMs, video-centric LLMs, and comprehensive LLMs capable of handling both images and videos. For image understanding, we utilized Cambrian-1 [67] and Pixtral-12B [2]. For video understanding, VideoCCAM [20] and LLaVA-Video [84] are employed. Comprehensive models included LLaVA-OneVision [32], MiniCPM-V 2.6 [75], InternVL2.5 [9], and Qwen2.5-VL [65] which excel across various visual benchmarks. In the audio domain, we compared our work with state-of-the-art models such as SALMONN [63] and Qwen-2 Audio [11]. As an Omni-modal LLM, our model, *Ola*, was compared with state-of-the-art open-source omni-modal LLMs like Mini-Omni2 [72], VITA-1.5 [22], InternLM-XComposer2.5-OmniLive [81], which support image, audio, and text inputs. Additionally, LLaMA-Omni [72], an audio-text omni model, was noted for its strong speech generation

Table 2. **Analysis Results on Audio Benchmarks.** We report the WER rate on test-clean, test-other, dev-clean, dev-other subsets of LibriSpeech, and the scores on AIR-Bench. In the table, "−" indicates the model is capable of solving the tasks, while the result is lacking. "✗" indicates that the model is not capable of the task.

| Model | ASR on LibriSpeech↓ | | | | Chat on AIR-Bench | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | test-c | test-o | dev-c | dev-o | speech | sound | music | mix | avg |
| *Audio Models* | | | | | | | | | |
| SpeechGPT [80] | ✗ | ✗ | ✗ | ✗ | 1.57 | 0.95 | 0.95 | 4.13 | 1.90 |
| Whisper-small [58] | 4.4 | 10.1 | 4.6 | 10.3 | ✗ | ✗ | ✗ | ✗ | ✗ |
| SALMONN [63] | 2.1 | 4.9 | − | − | 6.16 | 6.28 | 5.95 | 6.08 | 6.12 |
| Qwen2-Audio [11] | 1.6 | 3.6 | 1.3 | 3.4 | 7.18 | 6.99 | 6.79 | 6.77 | 6.93 |
| *Omni-Modal LLMs* | | | | | | | | | |
| LLaMA-Omni [17] | ✗ | ✗ | ✗ | ✗ | 5.22 | 5.25 | 4.33 | 4.00 | 4.70 |
| Mini-Omni2 [72] | 4.8 | 9.8 | 4.7 | 9.4 | 3.58 | 3.54 | 2.62 | 3.08 | 3.20 |
| VITA-1.5 [23] | 3.3 | 7.2 | 3.4 | 7.5 | 4.83 | 5.48 | 4.91 | 2.92 | 4.54 |
| IXC2.5-OmniLive [81] | 2.5 | 5.7 | 2.6 | 5.8 | 1.60 | 1.77 | 1.74 | 1.58 | 1.67 |
| *Ola* (Pure audio) | 2.1 | 4.7 | 2.1 | 4.6 | 6.32 | 5.43 | 5.76 | 5.83 | 5.84 |
| *Ola* | **1.9** | **4.4** | **1.9** | **4.2** | **7.34** | **6.40** | **5.90** | **6.01** | **6.41** |

capabilities.

**Results.** We present the comprehensive results in Table 1, highlighting *Ola*'s competitive performance across major multi-modal benchmarks when compared to state-of-the-art specialist-modal LLMs. Specifically, in image benchmarks, *Ola* achieves 84.3% on MMBench-1.1 [40], 70.8% on MMStar [6], 57.0% on MMMU [77], 68.4% on MathVista [46], 86.1% on AI2D [29] and 827 on OCRBench [41], surpassing all the relative multi-modal LLMs in similar number of parameters. In video benchmarks, *Ola* attains an impressive 68.4% on VideoMME [21], showcasing its robust capability to handle both video and audio inputs simultaneously, and setting a new state-of-the-art performance among 7B models on the VideoMME benchmark. *Ola* also maintains a leading position compared to mainstream video LLMs including LLaVA-Video [84] and VideoCCAM [20] on LongVideoBench [69] and MVBench [34]. In audio benchmarks, *Ola* demonstrates strong audio-speech recognition and conversational abilities, with 3.1% mean WER rate on LibriSpeech [54] and 6.41 mean score on AIR-Bench [74], outperforming existing omni-modal LLMs, including LLaMA-Omni [17], which focuses on audio understanding. These results indicate a significant advancement over current omni-modal LLMs, underscoring the effectiveness of *Ola*'s training approach.

### 4.3. Analysis

For the analysis part, we report the detailed results on audio benchmarks to illustrate the fine-grained performance. We also demonstrate our designs on training and cross-modal training data with ablations on critical benchmarks. At last, we perform qualitative showcases of *Ola*.

**Analysis on Audio Benchmarks.** To demonstrate the effectiveness of our approach on audio and speech tasks, we conducted experiments using the LibriSpeech [54] and AIR-Bench [74] datasets, and we illustrate our results on Tab. 2.

Table 3. **Analysis on Omni-Modal Training.** We conduct analysis for the performance before/after omni-modal learning, and show the performance gain with audio inputs on mainstream video benchmarks. The highlighted row indicates the final accepted strategy.

| Omni-Stage Training | Audio | Subtitle | VideoMME | | | |
|---|---|---|---|---|---|---|
| | | | Short | Medium | Long | Overall |
| ✗ | ✗ | ✗ | 75.9 | 61.2 | 54.3 | 63.8 |
| ✓ | ✗ | ✗ | 76.4 | 61.9 | 54.8 | 64.4 |
| ✓ | ✗ | ✓ | 78.4 | 66.6 | 56.4 | 67.1 |
| ✓ | ✓ | ✗ | 78.7 | 68.3 | 58.3 | 68.4 |
| ✓ | ✓ | ✓ | 78.8 | 68.8 | 60.3 | 69.3 |

Table 4. **Analysis on Cross-Modal Training Data.** We analyze our data mixture for the cross-modal video-audio alignment data about sources from academic or open-ended videos. The highlighted row indicates the final accepted strategy. The experiment is conducted on a subsampled training set.

| Acadamic | Open-End | MMMU | VideoMME | LongVideo | LibriSpeech↓ |
|---|---|---|---|---|---|
| ✗ | ✗ | 48.2 | 59.0 | 56.4 | 4.5 |
| ✓ | ✗ | **48.3** | 64.2 | 56.8 | 4.1 |
| ✓ | ✓ | 48.1 | **65.7** | **57.4** | **4.0** |

Table 5. **Analysis on Progressive Modality Learning.** We evaluate the basic performance on image and video understanding for the intermediate models during the training stage. The highlighted row indicates the final accepted strategy.

| Stage1 | Stage2 | Stage3 | MMBench-1.1 | MMMU | OCRBench | VideoMME |
|---|---|---|---|---|---|---|
| ✓ | ✗ | ✗ | 83.5 | **57.5** | 832 | ✗ |
| ✓ | ✓ | ✗ | 83.8 | 57.2 | 820 | 63.8 |
| ✓ | ✓ | ✓ | **84.3** | 57.0 | 827 | **68.4** |

Specifically, we report the Word Error Rate (WER) on the test-clean, test-other, dev-clean, and dev-other subsets of LibriSpeech. Additionally, we present the GPT-4-eval scores on speech, sound, music, and mix sub-metrics in AIR-Bench. Our model, *Ola*, is compared against state-of-the-art audio models and omni-modal LLMs.

Notably, *Ola* demonstrates a significant advantage over existing omni-modal models, achieving a 1.9 WER on the test-clean subset of LibriSpeech and a 6.41 average score on AIR-Bench. This is in contrast to the previous state-of-the-art omni-modal models, which achieved a 2.5 WER and a 4.70 score, respectively. *Ola*'s performance is even approaching that of audio-specific models, highlighting its strong universality.

Furthermore, we evaluated *Ola* under two situations. *Ola* (Pure audio) indicates that we omit video-audio data in stage 3 and replace the same amount of data with pure audio inputs, where we can observe a consistent performance gain with cross-modal joint learning. Despite the significant distribution gap between video audio and speech-relevant datasets, this improvement indicates the robust connections between video and speech modalities.

**Speech Understanding**

*Audio:* (Oh Annie, we're going to live now. I'm going to make you so happy.)

*Text:* Given the speaker's emotional state and the content of the speech, what could be the reason behind his happiness?

*Ola:* The speaker's happiness could be due to a significant positive event or achievement in his life. The phrase **"we're going to live now"** suggests a sense of liberation or newfound freedom, while **"I'm going to make you so happy"** indicates a deep emotional connection and the desire to share joy with someone else.

**Visual Understanding**

*Video:* 

*Audio:* (You know I was saying that representing my country has been always the greatest priority and honor for me whether it's Olympic Games or Davis Cup or any anywhere that I can carry proudly the Serbian flag ...)

*Text:* Who is the man in the video and why the man in the video is wearing the nation flag? explain in detail.

*Vision LLM:* The man in the video is a **player** who is wearing the **nation flag**. He is sitting at a table with microphones and water bottles, indicating that he is likely participating in a **press conference or interview** ...

*Ola:* The man in the video is a **professional tennis player named Djokovic** who is representing Serbia. He is wearing the **national flag of Serbia**, which signifies his loyalty and pride to his country. The fact that he has **won gold in the Olympic Games** further emphasizes his commitment to representing Serbia at the highest level of international competition ...

Figure 5. **Generative results on speech and visual understanding tasks.** We illustrate results on speech and video understanding and show the strong ability of omni-modal *Ola* compared with conventional vision-language models.

**Effectiveness of Omni-Modal Training.** In exploring the relationships between video and audio, we examined the effectiveness of omni-modal training and its impact on audio within videos, both during training and in benchmark outcomes. The analysis results are shown in Tab. 3. By comparing results before and after omni-modal training (i.e., stage 3 of the progressive training strategy), we observed performance improvements from 63.8% to 64.4% on VideoMME [21]. Additionally, incorporating audio modalities alongside raw video resulted in significant performance gains, increasing scores from 64.4% to 68.4% on VideoMME [21]. These findings suggest that audio contains valuable information that enhances overall recognition performance. Notably, the multiple-choice accuracy for *Ola* with omni-modal training and audio inputs even surpasses the results with original text subtitles, with 68.4% overall performance compared with 67.1% overall performance. The results indicate that audio data may include more information beyond the original text information in certain benchmarks.

**Effectiveness of Progressive Modality Learning.** To evaluate the effectiveness of the proposed training strategy, we evaluate the basic performance of the intermediate model in each stage (Stage-1 for Ola-Image, Stage-2 for Ola-Video and Stage-3 for the final *Ola* model). Specifically, we adopt the representative MMBench-1.1 [40], MMMU [77], and OCRBench [41] for image performance and VideoMME [21] for video performance. Results are shown in Tab. 5. We can observe that the progressive modality training from image, video to audio can maximally preserve the previously learned capability. Additionally, we can also observe a performance gain on MMBench-1.1 and VideoMME, revealing the superiority of joint learning.

**Ablations on Cross-Modal Training Data.** In our implementations, we collect the cross-modal video-audio data for modality alignment from multiple sources including academic datasets and open-ended videos from YouTube. While the data distribution and the processing pipeline vary for the two sources, we conduct ablation analysis on the combination of dual video sources. Results are shown in Tab. 4 Our baseline model excluded video-audio training, focusing solely on audio-relevant data in Stage 3. Results indicate that video-audio training minimally affects image benchmarks, suggesting stable image understanding post-text-image training. For video benchmarks, we observed consistent performance improvements: 59.0% without video training, rising to 64.2% with academic data and 65.7% with open-ended data in VideoMME [21]. Furthermore, ASR performance on LibriSpeech [54] improved with video-audio data, likely due to the challenging subtitling tasks in complex environments, enhancing speech recognition capabilities.

**Showcases.** We present the qualitative generation results of the *Ola* model on both speech and visual understanding tasks in Fig. 5. For speech understanding, we utilize sources from the AIR-Bench [74] benchmark, where *Ola* demonstrates precise speech recognition and effective emotional analysis, as well as reasoning about the questions posed. *Ola* tries to find out the answer for the instruction based on the audio inputs. In visual understanding tasks, we analyze an interview with a famous tennis player after the Olympic Games. We compare our approach with the state-of-the-art vision-language model [32]. The large vision-language models without audio inputs exhibit significant information loss due to the absence of audio input and recognition capabilities for audio. In contrast, *Ola*, with its omni-modal inputs, provides more accurate responses regarding nationality, context, and background from the speaker's dialogue.

## 5. Conclusion

In this paper, we propose *Ola*, a comprehensive and powerful omni-modal language model that achieves competitive performance in image, video, and audio understanding tasks.

Our solution, based on a progressive modality alignment strategy, offers a natural, efficient, and competitive pipeline for training omni-modal models. Enhancements in architectural design with omni-modal inputs and streaming decoding, along with high-quality cross-modal video data preparation, further extend *Ola*'s capabilities. We hope our work inspires future research on more general AI models.

## Appendix

## A. More Details

We provide more details that are not implemented in the main paper. Specifically, we provide the detailed architecture for the model, more training details about *Ola*'s progressively modality alignment, and details of the data preparation procedure.

### A.1. Model Details

The visual encoder for *Ola* is based on the SigLIP-400M [79] backbone and is further fine-tuned for native-resolution visual inputs. The patch size for the visual encoder is set to 16, the hidden dimension is 1152, and the MLP hidden dimension is 4304. The SigLIP-400M model consists of 27 transformer blocks and 16 attention heads. The audio encoder for *Ola* is built on the whisper-v3-large [58] model. The length of the input audio tensor for the whisper model is fixed at 480,000; therefore, we chunk the entire audio tensor into pieces and concatenate the audio features. The mel size for the whisper-v3-large model is set to 128, and the hidden dimension for the speech features after the whisper-v3-large model is 1280.

For the connector layer, we utilize a 2-layer MLP for feature projection. We initialize two separate MLPs for visual and audio features, respectively. The input dimension matches that of the visual or audio encoder, and the output dimension matches the LLM dimension. For the *Local-Global Attention Pooling* layer, we use a predictor to calculate the score based on the concatenated features. Therefore, the dimension for the predictor is $2\times$ to $1\times$ dimension.

We integrate Qwen-2.5-7B [64] into the *Ola* large language model. The Qwen-2.5-7B model has a hidden LLM dimension of 3,584 and an intermediate size of 18,944. It consists of 28 transformer layers. The basic architecture for the speech decoder mirrors that of the Qwen-2.5-7B LLM, but we use only a 2-layer transformer block for the speech decoder. During generation, the maximum number of classification categories for the unit speech tensor is set to 1,000. We use a pre-trained speech vocoder to convert the unit speech tensor into speech waveforms. During inference, speech outputs are generated whenever a punctuation mark is detected, ensuring that the features of streaming outputs are preserved.

## A.2. Training Details

As stated in the main paper, the progressive modality alignment procedure is conducted in four stages. The first image-text stage involves adapter pre-training and supervised fine-tuning. The adapter pre-training stage is conducted on 808k image captioning data collected from the LAION datasets. During pre-training, we unfreeze the parameters for the connector while keeping other parameters frozen. We set the training batch size to 256 and the overall learning rate to 1e-3. The supervised fine-tuning stage is conducted on 7.3 million image-text data pairs. During image-text training, input images are maintained at their original aspect ratio, with the maximum image size restricted to 1536×1536. We set the batch size to 128 and the overall learning rate to 2e-5. The stage 1 experiment is conducted on 64 NVIDIA A800 GPUs.

Stage 2 integrates both image and video data for supervised fine-tuning, following most of the training strategies from Stage 1. The total amount of training data is 2.7 million, comprising 800k image-text pairs sampled from Stage 1 and 1.9M video data collected from open-source datasets. We set the training batch size to 256 and the overall learning rate to 2e-5. The model's maximum sequence length is set to 16k. The maximum number of frames is set to 64. The Stage 2 experiment is conducted on 64 NVIDIA A800 GPUs.

Stage 3 involves joint training in the audio domain. We conduct a projector alignment procedure similar to the image pre-training for initializing the speech adapter. During the speech adapter pre-training, we unfreeze the parameters of the speech adapter while freezing the other parameters. We set the batch size to 256 and the overall learning rate to 1e-3. The pre-training phase is conducted on the 370k LibriTTS [78] dataset. After pre-training, we integrate audio-video joint alignment by combining image data, video data, and pure audio data. We use 600k image data, 1.1M audio data, and 243k video-audio training data. The training batch size is set to 128, and the overall learning rate is set to 1e-5. We maintain the original audio data for inputs in the audio and video data and append the necessary prompts for instructions. Specifically, for the ASR tasks, we set the ASR prompt as *"Please give the ASR results of the given speech."* For the audio instruction tasks, we set the instruction-following prompt as *"Please directly answer the questions in the user's speech."* We maintain the maximum frame number for the video and set the maximum speech chunk number to 20. The Stage 3 experiment is conducted on 64 NVIDIA A800 GPUs.

## A.3. Data Collection Details

We provide details on collecting video-audio relevant data. Our data comes from two sources with high-quality raw videos: LLaVA-Video-178k [83], which contains 178k raw videos, and FineVideo [18], which contains 42k raw videos.

For the open-sourced video data from LLaVA-Video-178k, we first use the Whisper [58] model to generate subtitles. We find that the videos include content in other languages and videos without valid audio, so we design a filtering method for better results. Specifically, we first assess the ratio of English words in the generated subtitles and discard those with a lower ratio, indicating subtitles in other languages. We also discard extremely short subtitles. Then, we use a large language model, Qwen-2.5-72B, to further filter the subtitles. The model is asked to identify meaningless sentences with the following prompt: *"I will give you a subtitle generated from a video. Identify whether the subtitle is complete, fluent, and informative. Answer directly with yes or no and do not add other explanations."* After this procedure, we gathered 41k valid videos. For the videos in FineVideo, as they are already well-processed, we directly use the subtitles for the following steps. We utilize Qwen-2-72B to generate audio-relevant question-answer pairs based on the given videos and subtitles. The prompt for the vision-language model is: *"Please generate at least three questions and answers based on the information in the subtitle. You can refer to the video for additional context. The questions and answers must be highly relevant to the subtitle and video and should not include fabricated content."* We then generate 243k cross-modal video-audio data points from the 81k collected videos. This data is used for stage 3 training for omni-modal alignment.

## B. More Showcases

### B.1. Text and Audio Understanding

In this subsection, we provide more practical text-audio understanding samples for visualizations. The inputs of the text-audio understanding are a mixture of audio and text instructions, which can strongly test the cross-modal capability for *Ola* model. Results are shown in Fig. 6, we provide results on music-related, speech-related, and sound-related audio inputs and *Ola* excels at all the circumstances with a strong performance on mixed audio and text understanding.

### B.2. Video Understanding

In this subsection, we provide more results on video understanding tasks and provide comparisons with state-of-the-art vision LLM. Results are shown in Fig. 7. With the capability to recognize video, audio, and text jointly, *Ola* can gather more information from the video.

## References

[1] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023. 5

**Music-Related Data**

*Audio*: *[This is a low quality recording of a motown style song. The song seems to have been recorded using a cellphone. The lead female vocalist has a soulful and powerful voice. The song is upbeat and jovial.]*
*Text*: How is the female vocalist's voice described?

*Ola*: The female vocalist's voice is described as having a **strong, powerful, and emotive quality**.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*Audio*: *[A drone sound is playing in the lower to mid range as an intro before the song starts with a high pitched bell sound with a very long release time and acoustic drums that sounds like playing behind a closed door. This song may be playing in a club.]*
*Text*: What could be the purpose of introducing a drone sound in the lower to mid range in the beginning of a track that is meant to be played in a club?

*Ola*: The drone sound in the lower to mid range can serve several purposes in a club track. It can **create a foundation for the beat**, providing a steady and consistent rhythm that helps keep the energy high. The drone can also **add depth and texture to the track**, making it more engaging and dynamic. Additionally, it can **help to build tension and anticipation**, leading up to drops or other climactic moments in the track. Overall, the drone sound can enhance the overall feel and impact of the track in a club setting.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Speech-Related Data**

*Audio*: *Of course I did want to a little further up the coast you know get away from all the lights and people and everything. Is it midnight, do they always start at midnight? Is that what it is midnight? How you doing, huh? You okay? That's good.*
*Text*: Where did the speaker want to go?

*Ola*: The speaker wanted to go **a little further out** at the coast.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*Audio*: *Well the plan was to kind of sneak up on you over a period of a week or so. But uh- they take it for granted that we're all set.*
*Text*: Why might the speaker have intended to employ a gradual approach in his plan?

*Ola*: The speaker might have intended to employ a gradual approach to make the plan less noticeable and to avoid drawing attention to himself. By sneaking up on someone over a period of a week or so, he could have hoped to blend in more effectively and reduce the risk of being discovered.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Sound-Related Data**

*Audio*: *['Outdoors, bugs and flies and birds calling and eventually some clunking.', 'Various birds chirp as flies and crickets buzz.', 'Bugs, flies and birds are calling, eventually making a clunking noise outdoors.', 'Faint chattering of people in the far distance, buzzing of a flying Insect and a knock of wood.', 'People talk in the distance while a piece of hard metal gently strikes a hard surface.']*
*Text*: What type of sound does the insect make?

*Ola*: The insect makes a buzzing sound.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*Audio*: *['A washing machine is turned on and then feet walk across the wooden floor.', 'A person shoveled dirt or sand then walked upon a boardwalk.', 'The washer and dryer are being ran at the same time.', 'a person is shoveling dirt or sand then walks a boardwalk', 'Turned on washing machine followed by walking sounds on wood.']*
*Text*: What kind of floor is the person walking on?

*Ola*: The person is walking on a wooden floor.

Figure 6. **Showcases on Text and Audio Understanding.**

[2] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024. 6, 7

[3] 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. *NeurIPS*, 35: 23716–23736, 2022. 3

[4] Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang Zhang, Jing Shi, Shuang Xu, and Bo Xu. X-llm: Bootstrapping advanced large language models by treating multi-modalities as foreign languages. *arXiv preprint arXiv:2305.04160*, 2023. 3

[5] 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. 5

[6] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024. 7

[7] 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. 3

**Video Understanding**

**Video:**

**Audio:** *(Let's make our yogurt gravy base. So in my half cup of yogurt here, you can see it's nice and thick. We'll go ahead and add half a cup of water. And along with water, we'll add some spices. So for spices, half a teaspoon of red chili powder, half a teaspoon of turmeric … )*
**Text:** *Describe the steps of the process shown in the video.*

**Vision LLM:** The video begins with a dark-colored bowl on a countertop, where the person pours a thick, creamy liquid from a spoon into the bowl. The liquid appears to be a type of sauce or dressing, creating a smooth texture. The person continues to stir the mixture in the bowl, ensuring it is well-mixed and consistent in appearance...

**Ola:** The video shows the process of **making a yogurt** by adding ingredients to a bowl and mixing them together. The first step is to add 1 cup of yogurt, followed by 1 cup of water, and then 2 tablespoons of tomato paste. After that, 1 teaspoon of turmeric powder and 1 teaspoon of cumin powder are added. The mixture is stirred well to combine all the ingredients.

**Video:**

**Audio:** *(Firstly, congratulations Ready for quarterfinals and two goals. What do you say about today's performance? We are really happy with the results, with the qualification. We wanted to go to the quarterfinal and we did a nice game. It was tough because it's never easy to come and to win an away game. But it's what we did and now we are ready for the next one … )*
**Text:** *What is the soccer player in the video talking about?*

**Vision LLM:** The player is speaking into a microphone with the Champions League logo.

**Ola:** The soccer player in the video is expressing his gratitude for the support he received from his fans and teammates after a victory of the match. He also mentions that he will be back on the field next week.

Figure 7. **Showcases on Video Understanding.**

[8] Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. Beats: Audio pre-training with acoustic tokenizers. 2022. 4, 6

[9] 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. 6, 7

[10] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *CVPR*, pages 24185–24198, 2024. 2, 3

[11] 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. 2, 3, 6, 7

[12] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024. 5

[13] Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu. Insight-v: Exploring long-chain visual reasoning with multimodal large language models. *arXiv preprint arXiv:2411.14432*, 2024. 3

[14] Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: An audio captioning dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 736–740. IEEE, 2020. 5

[15] Zhihao Du, Qian Chen, Shiliang Zhang, Kai Hu, Heng Lu, Yexin Yang, Hangrui Hu, Siqi Zheng, Yue Gu, Ziyang Ma, et al. Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens. *arXiv preprint arXiv:2407.05407*, 2024. 4

[16] 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. 3

[17] Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. Llama-omni: Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*, 2024. 2, 3, 4, 6, 7

[18] Miquel Farré, Andi Marafioti, Lewis Tunstall, Leandro Von Werra, and Thomas Wolf. Finevideo. https://huggingface.co/datasets/HuggingFaceFV/finevideo, 2024. 6, 10

[19] Yassir Fathullah, Chunyang Wu, Egor Lakomkin, Junteng Jia, Yuan Shangguan, Ke Li, Jinxi Guo, Wenhan Xiong, Jay Mahadeokar, Ozlem Kalinli, et al. Prompting large language models with speech recognition abilities. In *ICASSP 2024-*

*2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 13351–13355. IEEE, 2024. 3

[20] Jiajun Fei, Dian Li, Zhidong Deng, Zekun Wang, Gang Liu, and Hui Wang. Video-ccam: Enhancing video-language understanding with causal cross-attention masks for short and long videos. *arXiv preprint arXiv:2408.14023*, 2024. 6, 7

[21] 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. 2, 3, 7, 8, 9

[22] Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, et al. Vita: Towards open-source interactive omni multimodal llm. *arXiv preprint arXiv:2408.05211*, 2024. 2, 3, 7

[23] Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Yangze Li, Zuwei Long, Heting Gao, Ke Li, et al. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. *arXiv preprint arXiv:2501.01957*, 2025. 6, 7

[24] GeminiTeam. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 2, 3

[25] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024. 7

[26] Jarvis Guo, Tuney Zheng, Yuelin Bai, Bo Li, Yubo Wang, King Zhu, Yizhi Li, Graham Neubig, Wenhu Chen, and Xiang Yue. Mammoth-vl: Eliciting multimodal reasoning with instruction tuning at scale. *arXiv preprint arXiv:2412.05237*, 2024. 5

[27] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *NeurIPS*, 36:20482–20494, 2023. 3

[28] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audiogpt: Understanding and generating speech, music, sound, and talking head. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 23802–23804, 2024. 3

[29] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *ECCV*, pages 235–251. Springer, 2016. 7

[30] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 119–132, 2019. 5

[31] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models?, 2024. 3, 5

[32] 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, 5, 6, 7, 9

[33] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742. PMLR, 2023. 3

[34] 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 *CVPR*, pages 22195–22206, 2024. 7

[35] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024. 4

[36] 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. 2, 3

[37] Benlin Liu, Yuhao Dong, Yiqin Wang, Yongming Rao, Yansong Tang, Wei-Chiu Ma, and Ranjay Krishna. Coarse correspondence elicit 3d spacetime understanding in multimodal language model. *arXiv preprint arXiv:2408.00754*, 2024. 3

[38] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *CVPR*, pages 26296–26306, 2024. 4, 5

[39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36, 2024. 3, 4

[40] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023. 2, 3, 7, 9

[41] 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. 7, 9

[42] 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. 3, 4, 6

[43] Zuyan Liu, Yuhao Dong, Yongming Rao, Jie Zhou, and Jiwen Lu. Chain-of-spot: Interactive reasoning improves large vision-language models. *arXiv preprint arXiv:2403.12966*, 2024. 3

[44] Zuyan Liu, Benlin Liu, Jiahui Wang, Yuhao Dong, Guangyi Chen, Yongming Rao, Ranjay Krishna, and Jiwen Lu. Efficient inference of vision instruction-following models with elastic cache. *arXiv preprint arXiv:2407.18121*, 2024. 3

[45] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 3

[46] 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. 7

[47] 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, 5

[48] Brian McFee, Thierry Bertin-Mahieux, Daniel PW Ellis, and Gert RG Lanckriet. The million song dataset challenge. In *Proceedings of the 21st International Conference on World Wide Web*, pages 909–916, 2012. 5

[49] Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multi-modal research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024. 5

[50] OpenAI. Openai gpt-3.5 api. *OpenAI API*, 2023. 3

[51] OpenAI. Gpt-4v(ision) system card. *OpenAI Blog*, 2023.

[52] OpenAI. Gpt-4 technical report. *ArXiv:abs/2303.08774*, 2023. 3, 7

[53] OpenAI. Hello gpt-4o — openai. *OpenAI Blog*, 2024. 2, 3

[54] 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, 5, 7, 9

[55] 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. 2

[56] QwenTeam. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 1

[57] QwenTeam. Qwen2-vl: To see the world more clearly. *Wwen Blog*, 2024. 2, 3, 6

[58] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. 4, 6, 7, 9, 10

[59] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Am-radio: Agglomerative vision foundation model reduce all domains into one. In *CVPR*, pages 12490–12500, 2024. 3

[60] Ruchit Rawal, Khalid Saifullah, Ronen Basri, David Jacobs, Gowthami Somepalli, and Tom Goldstein. Cinepile: A long video question answering dataset and benchmark. *arXiv preprint arXiv:2405.08813*, 2024. 5

[61] Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. Audiopalm: A large language model that can speak and listen. *arXiv preprint arXiv:2306.12925*, 2023. 3

[62] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 5

[63] 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, 5, 6, 7

[64] Qwen Team. Qwen2.5: A party of foundation models, 2024. 5, 6, 9

[65] Qwen Team. Qwen2.5-vl, 2025. 6, 7

[66] John Thickstun, Zaid Harchaoui, and Sham Kakade. Learning features of music from scratch. *arXiv preprint arXiv:1611.09827*, 2016. 5

[67] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024. 2, 5, 6, 7

[68] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 3

[69] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*, 2024. 7

[70] Jian Wu, Yashesh Gaur, Zhuo Chen, Long Zhou, Yimeng Zhu, Tianrui Wang, Jinyu Li, Shujie Liu, Bo Ren, Linquan Liu, et al. On decoder-only architecture for speech-to-text and large language model integration. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE, 2023. 3

[71] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 3

[72] Zhifei Xie and Changqiao Wu. Mini-omni2: Towards open-source gpt-4o model with vision, speech and duplex. *arXiv preprint arXiv:2410.11190*, 2024. 2, 3, 4, 6, 7

[73] Jingkang Yang, Yuhao Dong, Shuai Liu, Bo Li, Ziyue Wang, Haoran Tan, Chencheng Jiang, Jiamu Kang, Yuanhan Zhang, Kaiyang Zhou, et al. Octopus: Embodied vision-language programmer from environmental feedback. In *European Conference on Computer Vision*, pages 20–38. Springer, 2024. 3

[74] Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv, Zhou Zhao, Chang Zhou, et al. Air-bench: Benchmarking large audio-language models via generative comprehension. *arXiv preprint arXiv:2402.07729*, 2024. 3, 7, 9

[75] 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. 6, 7

[76] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024. 1

[77] 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 multi-modal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 3, 7, 9

[78] Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. Libritts: A corpus derived from librispeech for text-to-speech. *arXiv preprint arXiv:1904.02882*, 2019. 10

[79] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023. 4, 6, 9

[80] 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, 7

[81] Pan Zhang, Xiaoyi 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, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2.5-omnilive: A comprehensive multimodal system for long-term streaming video and audio interactions. *arXiv preprint arXiv:2412.09596*, 2024. 6, 7

[82] Ruohong Zhang, Liangke Gui, Zhiqing Sun, Yihao Feng, Keyang Xu, Yuanhan Zhang, Di Fu, Chunyuan Li, Alexander Hauptmann, Yonatan Bisk, et al. Direct preference optimization of video large multimodal models from language model reward. *arXiv preprint arXiv:2404.01258*, 2024. 5

[83] 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. 10

[84] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data, 2024. 2, 5, 6, 7