

BAICHUAN-OMNI TECHNICAL REPORT

Yadong Li^{1*} Haoze Sun^{1*} Mingan Lin^{1*} Tianpeng Li^{1*} Guosheng Dong^{1*}
Tao Zhang¹ Bowen Ding^{2,3} Wei Song^{2,3} Zhenglin Cheng^{2,3} Yuqi Huo¹
Song Chen¹ Xu Li¹ Da Pan¹ Shusen Zhang¹ Xin Wu¹ Zheng Liang¹
Jun Liu¹ Tao Zhang¹ Keer Lu¹ Yaqi Zhao¹ Yanjun Shen¹ Fan Yang¹
Kaicheng Yu² Tao Lin² Jianhua Xu^{1†} Zenan Zhou^{1†} Weipeng Chen¹
¹ Baichuan Inc. ² Westlake University ³ Zhejiang University
{xu.jianhua, zhou.zenan}@baichuan-inc.com

<https://github.com/westlake-baichuan-mlm/bc-omni>

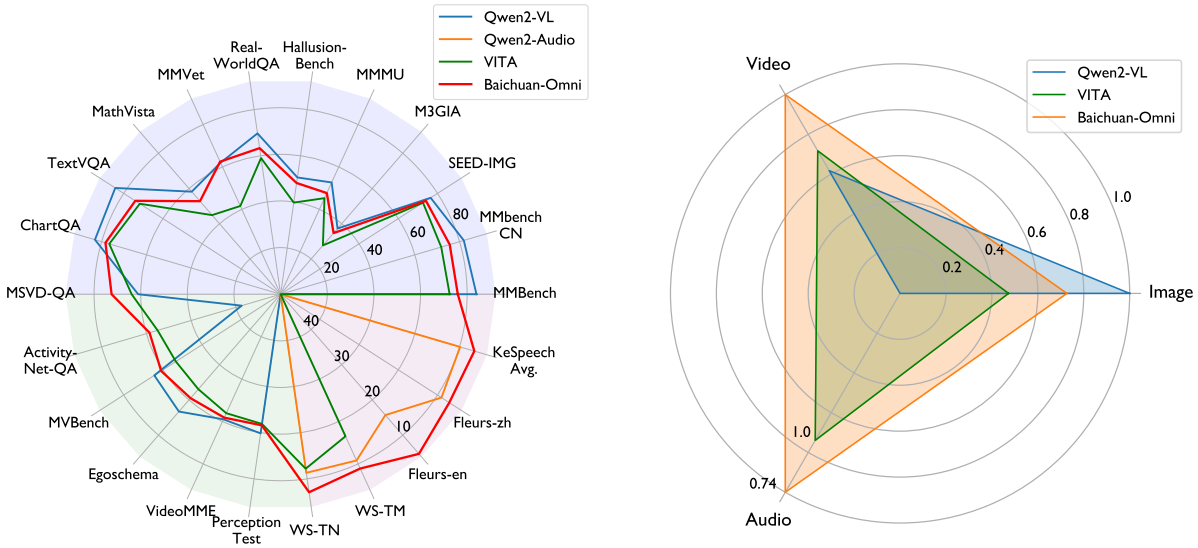


Figure 1: **Evaluation across image, video, and audio modalities.** (Left) Baichuan-omni covers more modalities than Qwen2 VL [79] and outperforms the current leading omni-modal model, VITA [24]. (Right) Average scores across benchmarks for all modalities. All the scores are normalized by $x_{\text{norm}} = (x - x_{\text{min}} + 10)/(x_{\text{max}} - x_{\text{min}} + 10)$.

ABSTRACT

The salient multimodal capabilities and interactive experience of GPT-4o highlight its critical role in practical applications, yet it lacks a high-performing open-source counterpart. In this paper, we introduce **Baichuan-omni**, the first open-source 7B Multimodal Large Language Model (MLLM) adept at concurrently processing and analyzing modalities of image, video, audio, and text, while delivering an advanced multimodal interactive experience and strong performance. We propose an effective multimodal training schema starting with 7B model and proceeding through two stages of multimodal alignment and multitask fine-tuning across audio, image, video, and text modal. This approach equips the language model with the ability to handle visual and audio data effectively. Demonstrating strong performance across various omni-modal and multimodal benchmarks, we aim for this contribution to serve as a competitive baseline for the open-source community in advancing multimodal understanding and real-time interaction.

*Equal Core Contributors.

†Corresponding author.

1 Introduction

The burgeoning field of artificial intelligence has witnessed a remarkable evolution, especially with the development of Large Language Models (LLMs) [1, 8, 105] and the subsequent emergence of Multimodal Large Language Models (MLLMs) [46, 63, 93], signifying a paradigm shift in how machines understand and interact with the world. The introduction of MLLMs like GPT-4o [63], characterized by their exceptional multimodal capabilities and enriched interactive experiences, has not only spotlighted the indispensable role of these technologies in real-world applications but also set a new benchmark for what is achievable in terms of human-computer interaction.

Despite the remarkable progress of MLLMs, current open-source solutions exhibit notable deficiencies, particularly in multimodal capabilities and the quality of user interaction experiences [24]. These shortcomings significantly impede the broader adoption and effectiveness of such models in diverse applications, from natural language processing [18, 68] to computer vision [84, 73] and beyond.

In response to these challenges, we introduce an omni-modal LLM **Baichuan-omni** alongside a multimodal training scheme designed to facilitate advanced multimodal processing and naturalistic user interactions. The architecture of Baichuan-omni is depicted in Figure 2. The scheme of Baichuan-omni is built upon three core components:

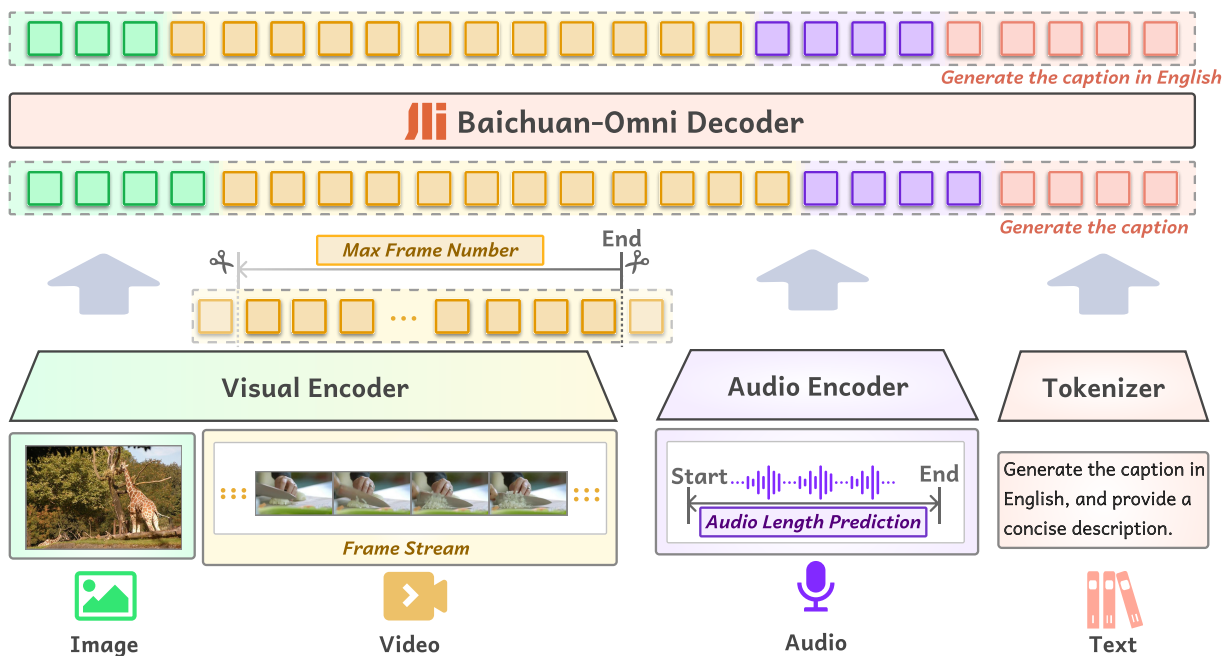


Figure 2: **Architecture of Baichuan-omni**. Baichuan-omni is designed to process both pure text/audio inputs and combinations of video/image with text/audio. In terms of interactivity, the model initially predicts the start and end of audio inputs. During this period, incoming images and videos are encoded into features and fed into the MLLM in a streaming fashion to calculate attention. The audio features are then input into the MLLM for inference following the end of the audio input, facilitating streaming input of audio and video.

Omni-Modal Data Construction. We utilize a substantial collection of high-quality, omni-modal data to train Baichuan-omni with a blend of open-source, synthetic, and internally annotated datasets. In the multimodal alignment pre-training phase, we curate a wide-ranging assortment of training corpora that encompasses image captions, interleaved data, OCR data, and image-text data. For audio alignment, we collect both open-source and in-house datasets for Automatic Speech Recognition (ASR) and Audio Question Answering (AQA). In the realm of video alignment, we acquire video data from both open-source and in-house sources. During the multimodal supervised fine-tuning phase, we compile and synthesize an extensive dataset that covers over 200 tasks and comprises 600,000 instances across pure text, audio, image-text, video-text, and image-audio interaction data.

Multimodal Alignment. During the pre-training phase for multimodal alignment, we meticulously align encoders and connectors across various modalities. Initially, we train the vision-language model using a substantial dataset of image-text pairs. This foundational training enables us to harness the visual capabilities developed during the image-text training to further train the video projector. Concurrently, we train the audio-language model utilizing Automatic

Speech Recognition (ASR) data. Building upon this robust foundation, we integrate high-quality image, audio, and video data to achieve comprehensive multimodal alignment.

Multitask Fine-tuning. For the omni-modal fine-tuning stage, we utilize a multi-task cross-modal interaction training corpus derived from a combination of open-source, synthetic, and internally annotated data. We select data for the final supervised fine-tuning (SFT) phase based on criteria that whether factual knowledge is already learned by the pre-trained model [27]. During this phase, we implement a packing technique to concatenate multiple samples, using the `cuseq_len` from `flash-attention2` for effective sample isolation. With this, multiple samples can be packaged into a large batch while ensuring each sample is correctly isolated during the computational process, preventing data confusion between different samples. This accelerates the training process and optimizes memory usage.

The contributions of this paper are summarized below:

- We introduce Baichuan-omni , an open-source, high-performance foundational omni-modal model capable of concurrently processing text, images, videos, and audio inputs. It also provides multilingual support for languages including English and Chinese. Our training framework features a comprehensive pipeline that includes the construction of omni-modal training data, multimodal alignment pre-training, and multimodal supervised fine-tuning, with a particular emphasis on enhancing omni-modal instruction-following capabilities.
- We explore early-stage research in natural multimodal human-computer interactions. Our approach initiates with the prediction of audio input boundaries, while simultaneously streaming and encoding incoming visual data into features. These features are then processed by a multimodal large language model (MLLM) for dynamic attention computation. Upon completion of the audio input, the corresponding features are fed into the MLLM for inference, thus facilitating the support for handling audio and video inputs. This integrated approach allows for real-time processing and enhances the interactive capabilities of the system.
- We have made our Baichuan-omni model, training code, and evaluation scripts publicly available, aiming at fostering progress within the research community. As pioneers in this field, we remain committed to furthering the development of multimodal foundational models and their interactions.

2 Related works

Recent advancements in Large Language Models (LLMs) have reshaped the AI landscape, paving the way for the emergence of Multimodal Large Language Models (MLLMs). These advanced models expand AI capabilities beyond text, allowing understanding and generation of content across multiple modalities, including images, audio, and video, signaling a significant leap in AI development.

Open-source MLLMs have demonstrated increasingly powerful capabilities, with efforts from both academia and industry fueling the rapid development of models. LLMs such as LLaMA [80, 81], MAP-Neo [101], Baichuan [89], Qwen [5, 90], and Mixtral [36] are trained on extensive text data, exhibiting strong capacities in natural language comprehension and task resolution through text generation. Vision-Language Models (VLMs) [43, 109, 102] have shown promising potential in addressing vision-focused issues, with representative models including LLaVA [51], DeepSeek-VL [54], the Qwen-VL series [6, 79], InternVL families [12, 11], and MiniCPM [32]. Additionally, Audio-Language Models (ALMs) [86, 17, 38] leverage audio-text pairs to perceive audio signals based on a singular audio encoder. Notable instances of these models encompass Qwen-Audio [15, 14], SALMONN [77], SpeechGPT [100], etc.

However, compared to proprietary models like GPT-4o [63], open-source models still exhibit substantial gaps in their capabilities for multimodal interactions, and there is a considerable scarcity of open-source models that effectively facilitate comprehensive multimodal interactions [24]. To address these, we propose Baichuan-omni , an open-source and capable MLLM which concurrently supports interactions across modalities including audio, image, video, and text.

3 Training

3.1 High-Quality Multimodal Data

For training an omni-modal model with strong ability, we build an extensive cross-modal dataset with high quality, including text, image-text, video-text, audio-text, and their interactions.

Image Data. Image data can be categorized into several types: Caption, Interleaved image-text, OCR data and Chart data [35]. From the perspective of sources, it is divided into Open-source data and Synthetic data. Regarding open-source data, we have collected major open-source datasets, including PIN-14M [83], MINT-1T [4], LAION-5B [70], OBELIC [39], etc. for Stage I training of Image-language branch (Detailed introduction in Section 3.2.1), and

Cauldron [40], Monkey [47], ArxivQA [45], TGDdoc [85], MM-Self-Instruct (Train split) [103], MMTAB [106], etc. for Stage II/III training of Image-language branch. These publicly available open-source datasets are subjected to a series of processing steps and careful sampling techniques within our data pipeline.

As for synthetic data, the purpose is to obtain higher quality data to enhance the performance of our models. One part is derived from books and papers, which are parsed to generate Interleaved image-text, OCR data and Chart data. It is highly complete and specialized, making it high-quality and knowledge intensive data. Another part involves training dedicated models to produce image captions. These captions describe the content of the images in detail from different perspective, belonging to high-quality caption data.

Video Data. Video dataset comprises a diverse array of publicly available resources, encompassing multiple tasks such as video classification, action recognition, and temporal localization. The video-text sources can be categorized into two main types: question-answering (QA) data and caption data.

For QA data, we incorporate: NExTVideo, introduced in LLaVA-NExT [104] and ActivityNet-QA (Train split) [95]. Our caption data sources include ShareGPT4Video [10], a large-scale dataset that leverages GPT-4 to generate rich, contextual captions for videos, and WebVid [7]. To further enrich our dataset, we have employed GPT-4o to generate diverse captions for videos collected from YouTube.

The sampling ratio for each dataset within our compilation is carefully determined based on the relative sizes of these datasets. This strategic approach ensures a balanced representation of various video types, tasks, and domains in our final dataset.

Audio Data. Considering the diversity of audio data, we extract audio from various media modalities, which includes different recording environments, languages, accents, and speakers. Guided by the principles in previous work [67], we posit that the variation in audio quality contributes to a robust speech understanding capability. To facilitate a more sophisticated classification and filtering procedure, we implemented a data processing pipeline comprising speaker voice recording, dialect recognition, accent recognition, sound effect detection, and quality assessment.

To enhance the quality of audio-text pairs derived from the dataset, we utilized an in-house ASR system along with several open-source models [67, 25, 75] to generate multiple transcript versions. These generated data are then refined through a model ensemble strategy for effective text filtering and error correction.

Text Data. In handling text corpus, we collected data from various domains such as web pages, books, academic papers, code, etc.. Following the data processing protocols proposed in previous works [19, 55], we implemented a selection process to enhance the diversity and quality of the dataset. The diversity criterion ensures broad coverage of topics and linguistic styles in the training corpus, accommodating various applications. High-quality processing removes redundancy and noise from the text data, increasing knowledge density.

Cross-Modal Interaction Data. To enhance the cross-modal interaction capabilities of our model, we synthesized a collection of visual-audio-text cross-modal interaction data, including both image-audio-text and video-audio-text datasets. For the image-text data, we segmented the textual data into a 1:3 ratio, converting the initial quarter of text into audio descriptions using text-to-speech (TTS) technology. Our audio encompasses 44 different timbres, ensuring a diversity of vocal tones. This setup is complemented by task prompts such as “Please listen to the following audio describing the content of the image. Your task is to supplement more information by integrating the image after listening”, aiming to predict the remaining three-quarters of the textual description. For the video-text data, we directly extracted the audio from the videos to serve as the cross-modal audio component.

3.2 Multimodal Alignment Pre-training

In this section, we will further illustrate the pre-training and alignment processes for the Image-Language, Video-Language, and Audio-Language branches.

3.2.1 Image-Language Branch

We utilize Siglip-384px [97] as the visual encoder, which processes a 384×384 image input and generates 182 tokens through a visual projector composed of a two-layer MLP and a 2×2 convolution layer serving as the pooling layer. To scale the input to arbitrary resolutions while preserving the intricate details of high-resolution images, we adopt AnyRes [50], which splits the image into grids and concatenates the features of a down-sampled image to provide global context. The training of our image-language branch is divided into three stages.

- **Stage I:** In the first stage, we train the visual projector to establish the initial alignment between image representations and text through image captioning task. During this phase, we freeze the LLM and the visual encoder, only training the visual projector with a learning rate of $1e - 3$.

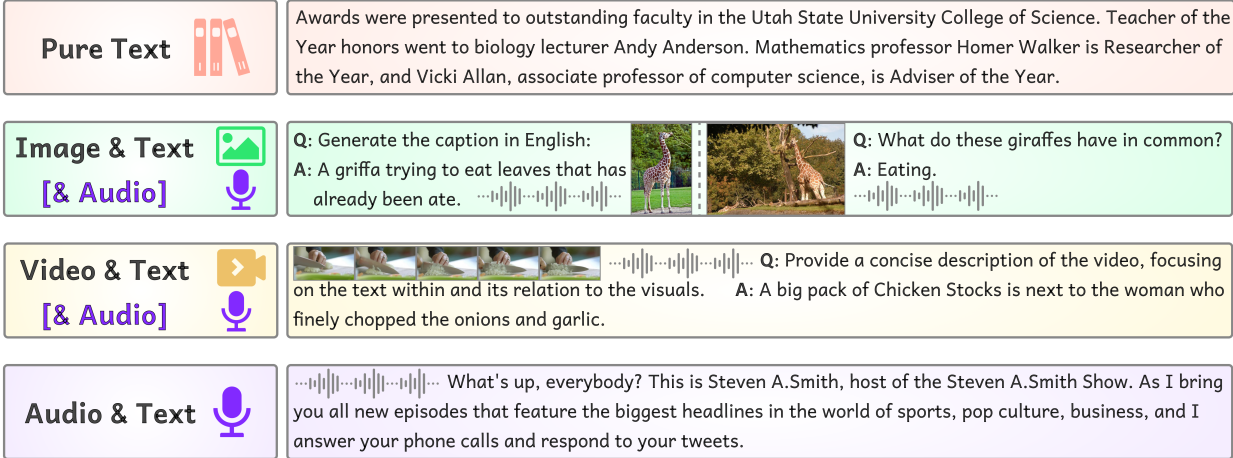


Figure 3: **Data illustration of Baichuan-omni** . We build an extensive cross-modal dataset, including text, image-text, video-text, audio-text, and their interactions. Our collection also features integrated image-audio-text and video-audio-text data.

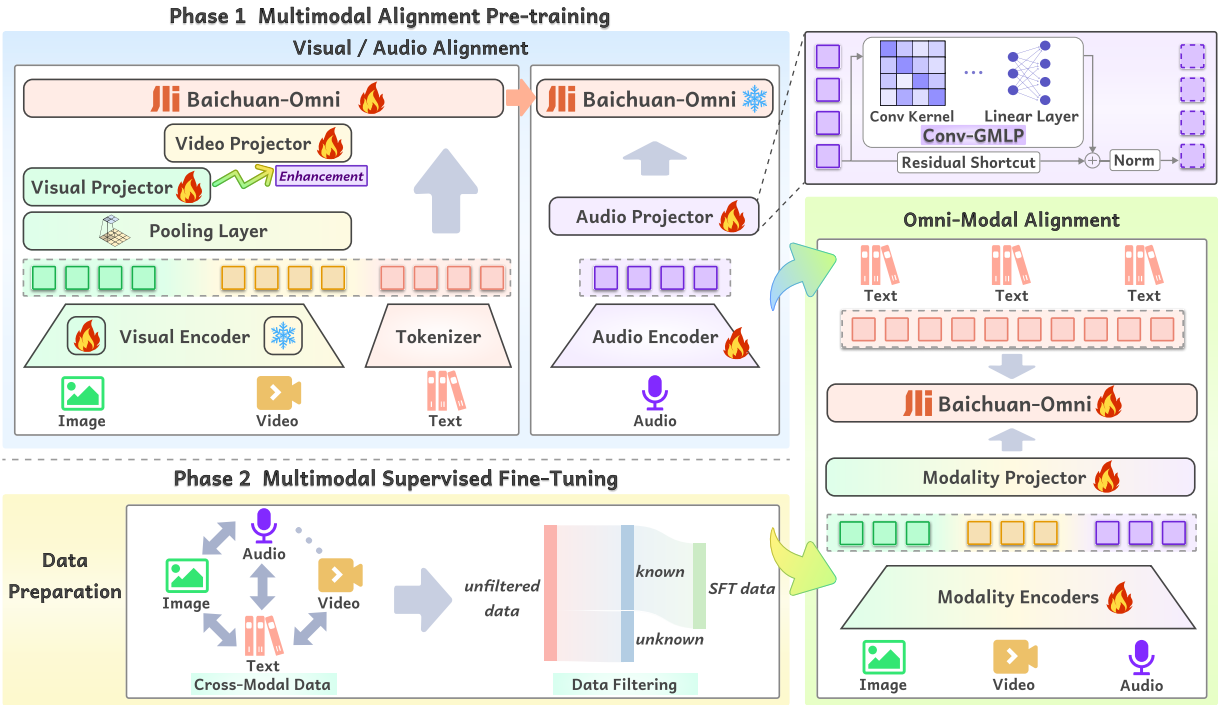


Figure 4: **Training Pipeline of Baichuan-omni** . During the pretraining phase, we initially train a vision-language model using extensive image-text data, followed by training an audio-language model with ASR data. Subsequently, we integrate high-quality images, audio, and video data for comprehensive multimodal alignment. In the fine-tuning phase, we synthesize a subset of cross-modal interaction data to blend with existing high-quality datasets. From this enriched dataset, we select a subset of data that the model is already capable of handling and proceed with multimodal multitask fine-tuning. This process aims to enhance the model’s adherence to omni-modal instructions.

- **Stage II:** In the second stage, we freeze the LLM and train both the projector and visual encoder with a smaller learning rate of $1e - 5$. In addition to general VQA tasks, we specifically synthesized 130k high-quality QA data for OCR and charts to enhance the model’s abstract visual understanding. We also introduced interleaved data and image

caption data in this stage, which help maintain and promote better alignment between image and text representations, mitigating shifts caused by changes in the image feature space after unfreezing the visual encoder.

- **Stage III:** Based on the second stage, we unfreeze the LLM and continue updating the parameters of all model components with a learning rate of $1e - 5$ to further enhance visual-language performance. In addition to VQA and image-caption pairs, we also introduce interleaved data and pure text data in this stage to better maintain the original capabilities of the LLM.

3.2.2 Video-Language Branch

Based on the visual capabilities acquired from the pre-training of the Image-Language Branch, we proceed to train the video projector using a frozen vision encoder (Siglip-384px, the same as that used in the Image-Language Branch) alongside an LLM (Large Language Model) backbone. This training process employs a low learning rate of $4e - 6$ to refine the alignment with the language modality.

During the training phase, the input video frames are sampled at a rate of 1 frame per second, with a maximum of 48 frames per video. Each input frame is resized to a maximum resolution of 384×768 pixels to maintain optimal quality and detail. Furthermore, a 2×2 convolution layer is applied prior to the video projector. This convolutional step serves to regulate the length of the video token sequence, ensuring a minimum of 182 tokens and a maximum of 546 tokens. This thoughtful configuration strikes a balance between performance and efficiency, facilitating effective model training while managing the computational load.

Rather than immediately proceeding with the pre-training of the Video-Language Branch using only pure video-text pairs, we have opted for a more nuanced two-stage approach. Initially, we leverage image-text pre-training data to strengthen the model’s visual understanding capabilities. After establishing a robust foundation, we incrementally integrate mixed image-text pairs and video-text pairs into the training regimen. This strategy has proven to yield superior results. By gradually enhancing the model’s visual competence, we provide valuable guidance for the video pre-training pipeline, allowing the model to better understand and integrate the complexities of video data in conjunction with language. This methodology underscores the importance of a comprehensive training strategy that incorporates diverse data modalities for improved alignment and performance.

3.2.3 Audio-Language Branch

The Audio-Language branch extends an LLM pre-trained on visual and video data by incorporating an audio encoder from the Whisper-large-v3 model [67] and a newly introduced audio projector.

The audio encoder processes the audio signal (30s, 128 mel-spectrum) into an audio representation in a 1280-channel feature space, while the audio projector (typically a linear projector [14] or MLP) maps that to the embedding space of LLM. Prior to projection, a pooling operation with a stride of n is traditionally used to down-sample the audio representation into fewer tokens (i.e., frames) for the downstream LLM. However, when we aggressively reduce the number of audio tokens, this simple pooling approach leads to a loss of audio information. In our approach, we replace the pooling with **Convolutional-Gated MLP** (Conv-GMLP), leveraging convolution layers for down-sampling to preserve more audio information.

Figure 5 illustrates the Conv-GMLP architecture, which functions similarly to a gated MLP [49] but replaces linear layers with convolutional ones. Each of the two convolutional layers reduces the sequence length of the audio representation by a factor of n , while proportionally expanding the feature space. In our projector, a residual shortcut is along with Conv-GMLP, enabling more efficient gradient back-propagation. Results in Section 4.5.3 demonstrate strong robustness in audio performance when setting the down-sampling rate³ n aggressively.

During training, the LLM remains frozen, and only the audio encoder and projector are trained using long audio-text sequences (up to 4K tokens). A cosine learning rate scheduler is employed to enhance performance.

3.2.4 Image-Video-Audio Omni-Alignment

The right part of Figure 4 illustrates the ‘Omni-Alignment’ stage, which follows the individual training of the Image-Language, Video-Language, and Audio-Language branches. During this stage, all modules are trained together on a mixture of high-quality image-text, video-text, and audio-text pairs to develop comprehensive multimodal understanding.

³The *down-sampling rate* is named as `avg_pooler` in our code configuration.

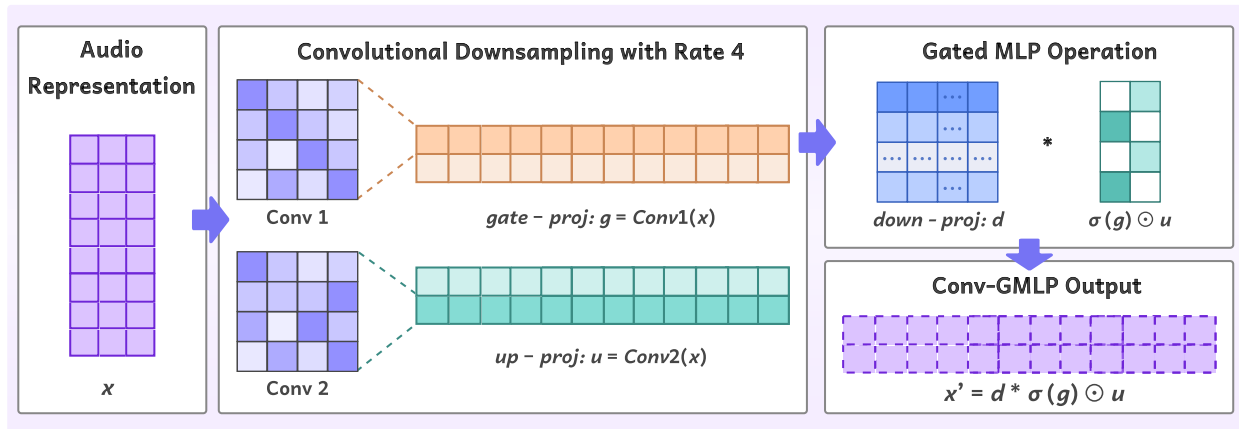


Figure 5: **Illustration of Conv-GMLP.** Conv-GMLP down-sampling is applied to the audio representation. With a down-sampling rate of 4, the output sequence length is reduced to a quarter of the input, while the number of feature channels increases fourfold.

3.3 Multimodal Supervised Fine-Tuning

In this section, we describe the multimodal supervised fine-tuning process aimed at improving the model’s ability to follow complex, multimodal instructions across various tasks. We leveraged a diverse set of open-source, synthetic, and internally annotated data, covering over 200 distinct tasks and comprising approximately 600K pairs across text, audio, image-text, video-text, and image-audio modalities.

Text-only Data. The text-only data covers a broad range of tasks, including knowledge-based question answering, mathematics, logical reasoning, code generation, text creation, information processing, persona-based tasks, and safety-related data. To further strengthen the model’s ability to handle complex, multi-step tasks, we included specialized datasets that feature intricate instructions, some of which contain a system message designed to structure more elaborate scenarios.

Image Understanding Data. For tasks involving image understanding, we primarily utilized the vFLAN dataset [9], focusing on its instruction-following data. Given the quality issues present in some of the samples, we employed a loss-based filtering method to clean the dataset:

1. We computed the loss for all vFLAN English instruction samples using the pretrained model and fit the resulting values to a Gaussian distribution.
2. Samples are removed if their loss values fell outside the range of $\mu \pm \sigma$.
 - (a) Samples with $loss < \mu - \sigma$ typically included trivial issues, such as cases where the prompt and response content are nearly identical.
 - (b) Samples with $loss > \mu + \sigma$ often had significant problems, such as reversed prompt-response pairs or hallucinations in the responses.

A subset of the cleaned vFLAN instruction data is then translated into Chinese, followed by manual re-annotation to ensure high-quality alignment. Alongside vFLAN, we incorporated several other open-source datasets, including synthdog-en/zh [37], handwritten OCR, street view OCR, reference grounding and grounded captioning duality tasks, and ImageInWords [26]. Most of these datasets are translated into Chinese. For ImageInWords, we ensured that if an image contained a recognizable entity, the corresponding caption explicitly referenced that entity by name (e.g., identifying a Samoyed dog by breed rather than simply labeling it as “dog”).

Although vFLAN covers 191 tasks, we found that it lacked variety in instruction types. To address this, we sampled data from our textual SFT dataset and rendered some of the prompts as images to increase the diversity of image-based instructions. Additionally, to enhance the model’s mathematical reasoning with images, we used the method from [108] to generate a large dataset of multimodal math problems involving images.

In experiments, we found that adding too much external world knowledge that the model inherently did not know resulted in diminishing performance returns. To mitigate this, we adopted the filtering method from [27] to exclude unknown data from the constructed SFT dataset.

Video Understanding Data. The video-text data is primarily sourced from the VideoInstruct100K dataset [57]. While each video in the dataset includes multiple instructions, the instructions tend to be relatively homogeneous, often focusing on simple video descriptions. To enhance the diversity of video-based tasks, we applied semantic deduplication to the instructions for each video and translated the dataset into Chinese, enriching the variety of video-based tasks for the model.

Audio Understanding Data. Most of the audio data is generated using TTS⁴, with prompts derived from text-only, image-text, and video-text datasets. To ensure the quality of the synthesized audio, we transcribed the generated audio using an ASR model and compared the transcriptions with the original prompts. Only those audio samples with accurate transcriptions are retained as final audio prompts. To further enrich the audio data, we included human-recorded audio samples that captured various dialects, accents, and background noises.

In addition to the general QA tasks, we also constructed a specific ASR dataset sourced from open-source data and internal logs. To improve training efficiency, we filtered out easily recognizable samples, focusing instead on more challenging audio data for supervised fine-tuning.

4 Experiment

4.1 Language Performance

4.1.1 Evaluation Benchmarks

We perform evaluations on 4 comprehensive benchmarks, including MMLU [31], CMMLU [42], AGIEval [107] and C-Eval [33]. MMLU includes 57 unique tasks consisting of multiple-choice questions across various fields of knowledge, including the humanities, social sciences, and hard sciences. CMMLU represents an extensive assessment framework tailored to assess the sophisticated knowledge and reasoning capabilities of LLMs within the context of Chinese language and culture. AGIEval is a human-centric benchmark crafted to evaluate the foundational models’ general cognitive and problem-solving abilities, based on official, public, and qualification tests designed for human participants. C-EVAL provides the comprehensive Chinese evaluation suite designed to evaluate the advanced knowledge and reasoning skills of LLMs within a Chinese context, encompassing 13,948 multiple-choice questions across 52 varied disciplines, from humanities to science and engineering. We conduct all the evaluations with a zero-shot measurement.

4.1.2 Major Performance

We compare Baichuan-omni with state-of-the-art proprietary multimodal models such as Gemini 1.5 Pro [69], and GPT-4o [63], as well as a series of competitive open-source LLMs and MLLMs such as VITA [24], MAP-Neo [101], Qwen1.5-Chat [5], Llama3-Instruct [3] and OLMo [29]. We list major results on comprehensive benchmarks in Table 1.

As shown in Table 1, our Baichuan-omni significantly outperforms open-source, general pure-text LLMs in comprehensive benchmarks. Compared to the open-source multimodal model VITA, Baichuan-omni demonstrates a substantial advantage in Chinese benchmarks, such as CMMLU (72.2% v.s 46.6%) and C-Eval (68.9% v.s 56.7%), and slightly surpasses VITA in AGIEval (47.7% v.s 46.2%).

4.2 Image Understanding

4.2.1 Evaluation Benchmarks

We evaluate Baichuan-omni on 13 representative vision-language benchmarks, including MMBench-EN, MMBench-CN [52], M3GIA [74], SEEDBench [41], RealWorldQA [87], MMMU [96], MathVista [56], MME [22], MMVet [94], TextVQA [72], OCRBench [53], ChartQA [60], and HallusionBench [30]. To ensure reproducible evaluation results, we use VLMEvalKit [20] uniformly for all evaluations. All evaluations are conducted in a zero-shot manner, adhering to the original setup of the models to ensure fair and consistent comparisons across all models and benchmarks.

⁴TTS tool: <https://github.com/2noise/ChatTTS>

Table 1: **Major results on comprehensive benchmarks.** *: Officially reported results. \diamond : Retrieved results from official leaderboard or recent papers. The rest unlabeled results are reproduced by ourselves.

Model	Comprehensive Tasks			
	MMLU (Acc.)	CMMLU (Acc.)	AGIEval (Acc.)	C-Eval (Acc.)
<i>Proprietary Models</i>				
GPT 4o	88.0 \diamond	78.3 \diamond	62.3 \diamond	86.0 \diamond
<i>Open-source Models (Pure text)</i>				
MAP-Neo (7B)	58.2	55.1	33.9	57.5
Qwen1.5-Chat (7B)	61.5	68.0	39.3	68.8
Llama3-Instruct (8B)	67.1	51.7	38.4	50.7
OLMo (7B)	28.4	25.6	19.9	27.3
<i>Open-source Models (Omni-modal)</i>				
VITA (8x7B)	71.0*	46.6	46.2*	56.7*
Baichuan-omni (7B)	65.3	72.2	47.7	68.9

4.2.2 Major Performance

We compare Baichuan-omni with state-of-the-art proprietary multimodal models such as Gemini 1.5 Pro [69], and GPT-4o [63], as well as a series of competitive open-source multimodal models such as VITA [24] and Qwen2-VL [79]. We list major results on VQA (Visual Question Answering) benchmarks and results on MCQ (Multi-choice & Yes-or-No Question) benchmarks in Table 2 and Table 3.

Table 2: **Major Results on Multi-choice benchmarks and Yes-or-No benchmarks.** *: Officially reported results. \diamond : Retrieved results from official leaderboard or recent papers. The rest unlabeled results are reproduced by ourselves.

Model	Multi-choice & Yes-or-No Question						
	MMBench (Acc.)	MMBench-CN (Acc.)	M3GIA (Acc.)	SEED-IMG (Acc.)	MME (Score)	MMMUM (val) (Acc.)	HallusionBench (Acc.)
<i>Proprietary Models</i>							
GPT-4o	83.4 \diamond	82.1 \diamond	59.8 \diamond	-	2328.7 \diamond	69.1 \diamond	55.0 \diamond
GPT-4o-mini	-	-	-	-	2003.4 \diamond	60.0 \diamond	46.1 \diamond
<i>Open-source Models (Vision-language)</i>							
Qwen2 VL (7B)	86.4	81.9	37.3	76.5	2326.8*	52.7	50.6*
MiniCPM-Llama3-V 2.5 (8B)	76.7	73.3	30.3	72.4	2024.6*	45.8*	42.5
<i>Open-source Models (Omni-modal)</i>							
VITA (8x7B)	74.7	71.4	27.7	72.6	2189.1	45.3	39.7*
Baichuan-omni (7B)	76.2	74.9	34.7	74.1	2186.9	47.3	47.8

As shown in Table 2 and Table 3, our Baichuan-omni comprehensively outperformed VITA-8*7b [24], which has 12B activated parameters, in multiple visual tasks, both on VQA benchmarks and MCQ benchmarks. Besides, we also demonstrates competitive performance comparable to, or even better than, open-source image-specialized models like MiniCPM-Llama3-V 2.5 [92]. Specifically, Baichuan-omni outperformed MiniCPM-Llama3-V 2.5 on most VQA tasks, including MMBench-CN, SEED-IMG, MME, HallusionBench and MMMUM which requires expert-level perception and reasoning. However, despite the advantage of incorporating an additional audio modality compared to Qwen2-VL [79], the performance gap between our model and Qwen2-VL in image tasks remains evident. Furthermore, it is worth noting that, beyond Qwen2-VL, the stark divide between open-source and closed-source models remains substantial.

Table 3: **Major Results on VQA benchmarks.** *: Officially reported results. \diamond : Retrieved results from official leaderboard or recent papers. The rest unlabeled results are reproduced by ourselves.

Model	Visual Question Answering					
	RealWorldQA (Acc.)	MMVet (Acc.)	MathVista-mini (Acc.)	TextVQA (val) (Acc.)	ChartQA (Acc.)	OCRBench (Acc.)
<i>Proprietary Models</i>						
GPT-4o	75.4 \diamond	69.1 \diamond	63.8 \diamond	-	85.7 \diamond	73.6 \diamond
GPT-4o-mini	67.1 \diamond	66.9 \diamond	52.4 \diamond	-	-	78.5 \diamond
<i>Open-source Models (Vision-language)</i>						
Qwen2 VL (7B)	69.7	62.0*	58.2*	84.3*	83.0*	84.5*
MiniCPM-Llama3-V 2.5 (8B)	63.5	52.0	54.3*	76.6	72.0	72.5
<i>Open-source Models (Omni-modal)</i>						
VITA (8x7B)	59.0	41.6*	44.9*	71.8	76.6	68.5*
Baichuan-omni (7B)	62.6	65.4	51.9	74.3	79.6	70.0

4.3 Video Understanding

4.3.1 Evaluation Benchmarks

We perform a thorough evaluation on general video understanding tasks (General VQA) and open-ended video question answering (Open-ended VQA) to comprehensively assess the video understanding capabilities of Baichuan-omni .

For general video understanding tasks, we select Perception-Test [65], MVBench [44], VideoMME [23], and EgoSchema [58] for long-form video-language understanding. We report top-1 accuracy for all benchmarks. For VideoMME, we report the results under the setting of "w/o subs". For open-ended video question answering part, we choose ActivityNet-QA [95] and MSVD-QA [88] as evaluation benchmarks. Following previous work [57], we utilize GPT to assess the quality of the response snippets. Specifically, we use GPT-3.5-Turbo to provide a "Yes-or-No" decision on the correctness of answers and a rating scaled from 0 to 5. We report the percentage of "Yes" responses as Accuracy and the average rating as Score.

We conduct all evaluations in a zero-shot way while avoiding the use of elaborate prompts. Besides, we follow the original setup of the models to be reproduced regarding the (maximum) number of frames, frame sampling rate, etc. they applied, ensuring fair and consistent comparisons across all models and benchmarks.

4.3.2 Major Performance

We compare Baichuan-omni with state-of-the-art multimodal proprietary models such as Gemini 1.5 Pro [69], GPT 4V [62], and GPT 4o [63], and a series of competitive open-source multimodal models such as VITA [24], Qwen2-VL [79], AnyGPT [98], VideoLLaMA 2 [13], VideoChat2 [44], LLaVA-NeXT-Video [104], and Video-LLaVA [48]. We list major results on general video understanding benchmarks in Table 4 and results on open-ended video question answering in Table 5.

Results on general video understanding benchmarks. As shown in Table 4, Baichuan-omni demonstrates competitive results over proprietary models on benchmarks like Egoschema and MVBench, and achieves strong performance across open-source multimodal models, which shows comprehensive video understanding capabilities of Baichuan-omni .

Compared to VITA, a MoE omni-modal LLM with about 12B activated parameters, Baichuan-omni (7B) outperforms it on all General Video QA benchmarks, and achieve an average improvement of about 4%. Additionally, Baichuan-omni excels a series of open-source models such as VideoLLaMA 2, VideoChat2, LLaVA-NeXT-Vide, and Video-LLaVA. Notably, Baichuan-omni also outperforms the proprietary model GPT 4V on MVBench (43.7%) and Egoschema (55.6%).

Results on open-ended video question answering benchmarks. The performance on Open-ended VQA is listed in Table 5. Baichuan-omni demonstrates SoTA performance (both Accuracy and Score) on ActivityNet-QA and MSVD-QA across all open-source models, such as the most recent SoTA multimodal models VITA and Qwen2 VL, and outperforms the proprietary model Gemini 1.5 Pro (56.7%) on ActivityNet-QA. The superior results indicate that Baichuan-omni is also effective in open-ended question answering, i.e., Baichuan-omni is more capable of generating informative and descriptive responses.

Table 4: **Major results on general VQA benchmarks: MVBench, Egoschema, VideoMME and Perception-Test.** max: Maximum number of sampling frames. *: Officially reported results. \diamond : Retrieved results from official leaderboard or recent papers. The rest unlabeled results are reproduced by ourselves.

Model	# Frames	General VQA			
		MVBench (Acc.)	Egoschema (Acc.)	VideoMME (Acc.)	Perception-Test (Acc.)
<i>Proprietary Models</i>					
Gemini 1.5 Pro	-	81.3 \diamond	63.2*	75.0 \diamond	-
GPT 4o	-	-	77.2*	71.9 \diamond	-
GPT 4V	-	43.7 \diamond	55.6*	59.9 \diamond	-
<i>Open-source Models (Vision-language)</i>					
Qwen2 VL (7B)	2 fps (max 768)	67.0* 64.4	66.7* 66.6	63.3* 59.0	62.3* 60.3
AnyGPT (8B)	48	33.2	32.1	29.8	29.1
VideoLLaMA 2 (7B)	16	54.6*	51.7*	46.6*	51.4*
VideoChat2 (7B)	16	51.1*	42.1 \diamond	33.7 \diamond	47.3 \diamond
LLaVA-NeXT-Video (7B)	32	46.5 \diamond	43.9 \diamond	33.7 \diamond	48.8 \diamond
Video-LLaVA (7B)	8	41.0 \diamond	38.4 \diamond	39.9 \diamond	44.3 \diamond
<i>Open-source Models (Omni-modal)</i>					
VITA (8x7B)	1 fps (max 32)	53.4	53.9	56.1	56.2
Baichuan-omni (7B)	1 fps (max 48)	60.9	58.8	58.2	56.8

Table 5: **Major results on ActivityNet-QA and MSVD-QA.** max: Maximum number of sampling frames. *: Officially reported results. The rest unlabeled results are reproduced by ourselves.

Model	# Frames	Open-ended VQA			
		ActivityNet-QA (Acc.)	MSVD-QA (Score)	ActivityNet-QA (Acc.)	MSVD-QA (Score)
<i>Proprietary Models</i>					
Gemini 1.5 Pro	-	56.7*	-	-	-
GPT 4o	-	61.9*	-	-	-
GPT 4V	-	59.5*	-	-	-
<i>Open-source Models (Vision-language)</i>					
Qwen2 VL (7B)	2 fps (max 768)	17.4	1.9	61.1	3.5
VideoLLaMA 2 (7B)	16	50.2*	3.3*	70.9*	3.8*
VideoChat2 (7B)	16	49.1*	3.3*	70.0*	3.9*
LLaVA-NeXT-Video (7B)	32	53.5*	3.2*	67.4	3.4
Video-LLaVA (7B)	8	45.3*	3.3*	70.7*	3.9*
<i>Open-source Models (Omni-modal)</i>					
VITA (8x7B)	1 fps (max 32)	55.0	3.5	63.9	3.7
Baichuan-omni (7B)	1 fps (max 48)	58.6	3.7	72.2	4.0

4.4 Audio Understanding

4.4.1 Evaluation Benchmarks

To validate the audio understanding capacity of Baichuan-omni, we present the evaluating results on benchmarks with three tasks:

- **Automatic Speech Recognition (ASR).** This is a fundamental task for audio-language model pre-training which directly transcribes the audio into the text. For ASR evaluation in the general scene, we report results on the Fleurs [16]

Chinese (*zh*) and English (*en*) test sets, as well as the WenetSpeech [99] *test_net* dataset. To assess performance in more challenging ASR scenarios, we include results from the WenetSpeech [99] *test_meeting* dataset and the KeSpeech [78] test set, which evaluate the model’s ASR capabilities in ‘Meeting’ and ‘Chinese dialect’ contexts. For WenetSpeech, we use both Word Error Rate (WER) and Character Error Rate (CER) as evaluation metrics, while for others, only WER is used.

- **Speech-to-Text Translation (S2TT).** The task aims to translate the audio signal in the source to the target language. We evaluate the model’s S2TT performance between Chinese and English using the zh2en and en2zh subsets of the Covost2 [82] dataset, with BLEU [64] scores as the evaluation metric.
- **AIR-Bench.** The goal of this benchmark is to evaluate the chat capabilities to follow instructions of the given audio. We evaluate chat performance on the chat benchmark [91] (test set), using Score as the metric.

4.4.2 Major Performance

Baichuan-omni is compared with the state-of-the-art baselines across ASR, S2TT and SER tasks, including the recent leading large audio-language model Qwen2-Audio-Instruct [14] and the large omni-modal language model VITA [24]. On top of that, the performance of the classical pre-trained audio language model, Whisper-large-v3 [67], is presented for ASR and the performance of SALMONN [77] is presented for S2TT.

Results on ASR benchmarks. Baichuan-omni exhibits a strong audio transcription capacity in Table 6. Baichuan-omni primarily targets the Chinese corpus. In the general Chinese ASR scene, Baichuan-omni has a 2.0% WER (2.6% CER) superiority on the Fleurs test-zh subset and 4.1% WER (4.2% CER) improvement on the WenetSpeech test_net when comparing with Qwen2-Audio-Instruct. The evaluation results on WenetSpeech consistently demonstrate the superiority of our model over VITA. Baichuan-omni achieves nearly a 50% improvement in the CER performance of VITA, both in test_net (7.1% v.s 12.2%) and test_meeting (8.9% v.s 16.5%) subsets. On the more challenging Chinese dialect benchmark, KeSpeech, our model maintains a comprehensive lead, with an average CER of 6.7% over the performance of all dialects. Notably, while our model excels in Chinese audio transcription, Baichuan-omni also maintains robust general ASR performance in English. We achieve 4.7% of WER, which exceeds Qwen2-Audio-Instruct by 11% WER.

Table 6: **Major results on Fleurs, WenetSpeech, and KeSpeech.** Test sets of WenetSpeech are evaluated with WER and CER, while other test sets are evaluated only with WER. VITA’s evaluation results are officially reported in their paper [24], marked with *. The rest unlabeled results are reproduced by ourselves, and any performance divergence may be attributed to differences in decoding parameters.

Scene	Dataset	Model	Results WER (CER) ↓
General	Fleurs <i>test-zh</i> <i>test-en</i>	Whisper-large-v3 (1.55B)	12.4 7.2
		Qwen2-Audio-Instruct (7B)	9.0 15.7
		Baichuan-omni (7B)	7.0 4.7
	WenetSpeech <i>test_net</i>	Whisper-large-v3 (1.55B)	17.5 (18.5)
		Qwen2-Audio-Instruct (7B)	11.0 (11.3)
VITA (8x7B)		-(12.2*)	
Baichuan-omni (7B)	6.9 (7.1)		
Meeting	WenetSpeech <i>test_meeting</i>	Whisper-large-v3 (1.55B)	30.8 (31.7)
		Qwen2-Audio-Instruct (7B)	10.7 (10.8)
		VITA (8x7B)	-(16.5*)
		Baichuan-omni (7B)	8.4 (8.9)
Chinese Dialect	KeSpeech <i>mandarin</i> <i>beijing</i> <i>southwest</i> <i>lan-yin</i> <i>zhongyuan</i> <i>northeast</i> <i>jiang-huai</i> <i>ji-lu</i> <i>jiao-liao</i>	Whisper-large-v3 (1.55B)	18.7 44.8 52.9 54.8 50.1 22.9 54.7 47.0 50.4
		Qwen2-Audio-Instruct (7B)	5.8 9.7 10.5 11.0 8.2 8.4 13.8 10.3 11.2
		Baichuan-omni (7B)	2.8 6.4 7.0 7.7 6.1 5.8 9.0 8.3 7.2

Results on S2TT and AIR-Bench benchmarks. In addition to ASR, Baichuan-omni excels in both S2TT and SER tasks. The evaluation results are summarized in Table 7. Notably, when translating from English to Chinese on the Covost-2 en2zh test set, Baichuan-omni outperforms Qwen2-Audio-Instruct by approximately 7 BLEU points. For the reverse translation, from Chinese to English, our model’s performance on the Covost-2 zh2en test set is comparable to that of Qwen2-Audio-Instruct. On the AirBench, Baichuan-omni scores 7.42 and 7.26 for speech and sound, respectively, outperforming Qwen2-Audio-Instruct and showcasing Baichuan-omni’s superior ability to generate realistic human speeches and sounds.

Table 7: **Major results on Covost2 and AirBench.** \diamond represents the results from the official leaderboard or recent papers. The rest unlabeled results are reproduced by ourselves, and any performance divergence may be attributed to differences in decoding parameters.

Task	Dataset	Model	Metrics	Results
S2TT	Covost-2 <i>zh2en en2zh</i>	SALMONN (7B)	BLEU \downarrow	- 33.1 \diamond
		Qwen2-Audio-Instruct (7B)		23.3 34.1
		Baichuan-omni (7B)		22.1 40.2
AIR-Bench	Chat Benchmark <i>speech sound music mix-audio</i>	Qwen2-Audio-Instruct (7B)	Score \uparrow	7.18 6.99 6.79 6.77
		VITA (8x7B)		6.40 6.59 6.59 5.94
		Baichuan-omni (7B)		7.42 7.26 6.12 5.76

4.5 Ablation Study

4.5.1 Image-Language Branch

Visual encoder. To compare the performance of different visual encoders in Baichuan-omni, we conducted experiments across various vision encoders with differing parameter sizes, input resolutions, and output token counts. We selected five mainstream vision encoders: OpenAI’s CLIP series [66], Google’s Siglip series [97], Apple’s DFN series [21], OpenGVLab’s InternViT series [11], and BAAI’s EVA series [76], totaling 14 models. All models are trained with contrastive learning, with parameters ranging from 300M (ViT-L) to 18B. The training data used during the pre-training of the visual encoders varied from 400M to 10B, with input resolutions spanning from 224 \times 224 to 448 \times 448 and output token counts from 256 to 1024. All comparative experiments are conducted under the same experimental conditions, specifically using a batch size of 8 and the same data for IFT training (with a data ratio of *Caption: Interleaved data: Pure text* set at 0.45: 0.45: 0.1).

Table 8: **Comparative study across various vision encoders with differing parameter sizes and input resolutions.** We evaluated the model on 10 benchmarks, including SEEDBench2 [41], TextCaps (val) [71], TextVQA (val) [72], OCRBench [53], OCRBench (CN), OKVQA [59], Nocaps [2], VQAv2 [28], DocVQA (val) [61], and GQA [34]. We compiled the Average Performance of the model across these 10 benchmarks. Additionally, based on the specific tasks targeted by these benchmarks or their subcategories, we calculated the average scores of the model in six areas: OCR, Nature Image Understanding (NIU), Spatial, Chart, Common Sense Knowledge, and Video.

Model	Params.	Resolution	OCR	NIU	Spatial	Chart	Common Sense	Video	Avg.
siglip-so400m-patch14-384	428M	384 px	44.67	56.91	41.70	25.00	51.57	40.63	43.80
clip-vit-large-patch14-336	304M	336 px	30.19	55.60	41.21	15.63	48.05	37.50	39.51
dfn5b-clip-ViT-H-14-378	631M	378 px	29.05	54.75	37.50	21.88	49.22	34.38	39.14
InternViT-6B-224px	5.9B	224 px	14.17	40.60	29.98	15.63	40.63	34.38	29.99
InternViT-6B-448px-v1-5	5.5B	448 px	17.49	46.97	35.06	18.75	41.02	31.25	32.80
eva-clip-8b	7.5B	224 px	28.86	56.61	40.92	25.00	49.22	40.63	41.51
eva-clip-8b-448	7.5B	448 px	<u>32.66</u>	58.09	43.26	21.88	<u>49.61</u>	37.50	<u>41.86</u>

As shown in Table 8, while increasing the resolution does lead to performance improvements (eva-448 vs eva-224, InternViT-6B-224px vs InternViT-6B-448px), the number of encoder parameters does not exhibit a direct relationship with the metrics. Overall, siglip-so400m-patch14-384 [97] achieved the highest average score and excelled in four out

of the six tasks, particularly demonstrating outstanding performance in OCR. Considering these factors along with efficiency issues, we ultimately selected siglip-so400m-patch14-384 as the visual encoder for our Baichuan-omni .

we further explored the impact of using AnyRes [50] on the model’s visual-language performance. We found that using AnyRes [50] results in a significant performance improvement compared to a fixed input of 384 pixels, particularly in tasks that rely on image details, such as visual document understanding, as shown in Table 9.

Table 9: **Comparative study on AnyRes.** Using AnyRes results in a significant improvement in visual document understanding and OCR.

Method	TextVQA (val)	DocVQA (val)	InfographicVQA	OCRBench
Baseline	66.48	72.61	47.54	76.92
Baseline+Anyres	69.13	87.48	62.80	78.44

Projector. Regarding the projector, we compared the following approaches: **(1) MLP:** Directly passes through a two-layer MLP, aligning the dimensions to that of the LLM, without reducing the number of image tokens. **(2) C-abs:** Passes through two convolutional layers and one pooling layer, aligning the dimensions to that of the LLM while reducing the number of tokens as needed (e.g., from 576 to 144). **(3) Concat:** Concatenates adjacent tokens and then processes them through a MLP, allowing for token reduction but increasing the number of parameters (as the MLP’s input dimension increases). **(4) Mean Pool:** Applies a convolutional layer with a stride of 2 for pooling before passing the tokens through a two-layer MLP, enabling token reduction while maintaining a consistent parameter count with the MLP.

In early experiments, we found that models trained with different projectors exhibited little difference in general image understanding, but shows disparities in Chinese OCR comprehension after adding 1 million pure Chinese OCR VQA data. The results showed that while the model with the C-abs projector struggled to learn Chinese OCR capabilities, the model with the MLP projector began fitting the data and demonstrated zero-shot capability after 0.75 epochs. Ultimately, we ranked the projectors as follows: MLP > Mean Pool > Concat > C-abs. On the other hand, to minimize the number of tokens per sub-image after the AnyRes operation (where MLP produces 729 tokens, while Mean Pool, Concat, and C-abs each produce 182 tokens), we chose Mean Pool as our visual projector.

4.5.2 Video-Language Branch

For video modality, we conduct an ablation study from three perspectives to thoroughly investigate the impact of various factors on model performance.

Number of frames. Within the constraints of the context length, we systematically adjust the frame sampling rate to control the maximum number of input video frames.

Resolution of vision encoder. We explore the effect of different vision encoder resolutions on the model’s ability to extract meaningful visual features. Our investigation spans from fixed resolutions (such as 384×384 pixels) to dynamic resolution approaches like AnyRes.

Video-language pre-training. We evaluate the model’s performance both with and without video-language pre-training. This comparison helps us quantify the benefits of leveraging large-scale multimodal datasets for pre-training, potentially enhancing the model’s ability to understand video-text relationships and generalize across various video-understanding tasks.

Table 10: **Ablation study on Video-language branch.** We analyze the influence of number of frames, resolution of vision encoder, and video-language pre-training.

w/ Pre-training	Resolution	# Frames	# Tokens	MVBench	VideoMME	ActivityNet-QA	Avg.
✓	384 px	64	45	50.5	56.9	56.8	54.7
✓	384 px	48	45	47.6	54.6	48.1	50.1
×	384 px	64	45	46.8	56.0	44.6	49.1
✓	AnyRes	48	182 - 546	60.9	58.2	58.6	59.2

As demonstrated in Table 10, the model’s performance in video comprehension is significantly impacted by the number of input frames processed. When the quantity of input frames is decreased from 64 to 48, there is a notable decline (from 54.7% to 50.1% average) in the model’s ability to understand and interpret video content.

When testing the model by inputting a total of 48 frames, it has been observed that the version of the model that utilizes AnyRes technology demonstrates superior performance compared to the version that operates with a fixed resolution set at 384×384 . This performance advantage is evident across various benchmarks, including MVBench, VideoMME, and ActivityNet-QA. In fact, the model with AnyRes enabled shows an average improvement of around 5% over its fixed-resolution counterpart.

Besides, from the first and third rows of the table, it shows that incorporating video-text pre-training has a pronounced impact on the model’s video understanding capability. For instance, in MVBench, the model without pre-training lags approximately 6% behind the one with pre-training.

Overall, we found that increasing the number of video frames, enhancing the resolution of the visual encoder, and incorporating video-text data during the pre-training stage all contribute to improving the model’s capabilities of understanding videos. We will leave the exploration of these three factors in situations where the input exceeds the context length (increasing number of frames and resolution) for future work.

4.5.3 Audio-Language Branch

The audio projector in the audio-language branch plays a key role in bridging the representations of audio and natural language modalities. Notably, our newly introduced projector with Conv-GMLP demonstrates the great performance robustness of the feature down-sampling rate.

For analysis, we measure the average WER on all our ASR benchmarks across Fleurs, WenetSpeech and KeSpeech by training a 1.5B audio-language model with three different down-sampling rates 2, 4, 8. To simulate the actual training of the audio branch in Baichuan-omni, we only train our audio encoder and projector while keeping the LLM frozen. This setup is consistent with the configuration described in Section 3.2.3.

From the Figure 6, we observe that when the down-sampling rate is set to 2, the audio-language model achieves the best ASR performance, with an average WER of 7.7%. When the down-sampling rate is adjusted to 4 and 8, there is a slight degradation in ASR performance, but the decrease is minimal (ranging from 0.3% to 0.6%). Surprisingly, despite the greater degree of down-sampling, the model with a rate of 8 outperforms the one with a rate of 4 (8.0% vs. 8.3%). This highlights the exceptional sequence compression capability of the Conv-GMLP.

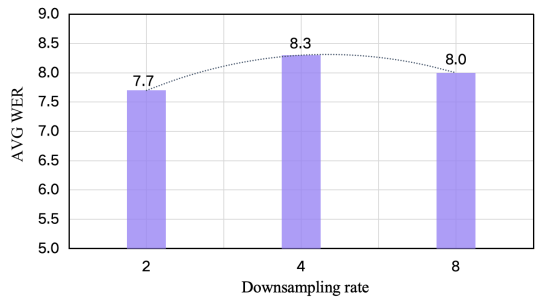


Figure 6: **Ablation study on down-sampling rate.** Average WER across multiple test sets (Fleurs zh/en, WenetSpeech net/meeting, and KeSpeech) using various down-sampling rates of Conv-GMLP.

4.5.4 Multimodal Supervised Fine-Tuning

Table 11 and Table 12 compare the performance of Baichuan-omni on various image and video benchmarks with and without multimodal supervised fine-tuning (SFT). The results indicate that the model exhibits superior overall performance after undergoing multimodal SFT compared to the version that only undergoes instruction fine-tuning (IFT). This improvement can be attributed to the use of high-quality, diverse instructions and our SFT data construction method, which avoid compromising the base model’s capabilities. (See Section 3.3 for more details.)

Table 11: **Performances of Baichuan-omni on image tasks before and after the supervised fine-tuning stage.** Generally, the model’s performance has improved across most image benchmarks.

Method	Multi-choice Question				Visual Question Answering			
	MMBench (Acc.)	MMBench-CN (Acc.)	MMMU (Acc.)	SEED-IMG (Acc.)	ChartQA (Acc.)	MathVista (Acc.)	MMVet (Acc.)	RealWorldQA (Acc.)
IFT	75.6	69.3	48.3	73.0	76.0	51.6	55.0	62.9
SFT	76.2	74.9	47.3	74.1	79.6	51.9	65.4	62.6

5 Conclusion

In this work, we have open-sourced Baichuan-omni as a step toward developing a truly omni-modal LLM that encompasses all human senses. With omni-modal pretraining and fine-tuning using high-quality omni-modal data,

Table 12: **Performances on video tasks before and after the supervised fine-tuning stage.** Multimodal supervised fine-tuning brings significant improvements for the vast majority of video benchmarks.

Method	General VQA				Open-ended VQA			
	Egoschema (Acc.)	MVBench (Acc.)	VideoMME (Acc.)	Perception (Acc.)	ActivityNet-QA (Acc.)	(Score)	MSVD-QA (Acc.)	(Score)
IFT	54.0	61.3	56.3	56.9	55.4	3.6	66.6	3.8
SFT	58.8	60.9	58.2	56.8	58.6	3.7	72.2	4.0

this version of Baichuan-omni has achieved leading levels in integrating comprehension across video, image, text, and audio. Despite its promising performance, there remains significant room for improvement in the foundational capabilities across each individual modality. This include (1) enhancing text extraction capabilities, (2) supporting longer video understanding, (3) developing an end-to-end TTS system integrated with LLMs, and (4) improving the ability to comprehend not only human voices but also natural environmental sounds, such as flowing water, bird calls, and collision noises, among others.

We anticipate efforts from both academia and industry in the field, and we believe that the expansion of model modalities, along with simultaneous advancements in enhancing model capabilities, will ultimately bring the dream of Artificial General Intelligence closer to reality.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. Nocaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8948–8957, 2019.
- [3] AI@Meta. Llama 3 model card, 2024.
- [4] Anas Awadalla, Le Xue, Oscar Lo, Manli Shu, Hannah Lee, Etash Kumar Guha, Matt Jordan, Sheng Shen, Mohamed Awadalla, Silvio Savarese, Caiming Xiong, Ran Xu, Yejin Choi, and Ludwig Schmidt. Mint-1t: Scaling open-source multimodal data by 10x: A multimodal dataset with one trillion tokens. *arXiv preprint arXiv:2406.11271*, 2024.
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [6] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [7] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1728–1738, 2021.
- [8] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [9] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*, 2024.
- [10] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Bin Lin, Zhenyu Tang, Li Yuan, Yu Qiao, Dahua Lin, Feng Zhao, and Jiaqi Wang. Sharegpt4video: Improving video understanding and generation with better captions, 2024.
- [11] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- [12] 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.

- [13] 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.
- [14] 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.
- [15] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- [16] Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. Fleurs: Few-shot learning evaluation of universal representations of speech. *arXiv preprint arXiv:2205.12446*, 2022.
- [17] Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. Pengi: An audio language model for audio tasks. *Advances in Neural Information Processing Systems*, 36:18090–18108, 2023.
- [18] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Guosheng Dong, Da Pan, Yiding Sun, Shusen Zhang, Zheng Liang, Xin Wu, Yanjun Shen, Fan Yang, Haoze Sun, Tianpeng Li, et al. Baichuanseed: Sharing the potential of extensive data collection and deduplication by introducing a competitive large language model baseline. *arXiv preprint arXiv:2408.15079*, 2024.
- [20] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models, 2024.
- [21] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023.
- [22] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2024.
- [23] 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.
- [24] 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.
- [25] Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie Yan. Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition. In *INTERSPEECH*, 2022.
- [26] Roopal Garg, Andrea Burns, Burcu Karagol Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa Onoe, Andrew Bunner, Ranjay Krishna, Jason Baldrige, and Radu Soricut. Imageinwords: Unlocking hyper-detailed image descriptions. *arXiv preprint arXiv:2405.02793*, 2024.
- [27] Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. Does fine-tuning llms on new knowledge encourage hallucinations?, 2024.
- [28] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [29] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*, 2024.
- [30] 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.
- [31] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.

- [32] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024.
- [33] Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [34] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [35] Anil K Jain. Image data compression: A review. *Proceedings of the IEEE*, 69(3):349–389, 1981.
- [36] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [37] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pages 498–517. Springer, 2022.
- [38] Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv preprint arXiv:2402.01831*, 2024.
- [39] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text documents, 2023.
- [40] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024.
- [41] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023.
- [42] Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. *arXiv preprint arXiv:2306.09212*, 2023.
- [43] 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 *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [44] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024.
- [45] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal ArXiv: A dataset for improving scientific comprehension of large vision-language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14369–14387, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [46] 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.
- [47] Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024.
- [48] 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.
- [49] Hanxiao Liu, Zihang Dai, David R. So, and Quoc V. Le. Pay attention to mlps. *arXiv preprint arXiv:2105.08050*, 2021.
- [50] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024.

- [51] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [52] 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.
- [53] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, 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*, 2024.
- [54] 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.
- [55] Keer Lu, Zheng Liang, Xiaonan Nie, Da Pan, Shusen Zhang, Keshi Zhao, Weipeng Chen, Zenan Zhou, Guosheng Dong, Wentao Zhang, et al. Datasculpt: Crafting data landscapes for llm post-training through multi-objective partitioning. *arXiv preprint arXiv:2409.00997*, 2024.
- [56] 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.
- [57] 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.
- [58] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023.
- [59] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019.
- [60] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- [61] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021.
- [62] OpenAI. GPT-4V(ision) system card. <https://openai.com/index/gpt-4v-system-card/>, 2023.
- [63] OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024.
- [64] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.
- [65] Viorica Pătrăucean, Lucas Smaira, Ankush Gupta, Adrià Recasens Contente, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, et al. Perception test: A diagnostic benchmark for multimodal video models. *arXiv preprint arXiv:2305.13786*, 2023.
- [66] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [67] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022.
- [68] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [69] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [70] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models, 2022.

- [71] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 742–758. Springer, 2020.
- [72] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019.
- [73] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009, 2023.
- [74] Wei Song, Yadong Li, Jianhua Xu, Guowei Wu, Lingfeng Ming, Kexin Yi, Weihua Luo, Houyi Li, Yi Du, Fangda Guo, et al. M3gia: A cognition inspired multilingual and multimodal general intelligence ability benchmark. *arXiv preprint arXiv:2406.05343*, 2024.
- [75] Tongyi SpeechTeam. Funaudiollm: Voice understanding and generation foundation models for natural interaction between humans and llms. *arXiv preprint arXiv:2407.04051*, 2024.
- [76] Quan Sun, Jinsheng Wang, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, and Xinlong Wang. Eva-clip-18b: Scaling clip to 18 billion parameters. *arXiv preprint arXiv:2402.04252*, 2023.
- [77] 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. In *The Twelfth International Conference on Learning Representations*, 2024.
- [78] Zhiyuan Tang, Dong Wang, Yanguang Xu, Jianwei Sun, Xiaoning Lei, Shuaijiang Zhao, Cheng Wen, Xingjun Tan, Chuandong Xie, Shuran Zhou, Rui Yan, Chenjia Lv, Yang Han, Wei Zou, and Xiangang Li. 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.
- [79] Qwen Team. Qwen2-VL: To See the World More Clearly. *Qwen*, August 2024.
- [80] 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 preprint arXiv:2302.13971*, 2023.
- [81] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutvi Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [82] Changhan Wang, Anne Wu, and Juan Pino. Covost 2 and massively multilingual speech-to-text translation, 2020.
- [83] Junjie Wang, Yin Zhang, Yatai Ji, Yuxiang Zhang, Chunyang Jiang, Yubo Wang, Kang Zhu, Zekun Wang, Tiezhen Wang, Wenhao Huang, Jie Fu, Bei Chen, Qunshu Lin, Minghao Liu, Ge Zhang, and Wenhui Chen. Pin: A knowledge-intensive dataset for paired and interleaved multimodal documents, 2024.
- [84] Wenhui Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36, 2024.
- [85] Yonghui Wang, Wengang Zhou, Hao Feng, Keyi Zhou, and Houqiang Li. Towards improving document understanding: An exploration on text-grounding via mllms. *arXiv preprint arXiv:2311.13194*, 2023.
- [86] Haibin Wu, Xuanjun Chen, Yi-Cheng Lin, Kai-wei Chang, Ho-Lam Chung, Alexander H Liu, and Hung-yi Lee. Towards audio language modeling—an overview. *arXiv preprint arXiv:2402.13236*, 2024.
- [87] x.ai. Grok-1.5 vision preview. <https://x.ai/blog/grok-1.5v>, 2024.
- [88] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653, 2017.
- [89] Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023.
- [90] 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.
- [91] Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv, Zhou Zhao, Chang Zhou, and Jingren Zhou. AIR-bench: Benchmarking large audio-language models via generative comprehension. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual*

- Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1979–1998, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [92] 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.
- [93] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- [94] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- [95] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134, 2019.
- [96] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- [97] 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.
- [98] 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. *arXiv preprint arXiv:2402.12226*, 2024.
- [99] Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, Di Wu, and Zhendong Peng. Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition, 2022.
- [100] 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, 2023.
- [101] Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang, Chenghua Lin, Chou Leuang Yu, Danny Pan, Esther Cheng, Jie Liu, Qunshu Lin, et al. Map-neo: Highly capable and transparent bilingual large language model series. *arXiv preprint arXiv:2405.19327*, 2024.
- [102] 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.
- [103] Wenqi Zhang, Zhenglin Cheng, Yuanyu He, Mengna Wang, Yongliang Shen, Zeqi Tan, Guiyang Hou, Mingqian He, Yanna Ma, Weiming Lu, et al. Multimodal self-instruct: Synthetic abstract image and visual reasoning instruction using language model. *arXiv preprint arXiv:2407.07053*, 2024.
- [104] 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, April 2024.
- [105] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [106] Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She, Zheng Lin, Wenbin Jiang, and Weiping Wang. Multimodal table understanding. *arXiv preprint arXiv:2406.08100*, 2024.
- [107] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364*, 2023.
- [108] Minxuan Zhou, Hao Liang, Tianpeng Li, Zhiyu Wu, Mingan Lin, Linzhuang Sun, Yaqi Zhou, Yan Zhang, Xiaoqin Huang, Yicong Chen, Yujing Qiao, Weipeng Chen, Bin Cui, Wentao Zhang, and Zenan Zhou. Mathscape: Evaluating mllms in multimodal math scenarios through a hierarchical benchmark, 2024.
- [109] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.