



MIO: A Foundation Model on Multimodal Tokens

Zekun Wang^{1,3,13}, King Zhu^{2,3}, Chunpu Xu⁴, Wangchunshu Zhou⁵, Jiaheng Liu^{3,12},
Yibo Zhang¹, Jiashuo Wang⁴, Ning Shi⁶, Siyu Li³, Yizhi Li^{3,8}, Haoran Que¹,
Zhaoxiang Zhang⁹, Yuanxing Zhang^{10,13}, Ge Zhang^{3,7}, Ke Xu¹,
Jie Fu^{11*}, Wenhao Huang^{2,3*}

¹Beihang University; ²01.AI; ³M-A-P ⁴The Hong Kong Polytechnic University; ⁵AIWaves;

⁶University of Alberta; ⁷University of Waterloo; ⁸University of Manchester; ⁹ Chinese Academy of Sciences

¹⁰Peking University; ¹¹Shanghai AI Lab; ¹²Nanjing University; ¹³Kuaishou Technology;

zenmoore@buaa.edu.cn

Abstract

In this paper, we introduce MIO, a novel foundation model built on multimodal tokens, capable of understanding and generating speech, text, images, and videos in an end-to-end, autoregressive manner. While the emergence of large language models (LLMs) and multimodal large language models (MM-LLMs) propels advancements in artificial general intelligence through their versatile capabilities, they still lack true any-to-any understanding and generation. Recently, the release of GPT-4o has showcased the remarkable potential of any-to-any LLMs for complex real-world tasks, enabling omnidirectional input and output across images, speech, and text. However, it is closed-source and does not support the generation of multimodal interleaved sequences. To address this gap, we present MIO, which is trained on a mixture of discrete tokens across four modalities using causal multimodal modeling. MIO undergoes a four-stage training process: (1) alignment pre-training, (2) interleaved pre-training, (3) speech-enhanced pre-training, and (4) comprehensive supervised fine-tuning on diverse textual, visual, and speech tasks. Our experimental results indicate that MIO exhibits competitive, and in some cases superior, performance compared to previous dual-modal baselines, any-to-any model baselines, and even modality-specific baselines. Moreover, MIO demonstrates advanced capabilities inherent to its any-to-any feature, such as interleaved video-text generation, chain-of-visual-thought reasoning, visual guideline generation, instructional image editing, etc.

1 Introduction

The advent of Large Language Models (LLMs) is commonly considered the dawn of artificial general intelligence (AGI) (OpenAI et al., 2023; Bubeck et al., 2023), given their generalist capabilities such as complex reasoning (Wei et al., 2022), role

playing (Wang et al., 2023b), and creative writing (Wang et al., 2024a). However, original LLMs lack multimodal understanding capabilities. Consequently, numerous multimodal LLMs (MM-LLMs) have been proposed, allowing LLMs to understand images (Li et al., 2023b; Alayrac et al., 2022), audio (Borsos et al., 2023; Rubenstein et al., 2023; Tang et al., 2023; Das et al., 2024), and other modalities (Lyu et al., 2023; Zhang et al., 2023d; Moon et al., 2023). These MM-LLMs typically involve an external multimodal encoder, such as EVA-CLIP (Sun et al., 2023b) or CLAP (Elizalde et al., 2022), with an alignment module such as Q-Former (Li et al., 2023b) or MLP (Liu et al., 2023) for multimodal understanding. These modules align non-textual-modality data features into the embedding space of the LLM backbone.

Another line of work involves building **any-to-any** and end-to-end MM-LLMs that can input and output non-textual modality data. Typically, there are four approaches: (1) Discrete-In-Discrete-Out (DIDO): Non-textual modality data is discretized using vector quantization techniques (van den Oord et al., 2017; Esser et al., 2020) and then fed into LLMs (Ge et al., 2023b; Zhan et al., 2024; Liu et al., 2024). (2) Continuous-In-Discrete-Out (CIDO): The LLM backbones intake densely encoded non-textual modality data features and generate their quantized representations (Diao et al., 2023; Team et al., 2023). (3) Continuous-In-Continuous-Out (CICO): The LLMs both understand and generate non-textual modality data in their densely encoded representations (Sun et al., 2023c,a; Dong et al., 2023; Zheng et al., 2023; Wu et al., 2023). (4) Autoregression + Diffusion (AR + Diff): The autoregressive and diffusion modeling are integrated in a unified LLM (Zhou et al., 2024; Xie et al., 2024; Li et al., 2024b). Although these works have succeeded in building MM-LLMs unifying understanding and generation, they exhibit some drawbacks, as illustrated in Table 1. For example, Emu1 (Sun

*Corresponding Authors.

| Models | Emu2 (Sun et al., 2023a) | SEED- LLaMA (Ge et al., 2023b) | AnyGPT (Zhan et al., 2024) | CM3Leon (Yu et al., 2023), Chameleon (Team, 2024) | Gemini (Reid et al., 2024) | Transfusion (Zhou et al., 2024) | MIO (ours) |
|-------------------|--------------------------------|---|----------------------------------|--|----------------------------------|---------------------------------------|---------------|
| I/O Consistency | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ |
| Uni. Bi. SFT | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Multi-Task SFT | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Speech I/O | ✗/✗ | ✗/✗ | ✓/✓ | ✗/✗ | ✓/✗ | ✗ | ✓/✓ |
| Video I/O | ✓/✓ | ✓/✓ | ✗/✗ | ✗/✗ | ✓/✗ | ✗ | ✓/✓ |
| Voice Output | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| MM. Inter. Output | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Modeling | CICO | DIDO | DIDO | DIDO | CIDO | AR+Diff | DIDO |

Table 1: The comparison between previous models and MIO (ours). **I/O Consistency** indicates whether the model ensures that the input and output representations for the same data remain consistent. **Uni. Bi. SFT** refers to whether the model undergoes a unified (Uni.) supervised fine-tuning (SFT) for both multimodal understanding and generation (Bi.=Bidirectional). **Multi-Task SFT** assesses whether the model undergoes a comprehensive SFT that includes diverse tasks, with at least visual question answering tasks. **MM. Inter. Output** evaluates whether the model supports the generation of multimodal interleaved (MM. Inter.) sequences. We refer readers to §1 for the definitions of the different modeling approaches.

et al., 2023c) and Emu2 (Sun et al., 2023a) explore the autoregressive modeling of three modalities: text, images, and videos. SEED-LLaMA (Ge et al., 2023b) proposes a new image quantizer aligned with LLMs’ embedding space and trains the MM-LLMs on images and videos. However, neither considers the speech modality, which is heterogeneous from visual modalities like videos and images. Although AnyGPT (Zhan et al., 2024) has explored settings involving four modalities, including text, image, speech, and music, it lacks video-related abilities, voice synthesis, and comprehensive multi-task supervised fine-tuning, leading to limited multimodal instruction-following and reasoning capabilities. Furthermore, AR + Diff approaches, such as Transfusion (Zhou et al., 2024), suffer from limited multimodal understanding capabilities because the multimodal inputs are noised for denoising modeling, and the image tokenizer used (VAE (Kingma and Welling, 2013)) is suitable for image generation rather than understanding.

Moreover, most of current MM-LLMs are dual-modal, combining text with another modality, such as images. Although previous works, such as Meta-Transformer (Zhang et al., 2023d) and Unified-IO 2 (Lu et al., 2023), have explored omni-multimodal understanding settings with more than two non-textual modalities, they lag significantly behind their dual-modal counterparts, especially in terms of multimodal instruction-following capabilities. Moreover, these MM-LLMs are typically focused on understanding only, neglecting the important aspect of multimodal generation. Several works

have enabled LLMs to call external tools to address this issue. For example, HuggingGPT (Shen et al., 2023) generates textual image descriptions for external diffusion models to synthesize images. GPT-4 (OpenAI et al., 2023) can utilize either an image generator like DALL-E 3 (Betker et al., 2024) or a text-to-speech (TTS) tool like Whisper (Radford et al., 2022) to support multimodal generation.¹ However, these methods are not end-to-end, relying on the text modality as an interface.

Recently, the release of GPT-4o has demonstrated the capabilities of any-to-any and end-to-end foundation models.² It is the first foundational model to accept multimodal tokens as inputs and generate multimodal tokens in a unified model while also demonstrating strong abilities in complex multimodal instruction-following, reasoning, planning, and other generalist capabilities. Furthermore, as the scaling up of LLMs in the community depletes high-quality language tokens, GPT-4o verifies a new source of data for LLM training: multimodal tokens. This approach suggests that the next generation AGI could derive more knowledge from multimodal tokens when language tokens are exhausted. However, GPT-4o is closed source and focuses on end-to-end support for speech I/O, image I/O, 3D generation, and video understanding. Its open-source “alternatives”, such as VITA (Fu et al., 2024), still lack the ability to *generate* data of all supported modalities, particularly for the generation of multimodal interleaved sequences.

¹<https://openai.com/index/chatgpt-can-now-see-hear-and-speak/>

²<https://openai.com/index/hello-gpt-4o/>

To address the aforementioned issues, we introduce MIO (Multimodal Input and Output, or Multimodal Interleaved Output), the first open-source any-to-any foundation model that unifies multimodal understanding and generation across four modalities—text, image, speech (with voice), and video, while enabling the generation of multimodal interleaved sequences. Specifically, MIO is built on discrete multimodal tokens that capture both semantic representations through contrastive loss and low-level features via reconstruction loss (Ge et al., 2023a; Zhang et al., 2023b) from raw multimodal data. Due to the consistent data format shared with textual corpora, the model can treat non-textual modalities as “foreign languages”, allowing it to be trained with next-token prediction. Note that since the representation of an image remains the same whether it is used as an input or an output, our model supports multimodal interleaved sequence generation, where an image functions for both understanding and generation. Moreover, we employ three-stage pre-training with an additional SFT stage to effectively train the model for modality scaling.

Our experimental results show that MIO, trained on a mixture of four modalities, demonstrates competitive performance compared to its dual-modal counterparts and previous any-to-any multimodal language model baselines. Additionally, MIO is the first model to demonstrate interleaved video-text generation, chain-of-visual-thought reasoning, and other emergent abilities relying on any-to-any and multimodal interleaved output features (*c.f.*, §E.5).

2 Method

Firstly, we elaborate on our modeling approach, which supports multimodal token input and output, as well as causal language modeling (CausalLM), in §2.1. Secondly, we describe our three-stage pre-training in §2.2. Thirdly, we provide details of our supervised fine-tuning on diverse multimodal understanding and generation tasks in §2.3.

2.1 Modeling

As illustrated in Figure 1, the framework of MIO involves three parts: (1) multimodal tokenization, (2) causal multimodal modeling, and (3) multimodal de-tokenization.

Multimodal Tokenization. In our work, we use SEED-Tokenizer (Ge et al., 2023a) as our image tokenizer and SpeechTokenizer (Zhang et al., 2023b)

as our speech tokenizer. SEED-Tokenizer encodes images using a ViT (Dosovitskiy et al., 2021) derived from BLIP-2 (Li et al., 2023b), and then converts the encoded features into fewer tokens with causal semantics via Q-Former (Li et al., 2023b). These features are subsequently quantized into discrete tokens that are well-aligned with the language model backbone’s textual space. The codebook size for these discrete image tokens is 8192. SEED-Tokenizer transforms each image into a 224x224 resolution and quantizes it into 32 tokens. We use two special tokens, <image> and </image>, to indicate the start and end of an image.

As for videos, we first apply specific frame-cutting methods to convert videos into image sequences. In our training data processing procedures, the number of frames for each video is dynamically determined by its duration, the length of its context, or its scene switching³ to (1) avoid exceeding the LLM backbone’s context window limit, and (2) capture complete but concise information of the video. Each frame is then tokenized in the same manner as an image.

In terms of speech, SpeechTokenizer (Zhang et al., 2023b) leverages an 8-layer RVQ (Lee et al., 2022) to tokenize speech into tokens with 8 codebooks, with each codebook derived from one layer. Since the first layer’s quantization output is distilled from HuBERT (Hsu et al., 2021), which encodes more semantic information, SpeechTokenizer can separate content tokens and timbre tokens from a quantized speech. The first-layer quantization is treated as content quantization, while the remaining layers’ quantization is treated as timbre quantization. SpeechTokenizer encodes speech into 50 tokens per second for each codebook, resulting in 400 tokens per second with all eight codebooks. To improve context efficiency, we drop the last four layers’ codebooks and only use the content codebook and the first three timbre codebooks. Our vocabulary size for the speech modality is $1024 \times 4 = 4096$.

Since the open-source pretraining-level speech data is collected from individuals with diverse voices, the timbre tokens exhibit a relatively random and noisy pattern, while the content tokens are more fixed-pattern and better aligned with the corresponding transcriptions. Given these priors in speech tokens, it is important to choose the proper interleaving mode of speech tokens (Copet

³<https://github.com/Breakthrough/PySceneDetect>

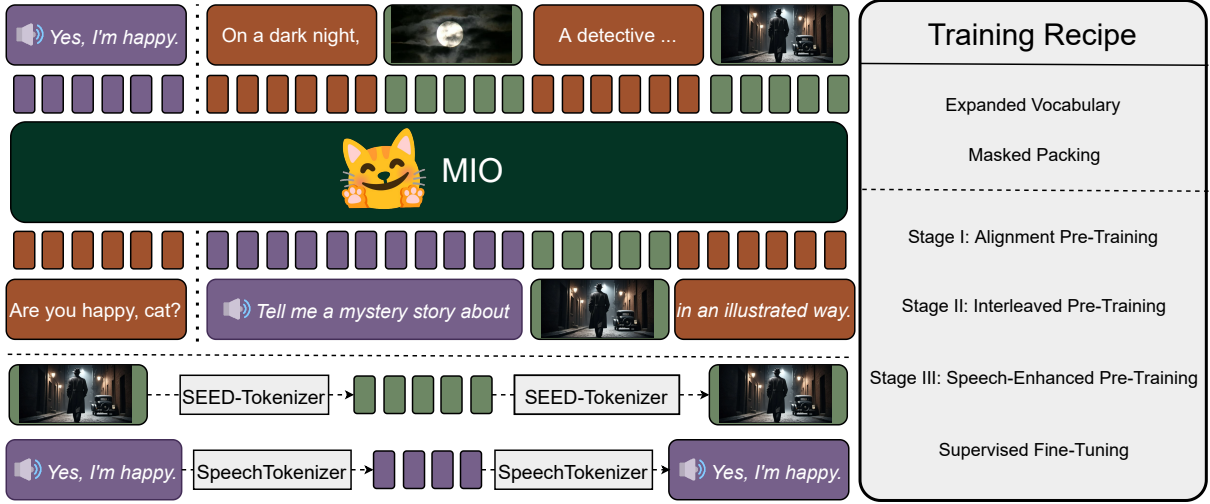


Figure 1: The framework of MIO and its training recipe.

et al., 2023). We denote the four codebooks as \mathcal{A} , \mathcal{B} , \mathcal{C} , and \mathcal{D} , where \mathcal{A} is the codebook for content tokens and the remaining three are for timbre tokens. For simplicity, assuming that we have only two tokens for each codebook in a tokenized speech sequence (i.e., a_1a_2 , b_1b_2 , c_1c_2 , and d_1d_2), there are two interleaving patterns for causal multimodal modeling: (1) sequential interleaving pattern: $a_1a_2b_1b_2c_1c_2d_1d_2$ and (2) alternating interleaving pattern: $a_1b_1c_1d_1a_2b_2c_2d_2$.

In our preliminary experiments, we observed that text-to-speech generation (TTS) training is difficult to converge when using the alternating interleaving pattern because the noisy and random timbre tokens ($b_1c_1d_1$) tend to mislead the continuations. Moreover, the speech-to-text understanding (ASR) performance improves much more slowly during training with the alternating interleaving pattern due to the sparsity of semantic information in the timbre tokens. Thus, we drop the timbre tokens for speech understanding and use the sequential interleaving pattern for speech generation. We use $\langle \text{spch} \rangle$ and $\langle / \text{spch} \rangle$ as special tokens to indicate the start and end of the speech token sequence.

Causal Multimodal Modeling. As illustrated in Figure 1, the speech and images, including video frames, are tokenized by SpeechTokenizer (Zhang et al., 2023b) and SEED-Tokenizer (Ge et al., 2023a), respectively. We add the 4096 speech tokens and 8192 image tokens to the LLM’s vocabulary. In addition, we introduce four new special tokens, namely $\langle \text{image} \rangle$, $\langle / \text{image} \rangle$, $\langle \text{spch} \rangle$, and $\langle / \text{spch} \rangle$, to the vocabulary. Consequently, the embedding layer of the LLM backbone and the

language modeling head are extended by $4096 + 8192 + 4 = 12292$ to support the embedding and generation of these new tokens. The image tokens contain *causal* semantics due to the use of a *Causal* Q-Former (Ge et al., 2023a), and the speech tokens are intrinsically causal due to their temporal nature. Therefore, these multimodal tokens are as suitable for autoregressive training as textual tokens, allowing us to unify the training objectives for understanding and generation of multimodal tokens into next-token-prediction with cross-entropy loss. The training objective is thus:

$$\mathcal{L} = - \sum_{t=1}^T \log P(x_t | x_{<t}; \theta) \quad (1)$$

where x_t denotes the discrete multimodal tokens, and θ denotes the LLM backbone parameters. We use Yi-6B-Base (AI et al., 2024) for initialization.

Furthermore, to eliminate the computational inefficiency caused by $\langle \text{pad} \rangle$ tokens, we use the masked packing strategy (Lu et al., 2023; Liu et al., 2024; Dehghani et al., 2023). The samples are concatenated along the sequence length until the context window is full. Then, we construct the causal attention mask for the tokens of each sample and mask out all the tokens of the other samples.

Multimodal De-Tokenization. After the generation of multimodal tokens, it is essential to use modality-specific decoders to reconstruct the images or speech from the codes. Specifically, for image tokens, we directly utilize SEED-Tokenizer’s decoder, which involves an MLP projection to convert the discrete codes into dense latents.

These latents condition an off-the-shelf diffusion model (Rombach et al., 2022) to generate the images in the pixel space (Ge et al., 2023a). The vanilla SpeechTokenizer (Zhang et al., 2023b) involves generating timbre tokens through a *non-autoregressive* model outside the language model, and then feeding the concatenated content and timbre tokens into the SpeechTokenizer decoder to synthesize speech. In our work, to inject the timbre priors into the multimodal language model itself, the timbre tokens are also generated by the *autoregressive* language model.

2.2 Pre-Training

As shown in Table 9, we use a three-stage strategy for pre-training, with each stage targeting different objectives. The three stages are: (1) Alignment Pre-training: This stage focuses on learning a multimodal representation more aligned with the language space. (2) Interleaved Pre-training: This stage aims to obtain a multimodal representation with richer contextual semantics. (3) Speech-enhanced Pre-training: This stage specifically enhances the model’s speech-related capabilities, while concurrently replaying data from other modalities. For more details on the pre-training data and its processing procedures, we refer the readers to Appendix A.

Stage I: Alignment Pre-Training. To fully leverage the superior capabilities of the pre-trained LLM backbone, it is essential to align the non-textual modality data representations with text. There are two types of pre-training data for image-text multimodal learning: (1) Image-text paired data: This data has well-aligned dependencies between images and text. (2) Image-text interleaved data: This data features more natural and contextual dependencies but is less aligned. Note that in our setting, video-text paired and interleaved data can be treated as image-text interleaved data, with videos being sequential images interleaved with text. Therefore, in this stage, we exclude the image-text interleaved data and video data to ensure the most aligned pattern between images and text.

Stage II: Interleaved Pre-Training. In this stage, we extend the data used for pre-training to include image-text interleaved data (including video-text data) as a novel image-text dependency pattern. The image-text interleaving pattern has a different nature compared to pairing patterns. Although (Li et al., 2023b) and (Sun et al., 2023c) argued

that interleaved image-text data mainly serves for *multimodal in-context learning*, we argue that it is also essential for context-aware image generation where images are generated based on specific context, rather than a precise description of the image content. For example, in image-text interleaved data, the text serves as the image’s preceding or continuing context, rather than its description. This pattern significantly differs from the previous descriptive image generation demonstrated in image-text paired data, where images are generated based on precise and detailed text that clearly describe the content of the images (Team et al., 2023). Therefore, context-aware image generation is essential for tasks like *chain-of-visual-thought reasoning* or *visual storytelling* (Team et al., 2023; Huang et al., 2016), where images are generated without textual descriptions. Due to the lack of benchmarks and evaluation metrics for context-aware image generation, we provide some demonstrations in §E.5 to showcase the potential of our model in visual storytelling, interleaved video-text generation, instructional image editing, chain-of-visual-thought reasoning, multimodal in-context learning, etc.

Moreover, in this stage, due to the extensive training on image-text paired data in Stage I, we can reduce its mixing ratio to the minimal essential scale for replay to avoid catastrophic forgetting. This allows us to increase the batch size for image-text interleaved data, video data, and speech data.

Stage III: Speech-Enhanced Pre-Training. The speech tokenizer that we use generates 200 tokens for each second of audio. Given that the duration of a speech sample can be 15 seconds, this results in around 3,000 tokens per sample. In comparison, the image tokenizer produces only 32 tokens per image. This creates a significant disparity in the number of tokens among different modalities. Consequently, our training data is dominated by speech tokens. If we mix all the different modalities according to their original proportions for training, the model would likely become overly focused on speech, at the expense of other modalities.

To address this issue, we implement a three-stage strategy that gradually increases the proportion of speech tokens. In Stage I, speech-text data accounts for 12.5% of the training tokens, which rises to 37.5% in Stage II, and finally reaches 75.0% in Stage III. This incremental increase in the proportion of speech tokens ensures that the model’s performance in non-speech modalities is not compro-

misled by the speech modality, while also allowing for the optimization of the model’s speech abilities.

Furthermore, we keep the data mixing ratio for other modalities of pre-training data at the minimal essential scales for replay, and we only use the high-quality subsets of them in this stage. This stage requires significantly fewer compute resources, due to the foundation laid in the previous stages.

We refer the reader to Appendix B for details about the hyperparameters and prompt templates.

2.3 Supervised Fine-Tuning

As shown in Table 10, our model undergoes comprehensive and systematic supervised fine-tuning (SFT) with 16 different tasks and 34 diverse open-source datasets. The chat template used for SFT is the same as that used for Yi-6B-Chat (AI et al., 2024), and only the assistant responses are supervised. We refer the reader to Appendix C for details about the hyperparameters and prompt templates.

3 Experiments

In this section, we first report our scores on MME-Unify (Xie et al., 2025), a widely used leaderboard for evaluating the multimodal performance of any-to-any MLLMs. Subsequently, we present our quantitative evaluation across various domains: image-related tasks (§3.2), speech-related tasks (§3.3), and video-related tasks (§3.4). Due to the lack of benchmarks for several advanced and emergent abilities of any-to-any multimodal LLMs, we provide qualitative demonstrations (§E.5) demonstrating these capabilities. The decoding hyperparameters and prompt templates are shown in Appendix D.

3.1 Scores on MME-Unify

We report our scores on a well-recognized third-party leaderboard for unified models, i.e., MME-Unify (Xie et al., 2025), in Table 2. According to the Overall metric on this leaderboard, our MIO model surpasses models such as Janus-Pro (Chen et al., 2025a) [2] and achieves multimodal unified modeling capabilities second only to Gemini-2.0-Flash, on the condition that MIO covers a broader range of tasks compared to Gemini-2.0-Flash, demonstrating MIO’s top-tier multimodal unified modeling capabilities. Furthermore, MIO exhibits impressive multimodal understanding capabilities, ranking highly among unified MLLMs according to the Understanding metric, and demon-

strates exceptionally leading performance in generation tasks, as evidenced by the Generation metric. Additionally, MIO extends its capabilities to include speech support, a feature not offered by many competing models. However, we observe that MIO demonstrates limited performance on the unified tasks. This is primarily because the unified tasks in MME-Unify involve single-choice questions with images as options, a type of data that MIO’s training lacks.

3.2 Image-Related Tasks

Image Understanding. We compare our models with Emu (Sun et al., 2023c), SEED-LLaMA (Ge et al., 2023b), AnyGPT (Zhan et al., 2024), Flamingo (Alayrac et al., 2022), Kosmos-1 (Huang et al., 2023), MetaLM (Hao et al., 2022), IDEFICS (Laurençon et al., 2023), CM3Leon (Yu et al., 2023), and InstructBLIP (Dai et al., 2023). We evaluate our models in diverse tasks: (1) image captioning on MS-COCO (Lin et al., 2014) Karpathy test split with CIDEr score (Vedantam et al., 2014) as the metric, (2) three visual question-answering benchmarks, i.e., VQAv2 (Goyal et al., 2016) (test-dev split), OK-VQA (Marino et al., 2019) (val split), and VizWiz (Gurari et al., 2018), with VQA accuracy as the metric, and (3) SEED-Bench (Li et al., 2023a), a comprehensive visual question-answering benchmark including 9 dimensions with MCQ accuracy as the metric. The scores for all baselines are copied from their reports. As shown in Table 3, our MIO-Instruct is ranked in the top group among all baselines, demonstrating its competitive image understanding performance. Although SEED-LLaMA achieved better scores, we additionally support the speech modality. It is noteworthy that MIO, with a size of approximately 7 billion parameters, outperforms several larger models such as Emu-14B and even IDEFICS-80B.

Image Generation. We compare our models with Emu (Sun et al., 2023c), SEED-LLaMA (Ge et al., 2023b), GILL (Koh et al., 2023), and AnyGPT (Zhan et al., 2024) for image generation. We use two benchmarks: MS-COCO (Lin et al., 2014) Karpathy test split and Flickr30K (Plummer et al., 2015). Following GILL (Koh et al., 2023) and SEED-LLaMA (Ge et al., 2023b), we use CLIP-I as the metric that evaluates the similarity between the generated and the ground-truth images with the CLIP image encoder (Radford et al., 2021). As shown in Table 4 and 13, the pre-trained and

| Model | Overall | Understanding | Generation | Unify |
|---------------------|---------|---------------|------------|-------|
| Gemini2.0-flash-exp | 45.6 | 65.2 | 29.8 | 40.7 |
| MIO-Instruct | 37.2 | 41.5 | 53.5 | 16.6 |
| SEED-LLaMA | 28.5 | 39.5 | 23.5 | 22.3 |
| Anole | 18.6 | 13.6 | 20.0 | 22.3 |
| VILA-U | 18.6 | 40.0 | 15.8 | - |
| Janus-Pro | 18.1 | 48.4 | 5.9 | - |
| Janus-Flow | 16.3 | 43.4 | 5.5 | - |
| Emu3 | 13.8 | 33.2 | 8.2 | - |
| Show-o | 12.7 | 31.0 | 7.3 | - |

Table 2: Scores on MME-Unify (Xie et al., 2025). “-” indicates that the model lacks the capability to complete the specified task. **Understanding** tasks involve single image perception & understanding, multiple & interleaved image-text understanding, and video perception & understanding. **Generation** tasks include conditional image-to-video generation, fine-grained image reconstruction, text-guided image editing, text-to-image generation, text-to-video generation, and video prediction. **Unified** tasks encompass image editing and explaining, common sense question answering, auxiliary lines, SpotDiff, and visual chain-of-thought reasoning. Arrange from top to bottom based on the overall score from highest to lowest.

| Models | Imagen | Speech | COCO(↑) | VQAv2(↑) | OKVQA(↑) | VizWiz(↑) | SEED Bench(↑) |
|---------------------|--------|--------|---------|----------|----------|-----------|---------------|
| Emu-Base (14B) | ✓ | ✗ | 112.4 | 52.0 | 38.2 | 34.2 | 47.3 |
| Emu-I (14B) | ✗ | ✗ | 120.4 | 57.2 | 43.4 | 32.2 | 58.0 |
| SEED-LLaMA-I (8B) | ✓ | ✗ | 124.5 | 66.2 | 45.9 | 55.1 | 51.5 |
| AnyGPT (8B) | ✓ | ✓ | 107.5 | - | - | - | - |
| Flamingo (9B) | ✗ | ✗ | 79.4 | 51.8 | 44.7 | 28.8 | 42.7 |
| Flamingo (80B) | ✗ | ✗ | 84.3 | 56.3 | 31.6 | - | - |
| Kosmos-1 (1.6B) | ✗ | ✗ | 84.7 | 51.0 | - | 29.2 | - |
| MetaLM (1.7B) | ✗ | ✗ | 82.2 | 41.1 | 11.4 | - | - |
| IDEFICS-I (80B) | ✗ | ✗ | 117.2 | 37.4 | 36.9 | 26.2 | 53.2 |
| CM3Leon (7B) | ✓ | ✗ | 61.6 | 47.6 | 23.8 | 37.6 | - |
| InstructBLIP (8.1B) | ✗ | ✗ | - | - | - | 34.5 | 58.8 |
| MIO-Instruct (7B) | ✓ | ✓ | 120.4 | 65.5 | 39.9 | 53.5 | 54.4 |

Table 3: Experimental results for image understanding abilities. “Imagen” denotes whether the model is capable of generating images. “Speech” denotes whether the model supports speech modality. “I” denotes the instruction tuned version. The metrics used are CIDEr for COCO, MCQ accuracy for the SEED Bench, and VQA accuracy for the other tasks, following the standard procedures.

| Models | MS-COCO(↑) | Flickr30K(↑) |
|--------------|------------|--------------|
| Emu-Base | 66.46 | 64.82 |
| SEED-LLaMA | 69.07 | 65.54 |
| SEED-LLaMA-I | 70.68 | 66.55 |
| GILL | 67.45 | 65.16 |
| AnyGPT | 65.00 | - |
| MIO-Base | 64.15 | 62.71 |
| MIO-Instruct | 67.76 | 68.97 |

Table 4: Image generation evaluation by CLIP-I score.

instruction-tuned model of MIO both have competitive image generation abilities. Note that beyond single image generation, our model can also exhibit multi-image generation capabilities such as generating visual stories, image sequences, and visual thoughts as illustrated in §E.5.

3.3 Speech-Related Tasks

We evaluate the speech understanding and generation abilities of MIO on ASR and TTS tasks. Wav2vec 2.0 (Baevski et al., 2020), Whisper Large V2 (Radford et al., 2023), and AnyGPT (Zhan et al., 2024) are the baselines for ASR tasks, while VALL-E (Wang et al., 2023a), USLM (Zhang et al., 2023b), and AnyGPT (Zhan et al., 2024) are the baselines for TTS tasks. The test set used for ASR evaluation is LibriSpeech (Panayotov et al., 2015), while the test set used for TTS evaluation is VCTK (Veaux et al., 2017) following AnyGPT (Zhan et al., 2024)’s practice. The Whisper medium model is used to transcribe the speech generated for the TTS task. The WER (word error rate) is computed by comparing the generated

| Models | ASR | Models | TTS |
|--------------|------|--------------|------|
| Wav2vec | 2.7 | VALL-E | 7.9 |
| Whisper | 2.7 | USLM | 6.5 |
| AnyGPT | 8.5 | AnyGPT | 8.5 |
| MIO-Base | 6.3 | MIO-Base | 12.0 |
| MIO-Instruct | 10.3 | MIO-Instruct | 4.2 |

Table 5: Speech ability evaluation by WER (\downarrow).

| Models | MSVDQA | MSRVTT-QA |
|---------------------|--------|-----------|
| Flamingo (9B) | 30.2 | 13.7 |
| BLIP-2 (4.1B) | 33.7 | 16.2 |
| InstructBLIP (8.1B) | 41.8 | 22.1 |
| Emu-Instruct (14B) | 32.4 | 14.0 |
| SEED-LLaMA-I (8B) | 40.9 | 30.8 |
| MIO-Instruct | 42.6 | 35.5 |

Table 6: Video understanding evaluation using accuracy.

transcribed text with the ground-truth transcription after text normalization⁴.

As shown in Table 5, our models exhibit speech performance comparable to the speech-specific baselines and outperform the AnyGPT baseline. It is important to note that although AnyGPT is capable of generating content tokens for speech, it lacks the ability to generate timbre tokens, which necessitates the use of an additional voice cloning model. In contrast, our models generate both content and timbre tokens, making the TTS tasks more challenging for our models compared to AnyGPT. Nonetheless, after instruction tuning, our model still achieves better TTS performance. More evaluations of the TTS and Speech-to-Speech generation performance are provided in Appendix E.3 and E.2.

3.4 Video-Related Tasks

We compare MIO with Flamingo (Alayrac et al., 2022), BLIP-2 (Li et al., 2023b), InstructBLIP (Dai et al., 2023), Emu (Sun et al., 2023c), and SEED-LLaMA (Ge et al., 2023b) for video understanding. The models are evaluated on MSVDQA (Chen and Dolan, 2011a) and MSRVTT-QA (Xu et al., 2017). The results are presented in Table 6. Our model achieves highest scores compared to all baselines. Due to the lack of video generation benchmarks in our setting, we provide examples in §E.5. These results demonstrate superior performance of our models in both video understanding and generation.

⁴<https://github.com/openai/whisper/blob/main/whisper/normalizers/english.py>

| Models | MMLU(\uparrow) |
|-----------------|--------------------|
| LLAMA-1-7B-Base | 33.0 |
| LLAMA-2-7B-Chat | 47.9 |
| SEED-LLAMA-8B-I | 36.1 |
| AnyGPT-Base | 26.4 |
| AnyGPT-Chat | 27.4 |
| MIO-Instruct | 45.7 |

Table 7: Language-only evaluation.

| Models | OmniBench(\uparrow) |
|---------------------|-------------------------|
| Gemini-1.5-Pro | 42.67 |
| Reka-Core-20240501 | 31.52 |
| AnyGPT (8B) | 17.77 |
| video-SALMONN (13B) | 34.11 |
| Unified-IO 2 (6.8B) | 34.24 |
| MIO-Instruct (7B) | 36.96 |

Table 8: Results for trimodal understanding.

3.5 Language-only Tasks

We evaluate our models on MMLU (Hendrycks et al., 2021). The baselines are two LLaMA variants (Touvron et al., 2023a,b), the instruction-tuned SEED-LLaMA (Ge et al., 2023b), and AnyGPT (Zhan et al., 2024). For the MMLU benchmark, we conduct zero-shot evaluation experiments using the official evaluation code. The experimental results are shown in Table 7. We can observe that our models have superior language-only performance compared with all any-to-any MM-LLM baselines and even surpass LLaMA-1-7B-Base, an advanced pure language model.

3.6 Ablation Studies

Generality for Trimodal Understanding. We evaluate our model using OmniBench (Li et al., 2024c), which incorporates text, image, and speech modalities as inputs, requiring the model to choose one of four options as the correct answer to determine accuracy. Although MIO acquires its multimodal understanding abilities via dual-modal training, the results in Table 8 indicate that MIO exhibits superior trimodal comprehension abilities.

Please refer to Appendix E.6 for more ablations including the effect of different image tokenizers.

4 Related Works

With the success of LLMs, MM-LLMs have emerged, extending LLMs to handle images, speech, and video (Liu et al., 2023; Li et al., 2024a, 2023b; Bai et al., 2023; Team, 2025; OpenAI, 2023;

Shi et al., 2025; Chen et al., 2025b). These models typically align image features with text embeddings. For instance, BLIP-2 (Li et al., 2023b) uses CLIP-ViT and a Q-Former for alignment, while LLaVA (Liu et al., 2023; Li et al., 2024a) uses linear projections. These models perform well in visual question answering and commonsense reasoning. Recent models like LLaSM (Shu et al., 2023) and Qwen2.5-VL (Team, 2025) extend to speech and video. However, most focus on understanding rather than generation, limiting their utility in fully multimodal tasks.

To address this, recent work explores any-to-any MM-LLMs that generate outputs across modalities without intermediate language. Key approaches include Discrete-In-Discrete-Out (DIDO), Continuous-In-Discrete-Out (CIDO), Continuous-In-Continuous-Out (CICO), and Autoregressive + Diffusion (AR + Diff), discussed in §1. DIDO is used in SEED-LLaMA (Ge et al., 2023b), AnyGPT (Zhan et al., 2024), and Chameleon (Team, 2024); CIDO in DaVinCi (Diao et al., 2023), Gemini (Team et al., 2023), and Unified-IO 2 (Lu et al., 2023); CICO in Emu (Sun et al., 2023c,a) and DreamLLM (Dong et al., 2023); AR + Diff in Transfusion (Zhou et al., 2024), Showo (Xie et al., 2024), and MAR (Li et al., 2024b).

However, these models have limitations. DreamLLM (CICO, (Dong et al., 2023)) and CIDO models suffer from inconsistencies between input and output forms for multimodal data, making it hard to generate interleaved multimodal sequences where an image functions in a coupled way as both input and output. Emu2 (CICO, (Sun et al., 2023a)) faces challenges with MSE loss for training continuous outputs, as well as with the uni-modal assumption of the Gaussian distribution in the MSE loss. Transfusion (AR + Diff, (Zhou et al., 2024)) applies noise to images from the input side to support multimodal generation with diffusion modeling, and relies on VAE (Kingma and Welling, 2013) rather than CLIP (Radford et al., 2021) features for denoising, which largely trade off the multimodal understanding abilities. To mitigate these issues, we adopt the DIDO approach. A comprehensive comparison of our models with other any-to-any MM-LLMs is presented in Table 1.

5 Conclusion

In conclusion, MIO represents an advancement in the realm of multimodal foundation models. By

employing a rigorous four-stage training process, MIO successfully integrates and aligns discrete tokens across text, image, video, and speech modalities. This comprehensive approach enables MIO to understand and generate multimodal content in an end-to-end, autoregressive manner, addressing the limitations of current multimodal large language models. Our experimental results showcase its competitive performance across a variety of benchmarks compared to the dual-modality baselines and other any-to-any multimodal large language models. With the any-to-any and multimodal interleaved output features, MIO exhibits novel emergent abilities such as interleaved video-text generation, chain-of-visual-thought reasoning, etc.

Limitations

While MIO demonstrates advancements in multimodal understanding and generation, it has certain limitations. First, the model’s performance is constrained by the quality and diversity of the training data, particularly for speech and video modalities, where open-source datasets may not fully capture the range of real-world scenarios. Second, the computational complexity of handling four modalities simultaneously requires substantial resources, potentially limiting scalability and accessibility for smaller research groups. Third, while MIO excels in multimodal interleaved sequence generation, the evaluation of such capabilities lacks standardized benchmarks, making it challenging to quantitatively compare with other models. Finally, the model’s ability to handle extremely long-context multimodal sequences or highly specialized tasks may be limited due to the fixed context window. Future work could address these by expanding dataset diversity, optimizing computational efficiency, and developing tailored evaluation metrics for advanced multimodal tasks.

Acknowledgements

This work was supported in part by the Jiangsu Science and Technology Major Project (BG2024031) and Nanjing University AI & AI for Science Funding (2024300540)

References

01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang,

- Shiming Yang, Tao Yu, and 13 others. 2024. Yi: Open foundation models by 01.ai. *arXiv preprint arXiv: 2403.04652*.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.
- R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber. 2020. Common voice: A massively-multilingual speech corpus. In *Proceedings of the 12th Conference on Language Resources and Evaluation (LREC 2020)*, pages 4211–4215.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *ArXiv preprint*, abs/2308.12966.
- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. 2021. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. 2024. [Improving image generation with better captions](#).
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. 2023. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv: 2311.15127*.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. 2023. [Audiolm: A language modeling approach to audio generation](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2523–2533.
- Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2023. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv: 2303.12712*.
- David Chen and William Dolan. 2011a. [Collecting highly parallel data for paraphrase evaluation](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 190–200, Portland, Oregon, USA. Association for Computational Linguistics.
- David L. Chen and William B. Dolan. 2011b. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics (ACL-2011)*, Portland, OR.
- Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Yujun Wang, and 2 others. 2021. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv: 2106.06909*.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. 2025a. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv: 2501.17811*.
- Xinlong Chen, Yuanxing Zhang, Yushuo Guan, Bohan Zeng, Yang Shi, Sihan Yang, Pengfei Wan, Qiang Liu, Liang Wang, and Tieniu Tan. 2025b. [Versavid-r1: A versatile video understanding and reasoning model from question answering to captioning tasks](#). *arXiv preprint arXiv: 2506.09079*.
- Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre D’efossez. 2023. [Simple and controllable music generation](#). *Neural Information Processing Systems*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *ArXiv preprint*, abs/2305.06500.
- Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv: 2307.08691*.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, A. Rudra, and Christopher R’e. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Neural Information Processing Systems*.

- Nilaksh Das, Saket Dingliwal, Srikanth Ronanki, Rohit Paturi, David Huang, Prashant Mathur, Jie Yuan, Dhanush Bekal, Xing Niu, Sai Muralidhar Jayanthi, Xilai Li, Karel Mundnich, Monica Sunkara, Sundararajan Srinivasan, Kyu J Han, and Katrin Kirchhoff. 2024. Speechverse: A large-scale generalizable audio language model. *arXiv preprint arXiv:2405.08295*.
- Mostafa Dehghani, Basil Mustafa, Josip Djolonga, J. Heek, Matthias Minderer, Mathilde Caron, A. Steiner, J. Puigcerver, Robert Geirhos, Ibrahim M. Alabdulmohsin, Avital Oliver, Piotr Padlewski, A. Gritsenko, Mario Luvci'c, and N. Houlsby. 2023. Patch n' pack: Navit, a vision transformer for any aspect ratio and resolution. *Neural Information Processing Systems*.
- Shizhe Diao, Wangchunshu Zhou, Xinsong Zhang, and Jiawei Wang. 2023. Write and paint: Generative vision-language models are unified modal learners. In *The Eleventh International Conference on Learning Representations*.
- Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, Xiangwen Kong, Xiangyu Zhang, Kaisheng Ma, and Li Yi. 2023. Dream-llm: Synergistic multimodal comprehension and creation. *arXiv preprint arXiv: 2309.11499*.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. *Preprint*, arXiv:2010.11929.
- Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. 2022. Clap: Learning audio concepts from natural language supervision. *arXiv preprint arXiv: 2206.04769*.
- Patrick Esser, Robin Rombach, and Björn Ommer. 2020. Taming transformers for high-resolution image synthesis. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12868–12878.
- Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. 2024. Llama-omni: Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*.
- Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, and 1 others. 2024. Vita: Towards open-source interactive omni multimodal llm. *arXiv preprint arXiv:2408.05211*.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wan-jun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. 2023. G-llava: Solving geometric problem with multi-modal large language model. *arXiv preprint arXiv: 2312.11370*.
- Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. 2023a. Planting a seed of vision in large language model. *arXiv preprint arXiv:2307.08041*.
- Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. 2023b. Making llama see and draw with seed tokenizer. *arXiv preprint arXiv:2310.01218*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *International Journal of Computer Vision*.
- Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. *arXiv preprint arXiv: 1802.08218*.
- Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. 2022. Language models are general-purpose interfaces. *ArXiv preprint*, abs/2206.06336.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *arXiv preprint arXiv: 2106.07447*.
- Shaohan Huang, Li Dong, Wenhui Wang, Y. Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, O. Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. 2023. Language is not all you need: Aligning perception with language models. *Neural Information Processing Systems*.
- Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. 2016. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1233–1239.
- Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. 2024. Mantis: Interleaved multi-image instruction tuning. *arXiv preprint arXiv: 2405.01483*.
- Wei Kang, Xiaoyu Yang, Zengwei Yao, Fangjun Kuang, Yifan Yang, Liyong Guo, Long Lin, and Daniel Povey. 2023. Libriheavy: a 50,000 hours asr corpus with punctuation casing and context. *Preprint*, arXiv:2309.08105.

- Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. 2023. Generating images with multimodal language models. *NeurIPS*.
- LAION. 2022. Laion coco: 600m synthetic captions from laion-2b-en. <https://laion.ai/blog/laion-coco/>.
- 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. 2023. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Preprint*, arXiv:2306.16527.
- Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive image generation using residual quantization. *Preprint*, arXiv:2203.01941.
- Bo Li*, Peiyuan Zhang*, Kaichen Zhang*, Fanyi Pu*, Xinrun Du, Yuhao Dong, Haotian Liu, Yuanhan Zhang, Ge Zhang, Chunyuan Li, and Ziwei Liu. 2024. Lmms-eval: Accelerating the development of large multimodal models.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. *Preprint*, arXiv:2307.16125.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ArXiv preprint*, abs/2301.12597.
- KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023c. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.
- Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. 2024b. Autoregressive image generation without vector quantization. *arXiv preprint arXiv: 2406.11838*.
- Yizhi Li, Ge Zhang, Yinghao Ma, Ruibin Yuan, Kang Zhu, Hangyu Guo, Yiming Liang, Jiaheng Liu, Jian Yang, Siwei Wu, Xingwei Qu, Jinjie Shi, Xinyue Zhang, Zhenzhu Yang, Xiangzhou Wang, Zhaoxiang Zhang, Zachary Liu, Emmanouil Benetos, Wenhao Huang, and Chenghua Lin. 2024c. Omnibench: Towards the future of universal omni-language models. *arXiv preprint arXiv: 2409.15272*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer.
- Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. 2024. World model on million-length video and language with ringattention. *arXiv preprint*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *ArXiv preprint*, abs/2304.08485.
- I. Loshchilov and F. Hutter. 2017. Decoupled weight decay regularization. *International Conference on Learning Representations*.
- Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. 2023. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. *Preprint*, arXiv:2312.17172.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 3195–3204.
- Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. 2019. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2630–2640.
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. 2019. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*.
- Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Tushar Nagarajan, Matt Smith, Shashank Jain, Chun-Fu Yeh, Prakash Murugesan, Peyman Heidari, Yue Liu, Kavya Srinet, Babak Damavandi, and Anuj Kumar. 2023. Anymal: An efficient and scalable any-modality augmented language model. *arXiv preprint arXiv: 2309.16058*.

- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mándry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. Gpt-4o system card. *arXiv preprint arXiv: 2410.21276*.
- OpenAI. 2023. Gpt-4v(ision) system card.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv: 2303.08774*.
- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24.
- Junting Pan, Keqiang Sun, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. 2023. *Journeydb: A benchmark for generative image understanding*. *Preprint*, arXiv:2307.00716.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. *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.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv: 2306.01116*.
- Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv: 2212.04356*.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. 2022. *High-resolution image synthesis with latent diffusion models*. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- 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, Hannah Muckenhirn, Dirk Padfield, James Qin, Danny Rozenberg, Tara Sainath, Johan Schalkwyk, Matt Sharifi, Michelle Tadmor Ramanovich, Marco Tagliasacchi, and 11 others. 2023. Audiopalm: A large language model that can speak and listen. *arXiv preprint arXiv: 2306.12925*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face. *arXiv preprint arXiv: 2303.17580*.
- Yang Shi, Jiaheng Liu, Yushuo Guan, Zhenhua Wu, Yuanxing Zhang, Zihao Wang, Weihong Lin, Jingyun Hua, Zekun Wang, Xinlong Chen, Bohan Zeng, Wentao Zhang, Fuzheng Zhang, Wenjing Yang, and Di Zhang. 2025. Mavors: Multi-granularity video representation for multimodal large language model. *arXiv preprint arXiv: 2504.10068*.
- Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang, Ruihua Zhang, Daochen Shi, Qiqi Xiang, and Yemin Shi. 2023. Llasmm: Large language and speech model. *arXiv preprint arXiv: 2308.15930*.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326.
- Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and 1 others. 2023a. Generative multimodal models are in-context learners. *arXiv preprint arXiv:2312.13286*.

- Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023b. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv: 2303.15389*.
- Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2023c. [Generative pretraining in multimodality](#). *Preprint*, arXiv:2307.05222.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. 2023. Salmonn: Towards generic hearing abilities for large language models. *arXiv preprint arXiv: 2310.13289*.
- Chameleon Team. 2024. Chameleon: Mixed-modal early-fusion foundation models. *arXiv e-prints*, pages arXiv–2405.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, and 1326 others. 2023. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv: 2312.11805*.
- Qwen Team. 2025. [Qwen2.5-vl](#).
- Teknium. 2023. [Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants](#).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *ArXiv preprint*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *ArXiv preprint*, abs/2307.09288.
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2017. [Neural discrete representation learning](#). *ArXiv*, abs/1711.00937.
- Christophe Veaux, Junichi Yamagishi, and Kirsten MacDonald. 2017. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2014. Cider: Consensus-based image description evaluation. *arXiv preprint arXiv: 1411.5726*.
- Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, and 1 others. 2023a. Neural codec language models are zero-shot text to speech synthesizers. *arXiv preprint arXiv:2301.02111*.
- Tiannan Wang, Jiamin Chen, Qingrui Jia, Shuai Wang, Ruoyu Fang, Huilin Wang, Zhaowei Gao, Chunzhao Xie, Chuou Xu, Jihong Dai, Yibin Liu, Jialong Wu, Shengwei Ding, Long Li, Zhiwei Huang, Xinle Deng, Teng Yu, Gangan Ma, Han Xiao, and 27 others. 2024a. Weaver: Foundation models for creative writing. *arXiv preprint arXiv: 2401.17268*.
- Wenbin Wang, Yang Song, and Sanjay Jha. 2024b. Globe: A high-quality english corpus with global accents for zero-shot speaker adaptive text-to-speech. *arXiv preprint arXiv: 2406.14875*.
- Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhao Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Chen, Jie Fu, and Junran Peng. 2023b. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. *arXiv preprint arXiv: 2310.00746*.
- Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. 2004. [Image quality assessment: from error visibility to structural similarity](#). *IEEE Transactions on Image Processing*, 13:600–612.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023. [Next-gpt: Any-to-any multimodal llm](#). *arXiv preprint arXiv: 2309.05519*.
- Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. 2024. [Show-o: One single transformer to unify multimodal understanding and generation](#). *arXiv preprint arXiv: 2408.12528*.
- Wulin Xie, Yi-Fan Zhang, Chaoyou Fu, Yang Shi, Bingyan Nie, Hongkai Chen, Zhang Zhang, Liang Wang, and Tieniu Tan. 2025. Mme-unify: A comprehensive benchmark for unified multimodal understanding and generation models. *arXiv preprint arXiv: 2504.03641*.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*, pages 1645–1653.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A large video description dataset for bridging video and language. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 5288–5296.

- Zhiyang Xu, Trevor Ashby, Chao Feng, Rulin Shao, Ying Shen, Di Jin, Qifan Wang, and Lifu Huang. 2023. [Vision-flan:scaling visual instruction tuning](#).
- Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, Candace Ross, Adam Polyak, Russell Howes, Vasu Sharma, Puxin Xu, Hovhannes Tamoyan, Oron Ashual, Uriel Singer, Shang-Wen Li, and 8 others. 2023. [Scaling autoregressive multi-modal models: Pretraining and instruction tuning](#). *Preprint*, arXiv:2309.02591.
- Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. 2021. [MERLOT: Multimodal neural script knowledge models](#). In *Advances in Neural Information Processing Systems*.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. [Sigmoid loss for language image pre-training](#). *IEEE International Conference on Computer Vision*.
- 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. 2024. [Anygpt: Unified multimodal llm with discrete sequence modeling](#). *ArXiv*, abs/2402.12226.
- Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023a. [Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities](#). *Preprint*, arXiv:2305.11000.
- Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. 2024. [Magicbrush: A manually annotated dataset for instruction-guided image editing](#). *Advances in Neural Information Processing Systems*, 36.
- Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and Xipeng Qiu. 2023b. [Spechtokenizer: Unified speech tokenizer for speech language models](#). *Preprint*, arXiv:2308.16692.
- Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2023c. [Llavar: Enhanced visual instruction tuning for text-rich image understanding](#). *arXiv preprint arXiv:2306.17107*.
- Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hongsheng Li, Yu Qiao, Wanli Ouyang, and Xiangyu Yue. 2023d. [Meta-transformer: A unified framework for multimodal learning](#). *arXiv preprint arXiv:2307.10802*.
- Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. 2023. [Mmicl: Empowering vision-language model with multi-modal in-context learning](#). *ArXiv preprint*, abs/2309.07915.
- Kaizhi Zheng, Xuehai He, and Xin Eric Wang. 2023. [Minigt-5: Interleaved vision-and-language generation via generative vokens](#).
- Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. 2024. [Transfusion: Predict the next token and diffuse images with one multi-modal model](#). *arXiv preprint arXiv:2408.11039*.
- Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. 2023. [Multimodal c4: An open, billion-scale corpus of images interleaved with text](#). *arXiv preprint arXiv:2304.06939*.

A Pre-training Data

Pre-training Data Sources. The pre-training data sources involve six types:

1. Image-text paired data: SBU (Ordenez et al., 2011), CC3M (Sharma et al., 2018), LAION-COCO (LAION, 2022), and JourneyDB (Pan et al., 2023), where JourneyDB only serves for image generation.
2. Language-only data: RefinedWeb (Penedo et al., 2023).
3. Image-text interleaved data: OBELICS (Laurençon et al., 2023), MMC4-core-ff (Zhu et al., 2023).
4. Video-text paired data: WebVid-10M (Bain et al., 2021).
5. Video-text interleaved data: HowTo-100M (Miech et al., 2019), Youtube-Temporal-180M (Zellers et al., 2021).
6. Speech-text paired data: Libriheavy (Kang et al., 2023).

Pre-training Data Processing. We have different data processing procedures for different data types illustrated in §A following Emu (Sun et al., 2023c) and Qwen-VL (Bai et al., 2023):

1. Image-text paired data: we remove pairs with more than 2:1 aspect ratio or smaller than 224×224 resolution of the image. We remove pairs with more than 0.27 CLIP scores. We remove non-English pairs. We randomly place the image or text at the forefront for generating captions based on images and vice versa.
2. Language-only data: we use the same data processing pipeline as used in Yi (AI et al., 2024).
3. Image-text interleaved data: we filter the data using a CLIP score threshold of 0.25, and follow the same procedure as illustrated in Emu (Sun et al., 2023c).
4. Video-text paired data: we randomly place the frames or text at the forefront for generating captions based on frames and vice versa. 60% of the pairs are text-to-video, while 40% of the pairs are video-to-text. We sample 4 to 8 frames of each video for training according to the text lengths.

5. Video-text interleaved data: We first use PySceneDetect to extract key frames from the video based on scene changes, following the practice of Stable Video Diffusion (Blattmann et al., 2023). Then, for each video clip between two key frames, we extract a central frame for textual caption generation with BLIP-2 (Li et al., 2023b). Additionally, the video clips between key frames are processed using ASR (automatic speech recognition) tools to extract subtitles. The ASR text and captions are then integrated and refined using Yi-34B-Chat (AI et al., 2024), resulting in a single text segment. These text segments, along with the key frames and central frames, form the video-text interleaved data.
6. Speech-text paired data: we remove speeches with more than 15 seconds.

B Pre-training Details

Hyperparameters. We enable Flash Attention (Dao et al., 2022; Dao, 2023) during pre-training. Gradient clipping is set to 1.0 for all stages. The maximum sequence length for training is 2800 tokens. We use a cosine learning rate scheduler with a peak learning rate of $3e-5$ and a warmup ratio of 0.03. The optimizer used is AdamW (Loshchilov and Hutter, 2017).

Prompt Templates. The prompt template is only necessary for paired datasets. For image-text paired data, we use the prompt templates of “{image} The caption of this image is: {caption}” and “Please generate an image of “{caption}”: {image}”. For video-text paired data: we use the prompt templates of “Please describe the following video: {image} {description}” and “Please generate a video for “{description}”: {video}”. For speech-text paired data: we use the prompt templates of “{speech} Transcribe this speech: {transcription}” and “Please generate a speech of “{transcription}”: {speech}” during Stage I and Stage II. While for Stage III, we change the ASR prompt template into “{speech} The transcription of this speech is: {transcription}”.

C Supervised Fine-Tuning Details

Supervised Fine-Tuning Data. As shown in Table 10, we use 16 tasks with 34 datasets for a comprehensive supervised fine-tuning.

| Pre-training Stage Objective | Stage I Multimodal Alignment | Stage II Multimodal Interleaving | Stage III Speech Enhancement |
|------------------------------|--|--|-------------------------------------|
| Image-Text Pair | SBU, CC3M, LAION-COCO, JourneyDB | SBU, CC3M, LAION-COCO, JourneyDB | CC3M LAION-COCO |
| Language-Only | RefinedWeb | RefinedWeb | RefinedWeb |
| Image-Text Inter | - | OBELICS, MMC4-core-ff | MMC4-core-ff |
| Video-Text Pair | - | WebVid-10M | WebVid-10M |
| Video-Text Inter | - | HowTo-100M, YT-Temporal- 180M | HowTo-100M, YT-Temporal- 180M |
| Speech-Text Pair | Libriheavy | Libriheavy | Libriheavy |
| GPUs | 128 A800-80GB | 128 A800-80GB | 8 A800-80GB |
| Training Steps | 24,800 | 12,800 | 32,200 |
| Batch Size | 12:2:0:2 | 2:2:6:6 | 2:1:1:12 |

Table 9: Pre-training stages and their details. We use “Inter” to denote “Interleaved” for short. We provide batch sizes for each data type per GPU in image-text pair data:language-only data:(image-text interleaved data + video data):speech-text pair data. See Appendix A and Appendix B for more details including pre-training data sources, data cleaning procedures, pre-training hyperparameters, etc.

Prompt Templates. The chat template is the same as used in Yi (AI et al., 2024). The system prompt is unified as: “You are MIO, an AI assistant capable of understanding and generating images, text, videos, and speech, selecting the appropriate modality according to the context.” except for speech generation and TTS whose system prompts are “You are MIO, an AI assistant capable of understanding images, text, videos, and speech, and generating speech. Please respond to the user with speech only, starting with <spch> and ending with </spch>.” to avoid randomness of the output modality.

Hyperparameters. Similar to pre-training (*c.f.*, Appendix B), we enable Flash Attention (Dao et al., 2022; Dao, 2023) during supervised fine-tuning. Gradient clipping is set to 1.0. The maximum sequence length for training is 2800 tokens. We use a cosine learning rate scheduler with a peak learning rate of $3e-5$ and a warmup ratio of 0.03. The optimizer used is AdamW (Loshchilov and Hutter, 2017).

D Evaluation Details.

Hyperparameters. The decoding strategies and hyperparameters are quite important for a superior performance. As shown in Table 11, we use

different sets of parameters for different output modalities.

Prompt Templates. The prompt templates used for evaluating pre-training checkpoints are the same as used during pre-training. For SFT checkpoint evaluation, we list the prompt templates in Table 12.

E More Experiments

E.1 Image Generation Evaluation

In Table 13, we compute two additional automatic metrics for evaluating image generation, *i.e.*, SSIM (Wang et al., 2004) and Aesthetic Predictor v2.5⁵ for the evaluation of structural integrity and aesthetics, respectively. SSIM (Structural Similarity Index Measure) evaluates the perceptual similarity between the generated images and the ground-truth images, focusing on luminance, contrast, and structure, with scores ranging from -1 (dissimilar) to 1 (identical). Aesthetic Predictor V2.5 is a SigLIP (Zhai et al., 2023)-based predictor that evaluates the aesthetics of an image on a scale from 1 to 10 (10 is the best). In addition, we randomly select 100 image descriptions from MS-COCO test

⁵<https://github.com/discus0434/aesthetic-predictor-v2-5?tab=readme-ov-file>

| Task | Dataset |
|--------------------------------|---|
| Language Only | OpenHermes (Teknum, 2023) |
| Multimodal ICL | MMICL (Zhao et al., 2023) |
| Multimodal CoT | ScienceQA (Lu et al., 2022) |
| Chart Understanding | Geo170K (Gao et al., 2023) |
| Instructional Image Generation | InstructPix2Pix (Brooks et al., 2023), MagicBrush (Zhang et al., 2024) |
| ASR | LibriSpeech (Panayotov et al., 2015), GigaSpeech (Chen et al., 2021), Common Voice (Ardila et al., 2020) |
| Video Dialogue | VideoChat2-IT (Li et al., 2023c) |
| Image QA | Vision-Flan (Xu et al., 2023), VizWiz (Gurari et al., 2018), LAION-GPT4V, LLaVAR (Zhang et al., 2023c), OCR-VQA (Mishra et al., 2019), VQA (Goyal et al., 2016), TextVQA (Singh et al., 2019), OK-VQA (Marino et al., 2019), Mantis-Instruct (Jiang et al., 2024) |
| Speech Generation | SpeechInstruct (Zhang et al., 2023a) |
| Speech Understanding | SpeechInstruct (Zhang et al., 2023a) |
| Image Captioning | Flickr30K (Plummer et al., 2015), MS-COCO (Lin et al., 2014) |
| Descriptive Image Generation | Flickr30K (Plummer et al., 2015), MS-COCO (Lin et al., 2014) |
| TTS | GigaSpeech (Chen et al., 2021), Common Voice (Ardila et al., 2020) |
| Video Generation | MSR-VTT (Xu et al., 2016), MSVD (Chen and Dolan, 2011b) |
| Video Understanding | MSR-VTT (Xu et al., 2016), MSVD (Chen and Dolan, 2011b), MSVD-QA (Chen and Dolan, 2011a), MSRVTQ-QA (Xu et al., 2017) |
| Visual Storytelling | VIST (Huang et al., 2016) |

Table 10: Supervised Fine-Tuning Data. “ICL” denotes In-Context Learning, and “CoT” denotes Chain of Thought.

set, and used each model to generate images accordingly for human preference evaluation. We ask 3 annotators to rank 3 images generated by the 3 models: “given the image description, which image is preferred?” The average ranking of MIO’s, AnyGPT’s, and Emu’s generated images are 1.2 (MIO), 2.9 (AnyGPT), 1.9 (Emu). MIO aligns the best with the human preference. The percentage agreement between the three annotators (calculated as the number of cases with identical rankings by all annotators divided by 100) is 82.3%, indicating a high consistency in the human evaluation.

E.2 Speech-to-Speech Evaluation

Since there is a lack of speech to speech evaluation benchmarks, we randomly sample some conversations from the moss-002-sft dataset⁶ and convert them into speech-to-speech format. Following the evaluation procedures outlined in LLaMA-Omni (Fang et al., 2024), we use the content score metric obtained from GPT-4o (OpenAI et al., 2024) to assess whether the model’s response effectively addresses the user’s instructions. The results are

⁶<https://huggingface.co/datasets/fnlp/moss-002-sft-data>

shown in Table 14.

Though the content score of MIO is slightly lower than LLaMA-Omni and AnyGPT, both LLaMA-Omni and AnyGPT first generate text replies and then convert these into voice. However, our model, MIO, is capable of directly generating speech responses to speech queries.

E.3 TTS Evaluation

We select two additional benchmarks, LibriSpeech test-clean (Panayotov et al., 2015) and GLOBE (Wang et al., 2024b), to evaluate the performance of TTS between our model and AnyGPT. For fair comparison, we don’t specify the input voice prompt during evaluation of MIO and AnyGPT. WER (Word Error Rate) and speaker similarity are employed as the automatic metrics. The results are shown in Table 16. The results show that MIO performs significantly better than AnyGPT on both WER and speaker similarity across both benchmarks.

Additionally, we conduct a human evaluation to assess the speech quality of the outputs from MIO and AnyGPT. In this evaluation, participants

| Output Modality | Text | Image | Speech | Video |
|---------------------------|-------|-------|--------|-------|
| Beam size | 5 | 1 | 1 | 1 |
| Do Sampling | False | True | True | True |
| Top-P | - | 0.7 | 0.7 | 0.7 |
| Repetition Penalty | 1.0 | 1.0 | 1.15 | 1.15 |
| Temperature | 1.0 | 1.0 | 1.0 | 1.0 |
| Guidance Scale | 1.0 | 1.0 | 1.0 | 1.0 |

Table 11: Decoding Hyperparameters.

| Task | Prompt Template |
|-------------------------|--|
| Image Captioning | Provide a one-sentence caption for the provided image. {image} |
| Image QA | (We use the prompt templates in LMMs-Eval (Li* et al., 2024)). |
| Image Generation | Please generate an image according to the given description. {description} |
| ASR | Please transcribe this speech.{speech_token} |
| TTS | Please generate a speech according to the given transcription. Start with <spch>. {transcription} |
| Text-only | The following are multiple choice questions (with answers) about {subject} {question} |
| Video QA | The goal is to use the visual information available in the image to provide an accurate answer to the question. This requires careful observation, attention to detail, and sometimes a bit of creative thinking.{video} Question: {question} Answer: |

Table 12: Prompt templates used for evaluating instruction-tuned models.

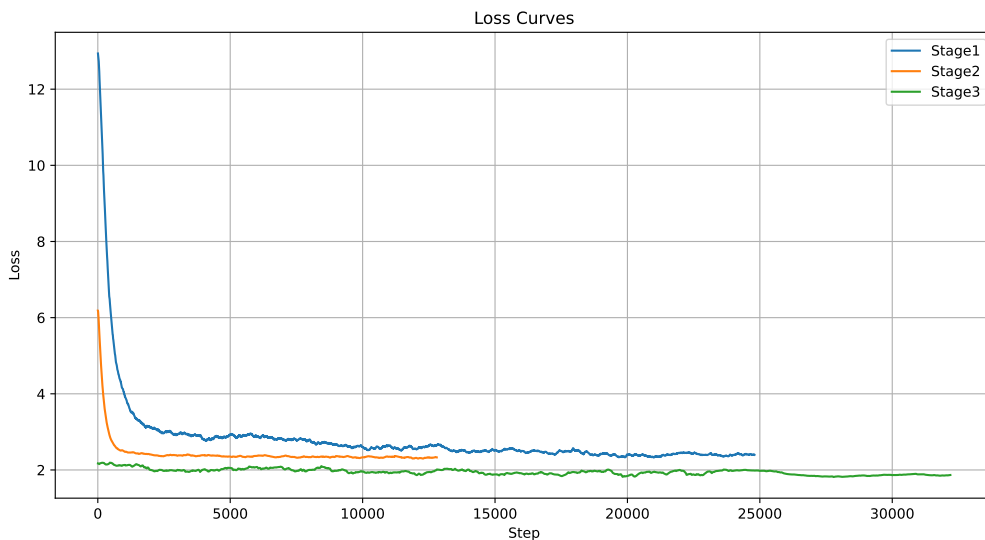


Figure 2: Loss curves of pretraing stages.

| Dataset Metric | MS-COCO | | Flickr30K | | MS-COCO Subset |
|-------------------|---------------------|--------------------------|---------------------|--------------------------|-------------------------------------|
| | SSIM (\uparrow) | Aesthetic (\uparrow) | SSIM (\uparrow) | Aesthetic (\uparrow) | Human Avg. Ranking (\downarrow) |
| Emu | 0.1749 | 3.733 | 0.1451 | 3.893 | 1.9 |
| AnyGPT | 0.1960 | 3.954 | 0.1585 | 4.251 | 2.9 |
| MIO | 0.2307 | 4.019 | 0.1727 | 4.326 | 1.2 |

Table 13: Image generation evaluation by SSIM, Aesthetic Predictor V2.5, and human preference.

| Model | Supported Workflow | Content Score (1-5 points) (\uparrow) |
|-----------------------------------|-----------------------|---|
| MIO | s2s | 1.4 |
| LLaMA-Omni (Fang et al., 2024) | s2t \rightarrow t2s | 2.4 |
| AnyGPT | s2t \rightarrow t2s | 1.8 |

Table 14: Speech-to-Speech performance. “s2s” means “speech-to-speech”, while “s2t” and “t2s” denote “speech-to-text” and “text-to-speech”, respectively.

| | |
|----------|-----|
| MIO Win | 54% |
| Tie | 25% |
| MIO Lose | 21% |

Table 15: Human evaluation for the TTS performance.

are provided with the target speech, the speech generated by AnyGPT, and the speech generated by our model. They are tasked with determining which one sounded more natural and closer to the target speech. Evaluators could choose one of the two generated speeches or indicate that they find them equally natural. Each evaluation is rated by three independent human evaluators, and we report the average scores. The results are shown in Table 15. MIO significantly outperforms AnyGPT in the human evaluation, consistent with the results from the automatic evaluation.

E.4 Loss Curves

We plot the loss curves for each stage in Figure 2. We can observe that when introducing a new data type (i.e., image-text interleaved data) in stage 2, the training loss suddenly increases. However, in the third pretraining stage, i.e., the speech-enhancement stage, the training loss transitions more smoothly. Despite the fluctuations in loss between stages, which do have some impact on downstream performance during the fluctuation periods, we find that with continued training, the model’s loss quickly recovers to its previous convergence level and continues optimizing effectively.



Figure 3: Comparing different image tokenizers for image generation within a controlled setting (limited to 3K training steps).

E.5 Demonstrations.

We illustrate the basic and advanced abilities of MIO in Figure 5 and 4. The basic abilities of MIO involve image understanding and generation, video understanding and generation, ASR, and TTS. The advanced abilities of MIO are based on its any-to-any and multimodal interleaved sequence generation features. These abilities involve visual storytelling (i.e., interleaved video-text generation), chain of visual thought, speech-in-speech-out, instructional image editing, visual guideline generation, etc. Figure 6 shows more demonstrations including multimodal chain of thought and in-context learning.

| Model | GLOBE | | LibriSpeech test-clean | |
|--------|---------|-----------------------|------------------------|-----------------------|
| | WER (↓) | Speech Similarity (↑) | WER (↓) | Speech Similarity (↑) |
| MIO | 9.8 | 67.8 | 10.3 | 75.1 |
| AnyGPT | 27.9 | 67.3 | 28.1 | 71.3 |

Table 16: More automatic evaluations for the TTS performance.

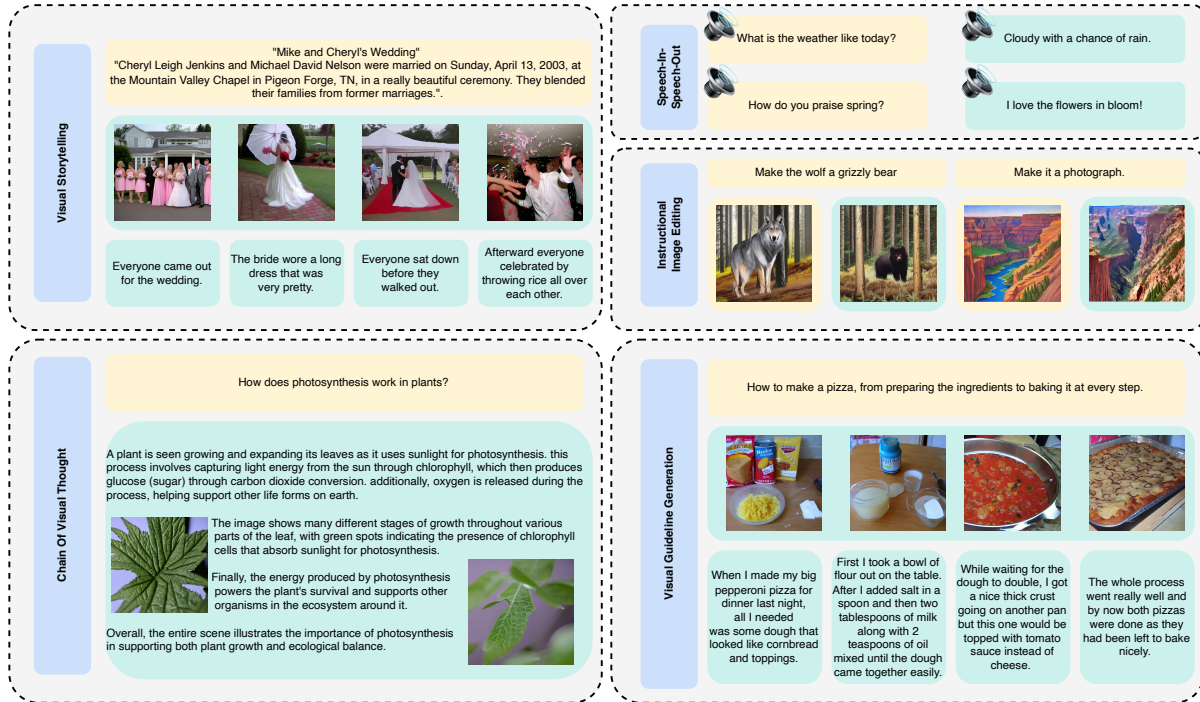


Figure 4: Demonstrations of MIO’s advanced abilities. Yellow : inputs; Green : outputs.

E.6 More Ablation Studies.

Effect of Different Image Tokenizers. The image tokenizer has a significant impact on image modality alignment. In Figure 3, we compare the image generation performance under a controlled setting after training for 3K steps in Stage I, using various image tokenizers. The image tokenizers for comparison include a VQGAN (Esser et al., 2020) with a vocabulary size of 1024 (VQGAN-1024), as well as the VQGAN-Gumbel with a vocabulary size of 8192 (VQGAN-8192)⁷. Our results indicate that the SEED-Tokenizer, which captures more semantic and higher-level image information, exhibits faster convergence. In contrast, both VQGAN tokenizers show slower convergence due to their lower-level image information.

⁷<https://github.com/CompVis/taming-transformers>

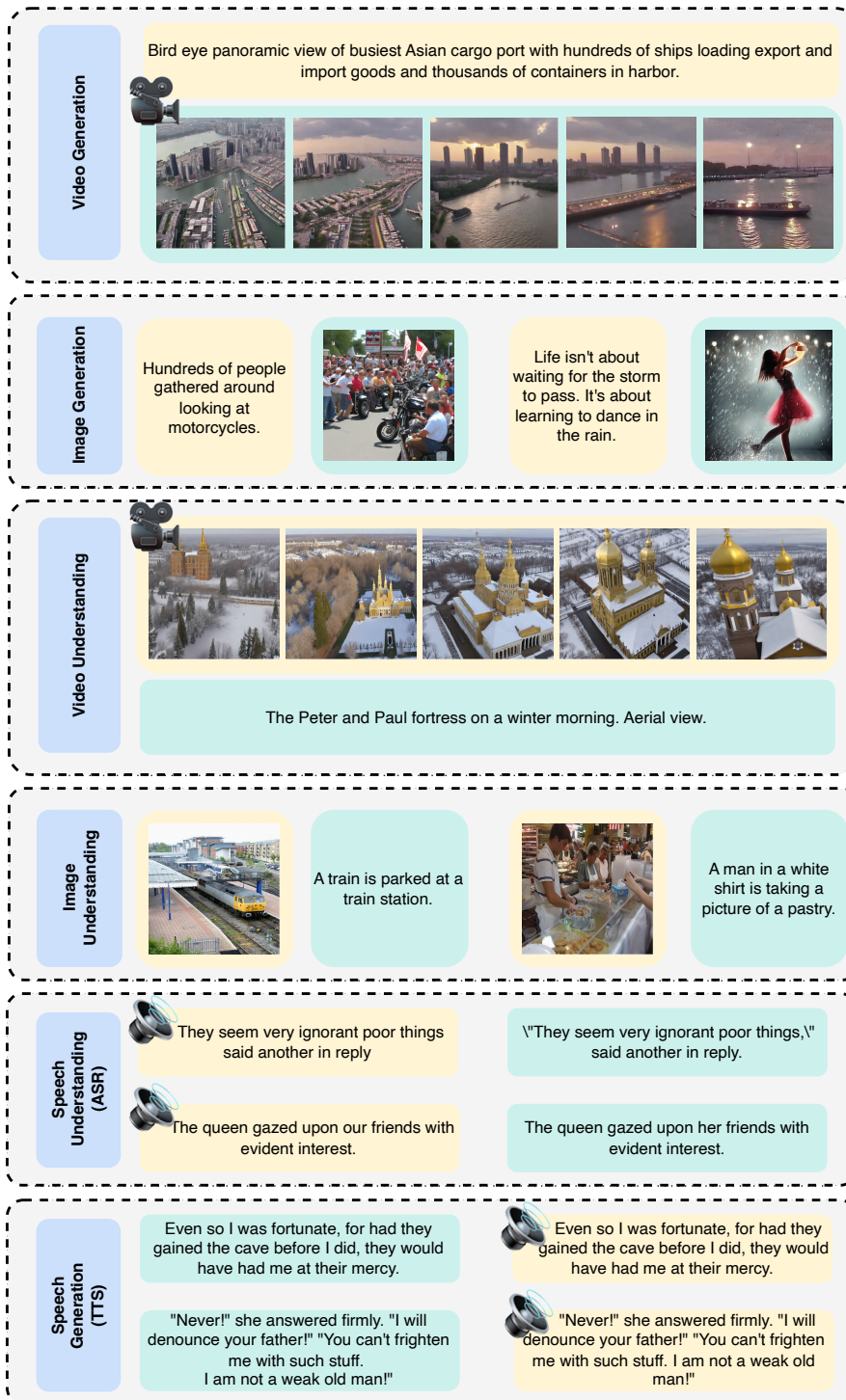


Figure 5: Demonstrations of MIO's basic abilities. Yellow : inputs; Green : outputs.

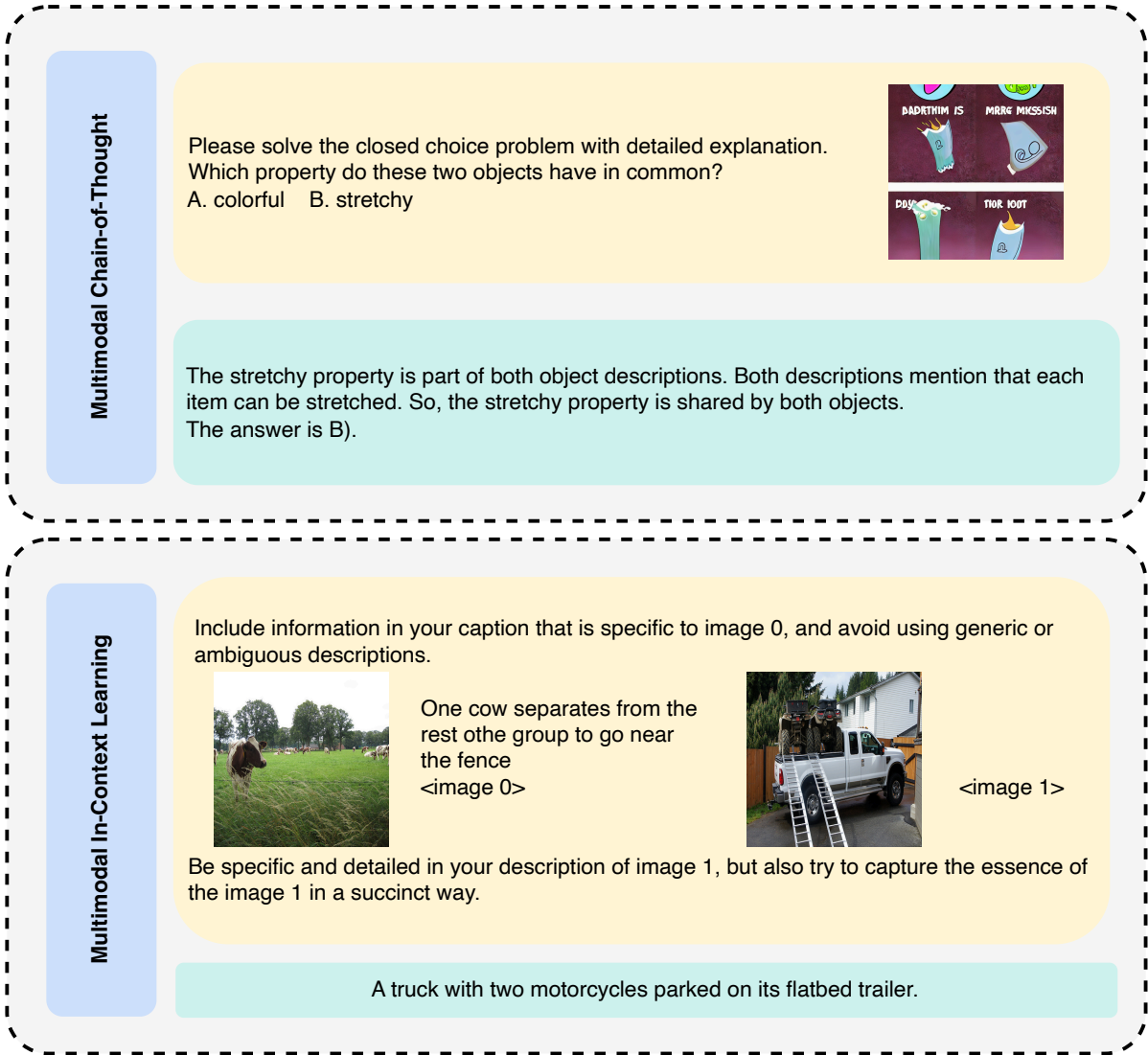


Figure 6: Multimodal Chain-of-Thought and Multimodal In-Context Learning Demos. Yellow : inputs; Green : outputs.